

# **HUMAN LANGUAGE TECHNOLOGY**

**Proceedings of a workshop held at  
Plainsboro, New Jersey  
March 21-24, 1993**

**Sponsored by:  
Advanced Research Projects Agency**

This document contains copies of reports prepared for the ARPA Human Language Technology Workshop. Included are reports from ARPA sponsored programs and other materials prepared for use at the workshop.

**APPROVED FOR PUBLIC RELEASE  
DISTRIBUTION UNLIMITED**

The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the Advanced research Projects Agency or the United States Government

**Distributed by:**  
**Morgan Kaufmann Publishers, Inc.**  
**340 Pine Street, 6th Floor**  
**San Francisco, CA 94104**  
**ISBN 1-55860-324-7**  
**Printed in the United States of America**

## TABLE OF CONTENTS

	<b>Page</b>
<b>Author Index.....</b>	<b>ix</b>
<b>Overview of the ARPA Human Language Technology Workshop, Madeleine Bates, Chairperson.....</b>	<b>3</b>
 <b>SESSION 1: Spoken Language Systems</b>	
<b>Session Summary, Alexander Rudnicky, Chair.....</b>	<b>5</b>
 <b>Benchmark Tests for the DARPA Spoken Language Program</b>	
David Pallett, Johathan Fiscus, William Fisher and John Garofolo; National Institute of Standards and Technology.....	7
 <b>Multi-Site Data Collection and Evaluation in Spoken Language Understanding</b>	
L. Hirschman, M. Bates, D. Dahl, W. Fisher, J. Garofolo, D. Pallett, K. Hunicke-Smith, P. Price, A. Rudnicky, and E. Tzoukermann.....	19
 <b>The HCRC Map Task Corpus: Natural Dialogue for Speech Recognition</b>	
Henry S. Thompson, Anne Anderson, Ellen Gurman Bard, Gwyneth Doherty-Sneddon, Alison Newlands, Cathy Sotillo; University of Edinburgh.....	25
 <b>A Portable Approach to Last Resort Parsing and Interpretation</b>	
Marcia C. Linebarger, Lewis M. Norton and Deborah A. Dahl; Paramax Systems Corporation.....	31
 <b>The Semantic Linker - A New Fragment Combining Method</b>	
David Stallard; Rusty Bobrow, BBN Systems and Technologies.....	37
 <b>Gemini: A Natural Language System for Spoken-Language Understanding</b>	
John Dowding, Jean Mark Gawron, Doug Appelt, John Bear, Lyn Cherny, Robert Moore and Doug Moran; SRI International.....	43
 <b>A Bilingual VOYAGER System</b>	
J. Glass, D. Goodine, M. Phillips, S. Sakai, S. Seneff and V. Zue; Massachusetts Institute of Technology .....	49
 <b>SESSION 2: Invited Overviews</b>	
<b>Session Summary, Madeleine Bates, Chair.....</b>	<b>55</b>
 <b>Survey of the Message Understanding Conferences</b>	
Beth Sundheim; NRaD, Nancy Chinchor; Science Applications International Corporation.....	56
 <b>Overview of TREC-1</b>	
Donna Harman; National Institute of Standards and Technology.....	61
 <b>SESSION 3: Continuous Speech Recognition</b>	
<b>Session Summary, Douglas B. Paul, Chair.....</b>	<b>67</b>

<b>Efficient Cepstral Normalization for Robust Speech Recognition</b> Fu-Hua Liu, Richard M. Stern, Xuedong Huang, Alejandro Acero; Carnegie Mellon University.....	69
<b>Comparative Experiments on Large Vocabulary Speech Recognition</b> Richard Schwartz, Tasos Anastasakos, Francis Kubala, John Makhoul, Long Nguyen, George Zavaliagkos ; BBN Systems and Technologies.....	75
<b>An Overview of the SPHINX-II Speech Recognition System</b> Xuedong Huang, Fileno Alleva, Mei-Yuh Hwang, Ronald Rosenfeld; Carnegie Mellon University.....	81
<b>Progressive-Search Algorithms for Large-Vocabulary Speech Recognition</b> Hy Murveit, John Butzberger, Vassilios Digalakis, Mitch Weintraub; SRI International.....	87
<b>Search Algorithms for Software-Only Real-Time Recognition with Very Large Vocabularies</b> Long Nguyen, Richard Schwartz, Francis Kubala, Paul Placeway; BBN Systems and Technologies.....	91
<b>Identification of Non-Linguistic Speech Features</b> Jean-Luc Gauvain and Lori F. Lamel; LIMSI-CNRS.....	96
<b>On the Use of Tied-Mixture Distributions</b> Owen Kimball and Mari Ostendorf; Boston University.....	102
<b>Adaptive Language Modeling Using the Maximum Entropy Principle</b> Raymond Lau, Ronald Rosenfeld, Salim Roukos; IBM Research Division.....	108
<b>Improved Keyword-Spotting Using SRI's DECIPHER (TM) Large-Vocabulary Speech- Recognition System</b> Mitchel Weintraub; SRI International.....	114
<b>Topic and Speaker Identification via Large Vocabulary Continuous Speech Recognition</b> Barbara Peskin, Larry Gillick, Yoshiko Ito, Stephen Lowe, Robert Roth, Francesco Scattone, James Baker, Janet Baker, John Bridle, Melvyn Hunt, Jeremy Orloff; Dragon Systems, Inc. ....	119
 <b>SESSION 4: Natural Language</b>	
<b>Session Summary, Robert C. Moore, Chair.....</b>	125
 <b>Heuristics for Broad-Coverage Natural Language Parsing</b> Michael C. McCord; IBM T. J. Watson Research Center.....	127
 <b>FASTUS: A System for Extracting Information from Text</b> Jerry R. Hobbs, Douglas Appelt, John Bear, David Israel, Megumi Kameyama, and Mabry Tyson; SRI International.....	133
 <b>Interpreting Temporal Adverbials</b> Chung Hee Hwang & Lenhart K. Schubert; University of Rochester.....	138
 <b>The Murasaki Project: Multilingual Natural Language Understanding</b> Chinatsu Aone, Hatte Blejer, Sharon Flank, Douglas McKee, Sandy Shinn; Systems Research and Applications.....	144
 <b>Validation of Terminological Inference in an Information Extraction Task</b> Marc Vilain; The MITRE Corporation.....	150

**SESSION 5: Discourse**

<b>Session Summary, Jerry R. Hobbs, Chair.....</b>	157
--	-----

<b>Development, Implementation and Testing of a Discourse Model for Newspaper Texts</b> Elizabeth D. Liddy; Syracuse University, Kenneth McVearry; Coherent Research, Inc., Woojin Paik; Syracuse University, Edmund Yu; Syracuse University, Mary McKenna; Syracuse University.....	159
---	-----

<b>Indexing and Exploiting a Discourse History to Generate Context-Sensitive Explanations</b> Johanna D. Moore; University of Pittsburgh.....	165
--	-----

<b>Generic Plan Recognition for Dialogue Systems</b> George Ferguson, James F. Allen; University of Rochester.....	171
---	-----

<b>Efficient Collaborative Discourse: A Theory and Its Implementation</b> Alan W. Biermann, Curry I. Guinn, D. Richard Hipp, Ronnie W. Smith; Duke University.....	177
---	-----

**SESSION 6: Machine Translation**

<b>Session Summary, Alex Waibel, Chair.....</b>	183
---	-----

<b>Building a Large Ontology for Machine Translation</b> Kevin Knight; USC/Information Sciences Institute.....	185
---	-----

<b>LINGSTAT: An Interactive, Machine-Aided Translation System</b> Jonathan Yamron, James Baker, Paul Bamberg, Haakon Chevalier, Taiko Dietzel, John Elder, Frank Kampmann, Mark Mandel, Linda Manganaro, Todd Margolis, Elizabeth Steele; Dragon Systems, Inc.....	191
---	-----

<b>An MAT Tool and Its Effectiveness</b> Robert Frederking, Dean Grannes, Peter Cousseau, Sergei Nirenburg; Carnegie Mellon University.....	196
--	-----

<b>But Dictionaries are Data Too</b> Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, Meredith J. Goldsmith, Jan Hajic, Robert L. Mercer, and Surya Mohanty; IBM T. J. Watson Research Center.....	202
--	-----

<b>Evaluation of Machine Translation</b> John S. White; PRC, Theresa A. O'Connell; PRC, Lynn M. Carlson; DoD.....	206
--	-----

<b>Recent Advances in Janus: A Speech Translation System</b> M. Woszczyna, N. Coccaro, A. Eisele, A. Lavie, A. McNair, T. Polzin, I. Rogina, C. P. Rose, T. Sloboda, M. Tomita, J. Tsutsumi, N. Waibel, A. Waibel, W. Ward; Carnegie Mellon University.....	211
--	-----

<b>A Speech to Speech Translation System Built from Standard Components</b> Manny Rayner, Hiyan Alshawi; SRI International, Cambridge, UK, Ivan Bretan; SICS, David Carter; SRI International, Cambridge, UK, Vassilios Digalakis; SRI International, Menlo Park, CA, Bjorn Gamback; SICS, Jann Kaja; Telia Research, Jussi Karlsgren; SICS, Bertil Lyberg; Telia Research, Steve Pulman; SRI International, Cambridge, UK, Patti Price; SRI International, Menlo Park, CA; and Christer Samuelsson; SICS.....	217
---	-----

**SESSION 7: Demonstrations**

<b>Chair: Hy Murveit; SRI International.....</b>	223
--	-----

	Page
<b>SESSION 8: Statistical Language Modeling</b>	
<b>Session Summary, Mitchell Marcus, Chair.....</b>	225
<b>Example-Based Correction of Word Segmentation and Part of Speech Labelling</b>	
Tomoyoshi Matsukawa, Scott Miller, Ralph Weischedel; BBN Systems and Technologies.....	227
<b>Measures and Models for Phrase Recognition</b>	
Steven Abney; Bell Communications Research.....	233
<b>Automatic Grammar Induction and Parsing Free Text: A Transformation-Based Approach</b>	
Eric Brill; University of Pennsylvania.....	237
<b>Prediction of Lexicalized Tree Fragments in Text</b>	
Donald Hindle; AT&T Bell Laboratories.....	243
<b>Hypothesizing Word Association from Untagged Text</b>	
Tomoyoshi Matsukawa; BBN Systems and Technologies.....	248
<b>Smoothing of Automatically Generated Selectional Constraints</b>	
Ralph Grishman and John Sterling; New York University.....	254
<b>Corpus-Based Statistical Sense Resolution</b>	
Claudia Leacock; Princeton University, Geoffrey Towell, Ellen Voorhees; Siemens University.....	260
<b>One Sense Per Collocation</b>	
David Yarowsky; University of Pennsylvania.....	266
<b>Augmenting Lexicons Automatically: Clustering Semantically Related Adjectives</b>	
Kathleen McKeown and Vasileios Hatzivassiloglou; Columbia University.....	272
<b>Semantic Classes and Syntactic Ambiguity</b>	
Philip Resnik; University of Pennsylvania.....	278
<b>SESSION 9: Government Panel</b>	
<b>Session Summary, Carol J. Van Ess-Dykema, Chair.....</b>	285
<b>Projected Government Needs in Human Language Technology and the Role of Researchers in Meeting Them</b>	
Helen M. Gigley; Naval Research Laboratory.....	287
<b>Language Research Sponsored by ONR</b>	
Susan Chipman; Office of Naval Research.....	290
<b>Technology Transfer: Problems and Prospects</b>	
Jesse Fussell; Department of Defense.....	295
<b>SESSION 10: The Lexicon</b>	
<b>Session Summary, Ralph Grishman, Chair.....</b>	299

	Page
<b>The Complex Syntax Project</b> Ralph Grishman, Catherine Macleod, Susanne Wolff; New York University.....	300
<b>A Semantic Concordance</b> George A. Miller, Claudia Leacock, Randee Tengi, Ross Bunker; Princeton University.....	303
<b>Interpretation of Proper Nouns for Information Retrieval</b> Woojin Paik, Elizabeth D. Liddy, Edmund Yu, Mary McKenna; Syracuse University.....	309
 <b>SESSION 11: Prosody</b>	
<b>Session Summary, Mari Ostendorf, Chair.....</b>	<b>315</b>
 <b>On Customizing Prosody in Speech Synthesis: Names and Addresses as a Case in Point</b> Kim Silverman; NYNEX Science and Technology, Inc.....	317
<b>Quantitative Modeling of Segmental Duration</b> Jan P. H. Van Santen; AT&T Bell Laboratories.....	323
<b>A Speech-First Model for Repair Detection and Correction</b> Christine Nakatani; Harvard University, Julia Hirschberg; AT&T Bell Laboratories .....	329
<b>Prosody/Parse Scoring and Its Application in ATIS</b> N. M. Veilleux and M. Ostendorf; Boston University.....	335
<b>Perceived Prosodic Boundaries and Their Phonetic Correlates</b> Rene Collier, Jan Roelof de Pijper and Angelien Sanderman; Institute for Perception Research.....	341
 <b>SESSION 12: Document Retrieval and Text Retrieval</b>	
<b>Session Summary, Karen Sparck Jones, Chair.....</b>	<b>347</b>
 <b>The Importance of Proper Weighting Methods</b> Chris Buckley; Cornell University.....	349
<b>Query Processing for Retrieval from Large Text Bases</b> John Broglio and W. Bruce Croft; University of Massachusetts.....	353
<b>An Overview of DR-LINK and Its Approach to Document Filtering</b> Elizabeth D. Liddy, Woojin Paik, Edmund S. Yu; Syracuse University and Kenneth A. McVearry; Coherent Research, Inc. ....	358
 <b>SESSION 13: New Directions</b>	
<b>Session Summary, Ralph Weischedel, Chair.....</b>	<b>363</b>
 <b>Mode Preference in a Simple Data-Retrieval Task</b> Alexander I. Rudnicky; Carnegie Mellon University.....	364

	Page
<b>A Simulation-Based Research Strategy for Designing Complex NL Systems</b> Sharon Oviatt, Philip Cohen, Michelle Wang, Jeremy Gaston; SRI International.....	370
<b>Speech and Text-Image Processing in Documents</b> Marcia A. Bush; Xerox Palo Alto Research Center.....	376
<b>SITE REPORTS.....</b>	<b>381</b>

## AUTHOR INDEX

Abney, S.	233	Chevalier, H.	191
Acero, A.	69	Chinchor, N.	56
Allan, J.	392	Chipman, S.	290
Allen, J.	171, 420	Coccaro, N.	211
Alleva, F.	81	Cohen, P.	370
Alshawi, H.	217	Collier, R.	341
Anastasakos, T.	75	Cousseau, P.	196
Anderson, A.	25	Cowie, J.	405
Aone, C.	144	Croft, W.	353, 418
Appelbaum, L.	383	Dahl, D.	19, 31
Appelt D.	43, 133	Danielson, D.	413
Baker, James	119, 191, 393	De Pijper, J.	341
Baker, Janet	119, 394	Della Pietra, S.	202
Bamberg, P.	191	Della Pietra, V.	202
Bard, E.	25	Dietzel, T.	191
Bates, M.	3, 19, 55, 387	Digalakis, V.	87, 217, 415
Bear, J.	43, 133	Doherty-Sneddon, G.	25
Bernstein, J.	412, 413	Dolan, C.	417
Biermann, A.	177	Dowding, J.	43
Blejer, H.	144	Eisele, A.	211
Bobrow, R.	37	Elder, J.	191
Bretan, I.	217	Ferguson, G.	171
Bridle, J.	119	Fiscus, J.	7
Brill, E.	237	Fisher, W.	7, 19
Broglio, J.	353	Flank, S.	144
Brown, P.	202, 397	Frederking, R.	196
Buckley, C.	349, 392	Fussell, J.	295
Bunker, R.	303	Gallant, S.	396
Bush, M.	376	Gamback, B.	217
Butzberger, J.	87	Garofolo, J.	7, 19
Caid, W.	396	Gaston, J.	370
Carlson, L.	206	Gauvain, J.	96
Carter, D.	217	Gawron, J.	43
Cherny, L.	43	Gigley, H.	287

Gillick, L.	119, 394	Leacock, C.	260, 303
Glass, J.	49	Lehnert, W.	417
Goldsmith, M.	202	Liddy, E.	159, 309, 358, 416
Goodine, D.	49	Linebarger, M.	31
Grannes, D.	196	Lippmann, R.	399
Grishman, R.	254, 299, 300, 407	Liu, F.	69
Guinn, C.	177	Lowe, S.	119
Hajic, J.	202	Lyberg, B.	217
Harman, D.	61	Macleod, C.	300
Hatzivassiloglou, V.	272	Makhoul, J.	75, 385, 387
Hindle, D.	243	Mandel, M.	191
Hipp, D.	177	Manganaro, L.	191
Hirschberg, J.	329	Marcus, M.	225, 419
Hirschman, L.	19, 401	Margolis, T.	191
Hobbs, J.	133, 157	Matsukawa, T.	227, 248
Hovy, E.	421	McCord, M.	127
Huang, X.	69, 81	McKee, D.	144
Hunicke-Smith, D.	19, 412	McKenna, M.	159, 309
Hunt, M.	119	McKeown, K.	272, 391
Hwang, C.	138	McNair, A.	211
Hwang, M.	81	McVearry, K.	159, 358
Israel, D.	133	Mercer, R.	202
Ito, Y.	119	Miller, G.	303, 409
Jacobs, P.	395	Miller, S.	227
Joshi, A.	419	Mohanty, S.	202
Kaja, J.	217	Moore, J.	165
Kameyama, M.	133	Moore, R.	43, 125, 414
Kampmann, F.	191	Moran, D.	43
Karlgren, J.	217	Murveit, H.	87, 415
Kimball, O.	102	Myaeng, S.	416
Knight, K.	185	Nakatani, C.	329
Kubala, F.	75, 91	Newlands, A.	25
Kwok, K.	410	Nguyen, L.	75, 91
Lamel, L.	96	Nirenburg, S.	196
Lau, R.	108	Norton, L.	31
Lavie, A.	211	O'Connell, T.	206

Orloff, J.	119	Sotillo, C.	25
Ostendorf, M.	102, 315, 335, 388, 389	Sparck Jones, K.	347
Oviatt, S.	370	Stallard, D.	37
Paik, W.	159, 309, 358	Steedman, M.	419
Pallett, D.	7, 19, 402	Steele, E.	191
Passonneau, R.	391	Sterling, J.	254
Paul D.	67, 400	Stern, R.	69
Peskin, B.	119	Strzalkowski, T.	408
Phillips, M.	49	Sundheim, B.	56, 403
Placeway, P.	91	Tengi, R.	303
Polzin, T.	211	Thompson, H.	25
Price, P.	19, 217, 388, 414	Tomita, M.	211
Pulman S.	217	Tong, R.	383
Pustejovsky, J.	405	Towell, G.	260
Rayner, M.	217	Tsutsumi, J.	211
Reddy, R.	390	Tyson, M.	133
Resnik, P.	278	Tzoukermann, E.	19
Rogina, I.	211	Van Ess-Dykema, C.	285
Rohlicek, J.	384, 389	Van Santen, J.	323
Rose, C.	211	Veilleux, N.	335
Rosenfeld, R.	81, 108	Vilain, M.	150
Roth, R.	119, 394	Voorhees, E.	260, 411
Roukos, S.	108, 398	Waibel, A.	183, 211
Rudnický, A.	5, 19, 364	Waibel, N.	211
Sakai, S.	49	Wang, M.	370
Salton, G.	392	Ward, W.	211
Samuelsson, C.	217	Webber, B.	419
Sanderman, A.	341	Weinstein, C.	400
Scattone, F.	119	Weintraub, M.	87, 114, 415
Schubert, L.	138, 420	Weischedel R.	227, 363, 386
Schwartz, R.	75, 91, 385	White, J.	206
Seneff, S.	49	Wilks, Y.	404, 405, 406
Shinn, S.	144	Wolff, S.	300
Silverman, K.	317	Woszczyna, M.	211
Sloboda, T.	211	Yamron, J.	191, 393
Smith, R.	177	Yarowsky, D.	266

- Yu, E. ....159, 309, 358  
Zavaliagkos, G. ....75  
Zue, V. ....49, 401