

XUEDONG HUANG | ALEX ACERO | HSIAO-WUEN HON

SPOKEN LANGUAGE PROCESSING

A Guide to Theory, Algorithm, and System Development



Foreword by Dr. Raj Reddy
Carnegie Mellon University

Spoken Language Processing



Spoken Language Processing

**A Guide to Theory, Algorithm,
and System Development**

Xuedong Huang

Alex Acero

Hsiao-Wuen Hon

Microsoft Research



**Prentice Hall PTR
Upper Saddle River, New Jersey 07458
www.phptr.com**

Library of Congress Cataloging-in-Publication Data

Huang, Xuedong.

Spoken language processing: a guide to theory, algorithm, and system development/
Xuedong Huang, Alex Acero, Hsiao-Wuen Hon.

p. cm.

Includes bibliographical references and index.

ISBN 0-13-022616-5

1. Natural language processing (Computer science) I. Acero, Alex. II. Hon,
Hsiao-Wuen. III. Title.

QA76.9.N38 H83 2001

00-050196

006.3'5—dc21

Editorial/production supervision: *Jane Bonnell*

Cover design director: *Jerry Votta*

Cover design: *Anthony Gemmellaro*

Manufacturing buyer: *Maura Zaldivar*

Development editor: *Russ Hall*

Acquisitions editor: *Tim Moore*

Editorial assistant: *Allyson Kloss*

Marketing manager: *Debby van Dijk*



© 2001 by Prentice Hall PTR
Prentice-Hall, Inc.
Upper Saddle River, New Jersey 07458

Prentice Hall books are widely used by corporations and government agencies for training, marketing, and resale.

The publisher offers discounts on this book when ordered in bulk quantities. For more information, contact Corporate Sales Department, Phone: 800-382-3419; FAX: 201-236-7141;

E-mail: corpsales@prenhall.com

Or write: Prentice Hall PTR, Corporate Sales Dept., One Lake Street, Upper Saddle River, NJ 07458.

Company and product names mentioned herein are the trademarks or registered trademarks of their respective owners.

All rights reserved. No part of this book may be reproduced, in any form or by any means, without permission in writing from the publisher.

Printed in the United States of America

ISBN 0-13-022616-5

Prentice-Hall International (UK) Limited, *London*

Prentice-Hall of Australia Pty. Limited, *Sydney*

Prentice-Hall Canada Inc., *Toronto*

Prentice-Hall Hispanoamericana, S.A., *Mexico*

Prentice-Hall of India Private Limited, *New Delhi*

Prentice-Hall of Japan, Inc., *Tokyo*

Pearson Education Asia Pte. Ltd.

Editora Prentice-Hall do Brasil, Ltda., *Rio de Janeiro*

Amazon/VB Assets

Exhibit 1012

Page 5

To Yingzhi, Angela, Christina, and Derek

To Donna and Nicolas

To Phen, Stephanie, and Jacqueline

Contents

FOREWORD	xxi
PREFACE	xxv
1. INTRODUCTION	1
1.1. MOTIVATIONS	2
<i>1.1.1. Spoken Language Interface.</i>	2
<i>1.1.2. Speech-to-Speech Translation.</i>	3
<i>1.1.3. Knowledge Partners.</i>	3
1.2. SPOKEN LANGUAGE SYSTEM ARCHITECTURE	4
<i>1.2.1. Automatic Speech Recognition.</i>	4
<i>1.2.2. Text-to-Speech Conversion .</i>	6
<i>1.2.3. Spoken Language Understanding .</i>	7
1.3. BOOK ORGANIZATION	8
<i>1.3.1. Part I: Fundamental Theory</i>	9
<i>1.3.2. Part II: Speech Processing</i>	9
<i>1.3.3. Part III: Speech Recognition</i>	9
<i>1.3.4. Part IV: Text-to-Speech Systems</i>	10
<i>1.3.5. Part V: Spoken Language Systems</i>	10
1.4. TARGET AUDIENCES.....	10
1.5. HISTORICAL PERSPECTIVE AND FURTHER READING.....	11

PART I: FUNDAMENTAL THEORY

2. SPOKEN LANGUAGE STRUCTURE	19
2.1. SOUND AND HUMAN SPEECH SYSTEMS	21
<i>2.1.1. Sound.</i>	21
<i>2.1.2. Speech Production .</i>	24
<i>2.1.3. Speech Perception.</i>	29

2.2. PHONETICS AND PHONOLOGY.....	36
2.2.1. <i>Phonemes</i>	36
2.2.2. <i>The Allophone: Sound and Context</i>	47
2.2.3. <i>Speech Rate and Coarticulation</i>	49
2.3. SYLLABLES AND WORDS.....	51
2.3.1. <i>Syllables</i>	51
2.3.2. <i>Words</i>	53
2.4. SYNTAX AND SEMANTICS.....	58
2.4.1. <i>Syntactic Constituents</i>	58
2.4.2. <i>Semantic Roles</i>	63
2.4.3. <i>Lexical Semantics</i>	64
2.4.4. <i>Logical Form</i>	67
2.5. HISTORICAL PERSPECTIVE AND FURTHER READING.....	68
3. PROBABILITY, STATISTICS, AND INFORMATION THEORY .	73
3.1. PROBABILITY THEORY.....	74
3.1.1. <i>Conditional Probability and Bayes' Rule</i>	75
3.1.2. <i>Random Variables</i>	77
3.1.3. <i>Mean and Variance</i>	79
3.1.4. <i>Covariance and Correlation</i>	82
3.1.5. <i>Random Vectors and Multivariate Distributions</i>	83
3.1.6. <i>Some Useful Distributions</i>	85
3.1.7. <i>Gaussian Distributions</i>	92
3.2. ESTIMATION THEORY	98
3.2.1. <i>Minimum/Least Mean Squared Error Estimation</i>	99
3.2.2. <i>Maximum Likelihood Estimation</i>	104
3.2.3. <i>Bayesian Estimation and MAP Estimation</i>	107
3.3. SIGNIFICANCE TESTING.....	113
3.3.1. <i>Level of Significance</i>	114
3.3.2. <i>Normal Test (Z-Test)</i>	115
3.3.3. <i>χ^2 Goodness-of-Fit Test</i>	116
3.3.4. <i>Matched-Pairs Test</i>	118
3.4. INFORMATION THEORY	120
3.4.1. <i>Entropy</i>	120
3.4.2. <i>Conditional Entropy</i>	123
3.4.3. <i>The Source Coding Theorem</i>	124
3.4.4. <i>Mutual Information and Channel Coding</i>	126
3.5. HISTORICAL PERSPECTIVE AND FURTHER READING.....	128
4. PATTERN RECOGNITION	133
4.1. BAYES' DECISION THEORY.....	134
4.1.1. <i>Minimum-Error-Rate Decision Rules</i>	135

4.1.2. <i>Discriminant Functions</i>	138
4.2. HOW TO CONSTRUCT CLASSIFIERS.....	140
4.2.1. <i>Gaussian Classifiers</i>	142
4.2.2. <i>The Curse of Dimensionality</i>	144
4.2.3. <i>Estimating the Error Rate</i>	146
4.2.4. <i>Comparing Classifiers</i>	148
4.3. DISCRIMINATIVE TRAINING.....	150
4.3.1. <i>Maximum Mutual Information Estimation</i>	150
4.3.2. <i>Minimum-Error-Rate Estimation</i>	156
4.3.3. <i>Neural Networks</i>	158
4.4. UNSUPERVISED ESTIMATION METHODS	163
4.4.1. <i>Vector Quantization</i>	163
4.4.2. <i>The EM Algorithm</i>	170
4.4.3. <i>Multivariate Gaussian Mixture Density Estimation</i>	172
4.5. CLASSIFICATION AND REGRESSION TREES	175
4.5.1. <i>Choice of Question Set</i>	177
4.5.2. <i>Splitting Criteria</i>	178
4.5.3. <i>Growing the Tree</i>	181
4.5.4. <i>Missing Values and Conflict Resolution</i>	182
4.5.5. <i>Complex Questions</i>	182
4.5.6. <i>The Right-Sized Tree</i>	184
4.6. HISTORICAL PERSPECTIVE AND FURTHER READING.....	190

PART II: SPEECH PROCESSING

5. DIGITAL SIGNAL PROCESSING	201
5.1. DIGITAL SIGNALS AND SYSTEMS.....	202
5.1.1. <i>Sinusoidal Signals</i>	203
5.1.2. <i>Other Digital Signals</i>	206
5.1.3. <i>Digital Systems</i>	206
5.2. CONTINUOUS-FREQUENCY TRANSFORMS.....	208
5.2.1. <i>The Fourier Transform</i>	208
5.2.2. <i>Z-Transform</i>	211
5.2.3. <i>Z-Transforms of Elementary Functions</i>	212
5.2.4. <i>Properties of the Z- and Fourier Transforms</i>	215
5.3. DISCRETE-FREQUENCY TRANSFORMS	216
5.3.1. <i>The Discrete Fourier Transform (DFT)</i>	218
5.3.2. <i>Fourier Transforms of Periodic Signals</i>	219
5.3.3. <i>The Fast Fourier Transform (FFT)</i>	222
5.3.4. <i>Circular Convolution</i>	227
5.3.5. <i>The Discrete Cosine Transform (DCT)</i>	228

5.4. DIGITAL FILTERS AND WINDOWS	229
5.4.1. <i>The Ideal Low-Pass Filter</i>	229
5.4.2. <i>Window Functions</i>	230
5.4.3. <i>FIR Filters</i>	232
5.4.4. <i>IIR Filters</i>	238
5.5. DIGITAL PROCESSING OF ANALOG SIGNALS	242
5.5.1. <i>Fourier Transform of Analog Signals</i>	243
5.5.2. <i>The Sampling Theorem</i>	243
5.5.3. <i>Analog-to-Digital Conversion</i>	245
5.5.4. <i>Digital-to-Analog Conversion</i>	246
5.6. MULTIRATE SIGNAL PROCESSING.....	248
5.6.1. <i>Decimation</i>	248
5.6.2. <i>Interpolation</i>	249
5.6.3. <i>Resampling</i>	250
5.7. FILTERBANKS	251
5.7.1. <i>Two-Band Conjugate Quadrature Filters</i>	251
5.7.2. <i>Multiresolution Filterbanks</i>	254
5.7.3. <i>The DFT as a Filterbank</i>	255
5.7.4. <i>Modulated Lapped Transforms</i>	258
5.8. STOCHASTIC PROCESSES	260
5.8.1. <i>Statistics of Stochastic Processes</i>	261
5.8.2. <i>Stationary Processes</i>	264
5.8.3. <i>LTI Systems with Stochastic Inputs</i>	267
5.8.4. <i>Power Spectral Density</i>	268
5.8.5. <i>Noise</i>	269
5.9. HISTORICAL PERSPECTIVE AND FURTHER READING	270
6. SPEECH SIGNAL REPRESENTATIONS.....	275
6.1. SHORT-TIME FOURIER ANALYSIS	276
6.1.1. <i>Spectrograms</i>	281
6.1.2. <i>Pitch-Synchronous Analysis</i>	283
6.2. ACOUSTICAL MODEL OF SPEECH PRODUCTION	283
6.2.1. <i>Glottal Excitation</i>	283
6.2.2. <i>Lossless Tube Concatenation</i>	284
6.2.3. <i>Source-Filter Models of Speech Production</i>	284
6.3. LINEAR PREDICTIVE CODING.....	288
6.3.1. <i>The Orthogonality Principle</i>	290
6.3.2. <i>Solution of the LPC Equations</i>	291
6.3.3. <i>Spectral Analysis via LPC</i>	292
6.3.4. <i>The Prediction Error</i>	300
6.3.5. <i>Equivalent Representations</i>	301
	303

6.4. CEPSTRAL PROCESSING	306
6.4.1. <i>The Real and Complex Cepstrum</i>	307
6.4.2. <i>Cepstrum of Pole-Zero Filters</i>	308
6.4.3. <i>Cepstrum of Periodic Signals</i>	311
6.4.4. <i>Cepstrum of Speech Signals</i>	312
6.4.5. <i>Source-Filter Separation via the Cepstrum</i>	314
6.5. PERCEPTUALLY MOTIVATED REPRESENTATIONS	315
6.5.1. <i>The Bilinear Transform</i>	315
6.5.2. <i>Mel-Frequency Cepstrum</i>	316
6.5.3. <i>Perceptual Linear Prediction (PLP)</i>	318
6.6. FORMANT FREQUENCIES	319
6.6.1. <i>Statistical Formant Tracking</i>	320
6.7. THE ROLE OF PITCH.....	324
6.7.1. <i>Autocorrelation Method</i>	324
6.7.2. <i>Normalized Cross-Correlation Method</i>	327
6.7.3. <i>Signal Conditioning</i>	329
6.7.4. <i>Pitch Tracking</i>	330
6.8. HISTORICAL PERSPECTIVE AND FURTHER READING.....	332
7. SPEECH CODING	337
7.1. SPEECH CODERS ATTRIBUTES	338
7.2. SCALAR WAVEFORM CODERS	340
7.2.1. <i>Linear Pulse Code Modulation (PCM)</i>	340
7.2.2. <i>μ-law and A-law PCM</i>	342
7.2.3. <i>Adaptive PCM</i>	344
7.2.4. <i>Differential Quantization</i>	345
7.3. SCALAR FREQUENCY DOMAIN CODERS.....	348
7.3.1. <i>Benefits of Masking</i>	349
7.3.2. <i>Transform Coders</i>	350
7.3.3. <i>Consumer Audio</i>	351
7.3.4. <i>Digital Audio Broadcasting (DAB)</i>	352
7.4. CODE EXCITED LINEAR PREDICTION (CELP).....	353
7.4.1. <i>LPC Vocoder</i>	353
7.4.2. <i>Analysis by Synthesis</i>	353
7.4.3. <i>Pitch Prediction: Adaptive Codebook</i>	356
7.4.4. <i>Perceptual Weighting and Postfiltering</i>	357
7.4.5. <i>Parameter Quantization</i>	358
7.4.6. <i>CELP Standards</i>	359
7.5. LOW-BIT RATE SPEECH CODERS	361
7.5.1. <i>Mixed-Excitation LPC Vocoder</i>	362
7.5.2. <i>Harmonic Coding</i>	363
7.5.3. <i>Waveform Interpolation</i>	367
7.6. HISTORICAL PERSPECTIVE AND FURTHER READING.....	371

PART III: SPEECH RECOGNITION

8. HIDDEN MARKOV MODELS.....	377
8.1. THE MARKOV CHAIN	378
8.2. DEFINITION OF THE HIDDEN MARKOV MODEL.....	380
8.2.1. <i>Dynamic Programming and DTW</i>	383
8.2.2. <i>How to Evaluate an HMM—The Forward Algorithm</i>	385
8.2.3. <i>How to Decode an HMM—The Viterbi Algorithm</i>	387
8.2.4. <i>How to Estimate HMM Parameters—Baum-Welch Algorithm</i>	389
8.3. CONTINUOUS AND SEMICONTINUOUS HMMS	394
8.3.1. <i>Continuous Mixture Density HMMs</i>	394
8.3.2. <i>Semicontinuous HMMs</i>	396
8.4. PRACTICAL ISSUES IN USING HMMS.....	398
8.4.1. <i>Initial Estimates</i>	398
8.4.2. <i>Model Topology</i>	399
8.4.3. <i>Training Criteria</i>	401
8.4.4. <i>Deleted Interpolation</i>	401
8.4.5. <i>Parameter Smoothing</i>	403
8.4.6. <i>Probability Representations</i>	404
8.5. HMM LIMITATIONS	405
8.5.1. <i>Duration Modeling</i>	406
8.5.2. <i>First-Order Assumption</i>	408
8.5.3. <i>Conditional Independence Assumption</i>	409
8.6. HISTORICAL PERSPECTIVE AND FURTHER READING.....	409
9. ACOUSTIC MODELING.....	415
9.1. VARIABILITY IN THE SPEECH SIGNAL	416
9.1.1. <i>Context Variability</i>	417
9.1.2. <i>Style Variability</i>	418
9.1.3. <i>Speaker Variability</i>	418
9.1.4. <i>Environment Variability</i>	419
9.2. HOW TO MEASURE SPEECH RECOGNITION ERRORS.....	419
9.3. SIGNAL PROCESSING—EXTRACTING FEATURES	421
9.3.1. <i>Signal Acquisition</i>	422
9.3.2. <i>End-Point Detection</i>	422
9.3.3. <i>MFCC and Its Dynamic Features</i>	424
9.3.4. <i>Feature Transformation</i>	426
9.4. PHONETIC MODELING—SELECTING APPROPRIATE UNITS	426
9.4.1. <i>Comparison of Different Units</i>	428
9.4.2. <i>Context Dependency</i>	429
9.4.3. <i>Clustered Acoustic-Phonetic Units</i>	430
9.4.4. <i>Lexical Baseforms</i>	432
	436

9.5. ACOUSTIC MODELING—SCORING ACOUSTIC FEATURES	439
9.5.1. <i>Choice of HMM Output Distributions</i>	439
9.5.2. <i>Isolated vs. Continuous Speech Training</i>	441
9.6. ADAPTIVE TECHNIQUES—MINIMIZING MISMATCHES.....	444
9.6.1. <i>Maximum a Posteriori (MAP)</i>	445
9.6.2. <i>Maximum Likelihood Linear Regression (MLLR)</i>	447
9.6.3. <i>MLLR and MAP Comparison</i>	450
9.6.4. <i>Clustered Models</i>	452
9.7. CONFIDENCE MEASURES: MEASURING THE RELIABILITY	453
9.7.1. <i>Filler Models</i>	453
9.7.2. <i>Transformation Models</i>	454
9.7.3. <i>Combination Models</i>	456
9.8. OTHER TECHNIQUES.....	457
9.8.1. <i>Neural Networks</i>	457
9.8.2. <i>Segment Models</i>	459
9.9. CASE STUDY: WHISPER.....	464
9.10. HISTORICAL PERSPECTIVE AND FURTHER READING	465
10. ENVIRONMENTAL ROBUSTNESS.....	477
10.1. THE ACOUSTICAL ENVIRONMENT	478
10.1.1. <i>Additive Noise</i>	478
10.1.2. <i>Reverberation</i>	480
10.1.3. <i>A Model of the Environment</i>	482
10.2. ACOUSTICAL TRANSDUCERS	486
10.2.1. <i>The Condenser Microphone</i>	486
10.2.2. <i>Directionality Patterns</i>	489
10.2.3. <i>Other Transduction Categories</i>	496
10.3. ADAPTIVE ECHO CANCELLATION (AEC)	497
10.3.1. <i>The LMS Algorithm</i>	499
10.3.2. <i>Convergence Properties of the LMS Algorithm</i>	500
10.3.3. <i>Normalized LMS Algorithm</i>	501
10.3.4. <i>Transform-Domain LMS Algorithm</i>	502
10.3.5. <i>The RLS Algorithm</i>	503
10.4. MULTIMICROPHONE SPEECH ENHANCEMENT.....	504
10.4.1. <i>Microphone Arrays</i>	505
10.4.2. <i>Blind Source Separation</i>	510
10.5. ENVIRONMENT COMPENSATION PREPROCESSING	515
10.5.1. <i>Spectral Subtraction</i>	516
10.5.2. <i>Frequency-Domain MMSE from Stereo Data</i>	519
10.5.3. <i>Wiener Filtering</i>	520
10.5.4. <i>Cepstral Mean Normalization (CMN)</i>	522
10.5.5. <i>Real-Time Cepstral Normalization</i>	525
10.5.6. <i>The Use of Gaussian Mixture Models</i>	525

10.6. ENVIRONMENTAL MODEL ADAPTATION	528
10.6.1. <i>Retraining on Corrupted Speech</i>	528
10.6.2. <i>Model Adaptation</i>	530
10.6.3. <i>Parallel Model Combination</i>	531
10.6.4. <i>Vector Taylor Series</i>	535
10.6.5. <i>Retraining on Compensated Features</i>	537
10.7. MODELING NONSTATIONARY NOISE.....	538
10.8. HISTORICAL PERSPECTIVE AND FURTHER READING	540
11. LANGUAGE MODELING	545
11.1. FORMAL LANGUAGE THEORY	546
11.1.1. <i>Chomsky Hierarchy</i>	547
11.1.2. <i>Chart Parsing for Context-Free Grammars</i>	549
11.2. STOCHASTIC LANGUAGE MODELS	554
11.2.1. <i>Probabilistic Context-Free Grammars</i>	554
11.2.2. <i>N-gram Language Models</i>	558
11.3. COMPLEXITY MEASURE OF LANGUAGE MODELS	560
11.4. N-GRAM SMOOTHING.....	562
11.4.1. <i>Deleted Interpolation Smoothing</i>	564
11.4.2. <i>Backoff Smoothing</i>	565
11.4.3. <i>Class N-grams</i>	570
11.4.4. <i>Performance of N-gram Smoothing</i>	573
11.5. ADAPTIVE LANGUAGE MODELS	574
11.5.1. <i>Cache Language Models</i>	574
11.5.2. <i>Topic-Adaptive Models</i>	575
11.5.3. <i>Maximum Entropy Models</i>	576
11.6. PRACTICAL ISSUES	578
11.6.1. <i>Vocabulary Selection</i>	578
11.6.2. <i>N-gram Pruning</i>	580
11.6.3. <i>CFG vs. N-gram Models</i>	581
11.7. HISTORICAL PERSPECTIVE AND FURTHER READING	584
12. BASIC SEARCH ALGORITHMS.....	591
12.1. BASIC SEARCH ALGORITHMS	592
12.1.1. <i>General Graph Searching Procedures</i>	593
12.1.2. <i>Blind Graph Search Algorithms</i>	597
12.1.3. <i>Heuristic Graph Search</i>	601
12.2. SEARCH ALGORITHMS FOR SPEECH RECOGNITION	608
12.2.1. <i>Decoder Basics</i>	609
12.2.2. <i>Combining Acoustic and Language Models</i>	610
12.2.3. <i>Isolated Word Recognition</i>	610
12.2.4. <i>Continuous Speech Recognition</i>	611

12.3.	LANGUAGE MODEL STATES	613
12.3.1.	<i>Search Space with FSM and CFG</i>	613
12.3.2.	<i>Search Space with the Unigram</i>	616
12.3.3.	<i>Search Space with Bigrams</i>	617
12.3.4.	<i>Search Space with Trigrams</i>	619
12.3.5.	<i>How to Handle Silences Between Words</i>	621
12.4.	TIME-SYNCHRONOUS VITERBI BEAM SEARCH.....	622
12.4.1.	<i>The Use of Beam</i>	624
12.4.2.	<i>Viterbi Beam Search</i>	625
12.5.	STACK DECODING (A [*] SEARCH)	626
12.5.1.	<i>Admissible Heuristics for Remaining Path</i>	630
12.5.2.	<i>When to Extend New Words</i>	631
12.5.3.	<i>Fast Match</i>	634
12.5.4.	<i>Stack Pruning</i>	638
12.5.5.	<i>Multistack Search</i>	639
12.6.	HISTORICAL PERSPECTIVE AND FURTHER READING	640
13.	LARGE-VOCABULARY SEARCH ALGORITHMS.....	645
13.1.	EFFICIENT MANIPULATION OF A TREE LEXICON	646
13.1.1.	<i>Lexical Tree</i>	646
13.1.2.	<i>Multiple Copies of Pronunciation Trees</i>	648
13.1.3.	<i>Factored Language Probabilities</i>	650
13.1.4.	<i>Optimization of Lexical Trees</i>	653
13.1.5.	<i>Exploiting Subtree Polymorphism</i>	656
13.1.6.	<i>Context-Dependent Units and Inter-Word Triphones</i>	658
13.2.	OTHER EFFICIENT SEARCH TECHNIQUES.....	659
13.2.1.	<i>Using Entire HMM as a State in Search</i>	659
13.2.2.	<i>Different Layers of Beams</i>	660
13.2.3.	<i>Fast Match</i>	661
13.3.	N-BEST AND MULTIPASS SEARCH STRATEGIES	663
13.3.1.	<i>N-best Lists and Word Lattices</i>	664
13.3.2.	<i>The Exact N-best Algorithm</i>	666
13.3.3.	<i>Word-Dependent N-best and Word-Lattice Algorithm</i>	667
13.3.4.	<i>The Forward-Backward Search Algorithm</i>	670
13.3.5.	<i>One-Pass vs. Multipass Search</i>	673
13.4.	SEARCH-ALGORITHM EVALUATION	674
13.5.	CASE STUDY—MICROSOFT WHISPER	676
13.5.1.	<i>The CFG Search Architecture</i>	676
13.5.2.	<i>The N-gram Search Architecture</i>	677
13.6.	HISTORICAL PERSPECTIVE AND FURTHER READING	681

PART IV: TEXT-TO-SPEECH SYSTEMS

14. TEXT AND PHONETIC ANALYSIS	689
14.1. MODULES AND DATA FLOW.....	690
14.1.1. <i>Modules</i>	692
14.1.2. <i>Data Flows</i>	694
14.1.3. <i>Localization Issues</i>	696
14.2. LEXICON	697
14.3. DOCUMENT STRUCTURE DETECTION	699
14.3.1. <i>Chapter and Section Headers</i>	700
14.3.2. <i>Lists</i>	701
14.3.3. <i>Paragraphs</i>	702
14.3.4. <i>Sentences</i>	702
14.3.5. <i>Email</i>	704
14.3.6. <i>Web Pages</i>	705
14.3.7. <i>Dialog Turns and Speech Acts</i>	705
14.4. TEXT NORMALIZATION	706
14.4.1. <i>Abbreviations and Acronyms</i>	709
14.4.2. <i>Number Formats</i>	712
14.4.3. <i>Domain-Specific Tags</i>	718
14.4.4. <i>Miscellaneous Formats</i>	719
14.5. LINGUISTIC ANALYSIS	720
14.6. HOMOGRAPH DISAMBIGUATION.....	724
14.7. MORPHOLOGICAL ANALYSIS	725
14.8. LETTER-TO-SOUND CONVERSION	728
14.9. EVALUATION.....	730
14.10. CASE STUDY: FESTIVAL	732
14.10.1. <i>Lexicon</i>	733
14.10.2. <i>Text Analysis</i>	733
14.10.3. <i>Phonetic Analysis</i>	735
14.11. HISTORICAL PERSPECTIVE AND FURTHER READING	735
15. PROSODY	739
15.1. THE ROLE OF UNDERSTANDING	740
15.2. PROSODY GENERATION SCHEMATIC	743
15.3. SPEAKING STYLE.....	744
15.3.1. <i>Character</i>	744
15.3.2. <i>Emotion</i>	744
15.4. SYMBOLIC PROSODY	744
15.4.1. <i>Pauses</i>	745
15.4.2. <i>Prosodic Phrases</i>	747
	749

<i>15.4.3. Accent</i>	751
<i>15.4.4. Tone</i>	753
<i>15.4.5. Tune</i>	757
<i>15.4.6. Prosodic Transcription Systems</i>	759
15.5. DURATION ASSIGNMENT.....	761
<i>15.5.1. Rule-Based Methods</i>	762
<i>15.5.2. CART-Based Durations</i>	763
15.6. PITCH GENERATION	763
<i>15.6.1. Attributes of Pitch Contours</i>	764
<i>15.6.2. Baseline F0 Contour Generation</i>	768
<i>15.6.3. Parametric F0 Generation</i>	774
<i>15.6.4. Corpus-Based F0 Generation</i>	778
15.7. PROSODY MARKUP LANGUAGES.....	783
15.8. PROSODY EVALUATION.....	784
15.9. HISTORICAL PERSPECTIVE AND FURTHER READING	785
16. SPEECH SYNTHESIS	793
16.1. ATTRIBUTES OF SPEECH SYNTHESIS.....	794
16.2. FORMANT SPEECH SYNTHESIS	796
<i>16.2.1. Waveform Generation from Formant Values</i>	797
<i>16.2.2. Formant Generation by Rule</i>	800
<i>16.2.3. Data-Driven Formant Generation</i>	803
<i>16.2.4. Articulatory Synthesis</i>	803
16.3. CONCATENATIVE SPEECH SYNTHESIS	804
<i>16.3.1. Choice of Unit</i>	805
<i>16.3.2. Optimal Unit String: The Decoding Process</i>	810
<i>16.3.3. Unit Inventory Design</i>	817
16.4. PROSODIC MODIFICATION OF SPEECH.....	818
<i>16.4.1. Synchronous Overlap and Add (SOLA)</i>	818
<i>16.4.2. Pitch Synchronous Overlap and Add (PSOLA)</i>	820
<i>16.4.3. Spectral Behavior of PSOLA</i>	822
<i>16.4.4. Synthesis Epoch Calculation</i>	823
<i>16.4.5. Pitch-Scale Modification Epoch Calculation</i>	825
<i>16.4.6. Time-Scale Modification Epoch Calculation</i>	826
<i>16.4.7. Pitch-Scale Time-Scale Epoch Calculation</i>	827
<i>16.4.8. Waveform Mapping</i>	827
<i>16.4.9. Epoch Detection</i>	828
<i>16.4.10. Problems with PSOLA</i>	829
16.5. SOURCE-FILTER MODELS FOR PROSODY MODIFICATION	831
<i>16.5.1. Prosody Modification of the LPC Residual</i>	832
<i>16.5.2. Mixed Excitation Models</i>	832
<i>16.5.3. Voice Effects</i>	834

16.6. EVALUATION OF TTS SYSTEMS	834
16.6.1. <i>Intelligibility Tests</i>	837
16.6.2. <i>Overall Quality Tests</i>	840
16.6.3. <i>Preference Tests</i>	842
16.6.4. <i>Functional Tests</i>	842
16.6.5. <i>Automated Tests</i>	843
16.7. HISTORICAL PERSPECTIVE AND FURTHER READING	844

PART V: SPOKEN LANGUAGE SYSTEMS

17. SPOKEN LANGUAGE UNDERSTANDING	853
17.1. WRITTEN VS. SPOKEN LANGUAGES	855
17.1.1. <i>Style</i>	856
17.1.2. <i>Disfluency</i>	857
17.1.3. <i>Communicative Prosody</i>	858
17.2. DIALOG STRUCTURE	859
17.2.1. <i>Units of Dialog</i>	860
17.2.2. <i>Dialog (Speech) Acts</i>	861
17.2.3. <i>Dialog Control</i>	866
17.3. SEMANTIC REPRESENTATION	867
17.3.1. <i>Semantic Frames</i>	867
17.3.2. <i>Conceptual Graphs</i>	872
17.4. SENTENCE INTERPRETATION	873
17.4.1. <i>Robust Parsing</i>	874
17.4.2. <i>Statistical Pattern Matching</i>	878
17.5. DISCOURSE ANALYSIS.....	881
17.5.1. <i>Resolution of Relative Expression</i>	882
17.5.2. <i>Automatic Inference and Inconsistency Detection</i>	885
17.6. DIALOG MANAGEMENT.....	886
17.6.1. <i>Dialog Grammars</i>	887
17.6.2. <i>Plan-Based Systems</i>	888
17.6.3. <i>Dialog Behavior</i>	892
17.7. RESPONSE GENERATION AND RENDITION	894
17.7.1. <i>Response Content Generation</i>	895
17.7.2. <i>Concept-to-Speech Rendition</i>	899
17.7.3. <i>Other Renditions</i>	901
17.8. EVALUATION.....	901
17.8.1. <i>Evaluation in the ATIS Task</i>	901
17.8.2. <i>PARADISE Framework</i>	903
17.9. CASE STUDY—DR. WHO	906
17.9.1. <i>Semantic Representation</i>	906
17.9.2. <i>Semantic Parser (Sentence Interpretation)</i>	908

<i>17.9.3. Discourse Analysis</i>	909
<i>17.9.4. Dialog Manager</i>	910
17.10.HISTORICAL PERSPECTIVE AND FURTHER READING	913
18. APPLICATIONS AND USER INTERFACES.....	919
18.1. APPLICATION ARCHITECTURE.....	920
18.2. TYPICAL APPLICATIONS	921
<i>18.2.1. Computer Command and Control</i>	921
<i>18.2.2. Telephony Applications</i>	924
<i>18.2.3. Dictation.....</i>	926
<i>18.2.4. Accessibility</i>	929
<i>18.2.5. Handheld Devices</i>	930
<i>18.2.6. Automobile Applications</i>	930
<i>18.2.7. Speaker Recognition.....</i>	931
18.3. SPEECH INTERFACE DESIGN.....	931
<i>18.3.1. General Principles</i>	931
<i>18.3.2. Handling Errors.....</i>	937
<i>18.3.3. Other Considerations</i>	941
<i>18.3.4. Dialog Flow</i>	942
18.4. INTERNATIONALIZATION	943
18.5. CASE STUDY—MiPAD	945
<i>18.5.1. Specifying the Application.....</i>	946
<i>18.5.2. Rapid Prototyping</i>	948
<i>18.5.3. Evaluation</i>	949
<i>18.5.4. Iterations</i>	951
18.6. HISTORICAL PERSPECTIVE AND FURTHER READING	952
INDEX.....	957

Foreword

*R*ecognition and understanding of spontaneous unrehearsed speech remains an elusive goal. To understand speech, a human considers not only the specific information conveyed to the ear, but also the context in which the information is being discussed. For this reason, people can understand spoken language even when the speech signal is corrupted by noise. However, understanding the context of speech is, in turn, based on a broad knowledge of the world. And this has been the source of the difficulty and over forty years of research.

It is difficult to develop computer programs that are sufficiently sophisticated to understand continuous speech by a random speaker. Only when programmers simplify the problem—by isolating words, limiting the vocabulary or number of speakers, or constraining the way in which sentences may be formed—is speech recognition by computer possible.

Since the early 1970s, researchers at AT&T, BBN, CMU, IBM, Lincoln Labs, MIT, and SRI have made major contributions in Spoken Language Understanding Research. In 1971, the Defense Advanced Research Projects Agency (DARPA) initiated an ambitious five-year, \$15 million, multisite effort to develop speech understanding systems. The goals were to develop systems that would accept continuous speech from many speakers, with minimal speaker adaptation, and operate on a 1000-word vocabulary, artificial syntax, and a

constrained task domain. Two of the systems, Harpy and Hearsay-II, both developed at Carnegie Mellon University, achieved the original goals and in some instances surpassed them.

During the last three decades I have been at Carnegie Mellon, I have been very fortunate to be able to work with many brilliant students and researchers. Xuedong Huang, Alex Acero, and Hsiao-Wuen Hon were arguably among the outstanding researchers in the speech group at CMU. Since then, they have moved to Microsoft and have put together a world-class team at Microsoft Research. Over the years, they have contributed standards for building spoken language understanding systems with Microsoft's SAPI/SDK family of products and pushed the technologies forward with the rest of the community. Today, they continue to play a premier leadership role in both the research community and in industry.

This new book, *Spoken Language Processing*, represents a welcome addition to the technical literature on this increasingly important emerging area of Information Technology. As we move from desktop PCs to personal digital assistants (PDAs), wearable computers, and Internet cell phones, speech becomes a central, if not the only, means of communication between the human and machine! Huang, Acero, and Hon have undertaken a commendable task of creating a comprehensive reference that covers theoretical, algorithmic, and systems aspects of the spoken language tasks of recognition, synthesis, and understanding.

The task of spoken language communication requires a system to recognize, interpret, execute, and respond to a spoken query. This task is complicated by the fact that the speech signal is corrupted by many sources: noise in the background, characteristics of the microphone, vocal tract characteristics of the speakers, and differences in pronunciation. In addition, the system has to cope with non-grammaticality of spoken communication and ambiguity of language. An effective system must strive to utilize all the available sources of knowledge—acoustics, phonetics and phonology, lexical, syntactic, and semantic structure of language, and task-specific context-dependent information.

Speech is based on a sequence of discrete sound segments that are linked in time. These segments, called phonemes, are assumed to have unique articulatory and acoustic characteristics. While the human vocal apparatus can produce an almost infinite number of articulatory gestures, the number of phonemes is limited. English as spoken in the United States, for example, contains 16 vowel and 24 consonant sounds. Each phoneme has distinguishable acoustic characteristics and, in combination with other phonemes, forms larger units such as syllables and words. Knowledge about the acoustic differences among these sound units is essential to distinguish one word from another, say, *bit* from *pit*.

When speech sounds are connected to form larger linguistic units, the acoustic characteristics of a given phoneme will change as a function of its immediate phonetic environment because of the interaction among various anatomical structures (such as the tongue, lips, and vocal chords) and their different degrees of sluggishness. The result is an overlap of phonemic information in the acoustic signal from one segment to the other. For example, the words, say, in *tea*, *tree*, *city*, *beaten*, and *steep*. This effect, known as coarticulation, can occur within a given word or across a word boundary. Thus, the word *this* will have very different acoustic properties in phrases such as *this car* and *this ship*.

This book is self-contained for those who wish to familiarize themselves with the current state of spoken language systems technology. However, a researcher or a professional in the field will benefit from a thorough grounding in a number of disciplines, including:

- *Signal processing*: Fourier Transforms, DFT, and FFT
- *Acoustics*: physics of sounds and speech, models of vocal tract
- *Pattern recognition*: clustering and pattern matching techniques
- *Artificial intelligence*: knowledge representation and search, natural language processing
- *Computer science*: hardware, parallel systems, algorithm optimization
- *Statistics*: probability theory, hidden Markov models, dynamic programming
- *Linguistics*: acoustic phonetics, lexical representation, syntax, and semantics

A newcomer to this field, easily overwhelmed by the vast number of different algorithms scattered across many conference proceedings, can find in this book a set of techniques that Huang, Acero, and Hon have found to work well in practice. This book is unique in that it includes both the theory and implementation details necessary to build spoken language systems. If you were able to assemble all the individual material that is covered in the book and put it on a shelf, it would be several times larger than this volume and yet you would be missing vital information. You would not have the material that is in this book that threads it all into one story, one context. If you need additional resources, the authors include extensive references to get that additional detail. *Spoken Language Processing* is very appealing both as a textbook and as a reference book for practicing engineers. Some readers familiar with a specific topic may decide to skip a few chapters; others may want to focus in other chapters. This is not a book that you will pick up and read once from cover to cover, but one you will keep near you for reference as long as you work in this field.

Raj Reddy
Dean, School of Computer Science
Carnegie Mellon University

Preface

*O*ur primary motivation in writing this book is to share our working experience to bridge the gap between the knowledge of industry gurus and newcomers to the spoken language processing community. Many powerful techniques hide in conference proceedings and academic papers for years before becoming widely recognized by the research community or the industry. We spent many years pursuing spoken language technology research at Carnegie Mellon University before we started spoken language R&D at Microsoft. We fully understand that it is by no means a small undertaking to transfer a state-of-the-art spoken language research system into a commercially viable product that can truly help people improve their productivity. Our experience in both industry and academia is reflected in the context of this book, which presents a contemporary and comprehensive description of both theoretic and practical issues in spoken language processing. This book is intended for people of diverse academic and practical backgrounds. Speech scientists, computer scientists, linguists, engineers, physicists, and psychologists all have a unique perspective on spoken language processing. This book will be useful to all of these special interest groups.

Spoken language processing is a diverse subject that relies on knowledge of many levels, including acoustics, phonology, phonetics, linguistics, semantics, pragmatics, and discourse. The diverse nature of spoken language processing requires knowledge in computer science, electrical engineering, mathematics, syntax, and psychology. There are a number of excellent books on the subfields of spoken language processing, including speech recognition, text-to-speech conversion, and spoken language understanding, but there is no single book that covers both theoretical and practical aspects of these subfields and spoken language interface design. We devote many chapters systematically introducing fundamental

theories needed to understand how speech recognition, text-to-speech synthesis, and spoken language understanding work. Even more important is the fact that the book highlights what works well in practice, which is invaluable if you want to build a practical speech recognizer, a practical text-to-speech synthesizer, or a practical spoken language system. Using numerous real examples in developing Microsoft's spoken language systems, we concentrate on showing how the fundamental theories can be applied to solve real problems in spoken language processing.

We would like to thank many people who helped us during our spoken language processing R&D careers. We are particularly indebted to Professor Raj Reddy at the School of Computer Science, Carnegie Mellon University. Under his leadership, Carnegie Mellon University has become a center of research excellence on spoken language processing. Today's computer industry and academia benefit tremendously from his leadership and contributions.

Special thanks are due to Microsoft for its encouragement of spoken language R&D. The management team at Microsoft has been extremely generous to the speech technology group. We are particularly grateful to Bill Gates, Nathan Myhrvold, Rick Rashid, Dan Ling, and Jack Breese for the great environment they have created for us at Microsoft Research. We would also like to thank Bob Muglia and Kai-Fu Lee for their leadership role in Microsoft's speech product development.

Scott Meredith helped us write a number of chapters in this book and deserves to be a co-author. His insight and experience in text-to-speech synthesis enriched this book a great deal. We also owe gratitude to many colleagues we worked with in the speech technology group of Microsoft Research. In alphabetic order, Jim Adcock, Bruno Alabiso, Fil Alleva, Eric Bidstrup, Antonio Bigazzi, Ciprian Chelba, Li Deng, James Droppo, Doug Duchene, Joshua Goodman, Mei-Yuh Hwang, Larry Israel, Derek Jacoby, Li Jiang, Yun-Cheng Ju, David Larson, Kevin Larson, Jingsong Liu, Ricky Loynard, Milind Mahajan, Peter Mau, John Merrill, Yunus Mohammed, Salman Mughal, Mike Plumpe, Scott Quinn, Bill Rockenbeck, Mike Rozak, Kevin Schofield, Roxana Teodorescu, Gina Venolia, Kuansan Wang, Ye-Yi Wang, and Shenzhi Zhang.

In addition, we want to thank Les Atlas, Jeff Bilmes, Alan Black, David Caulton, Eric Chang, Phil Chou, Dinei Florencio, Allen Gersho, Francisco Gimenez-Galanes, Hynek Hermansky, Henrique Malvar, Julian Odell, Mari Ostendorf, Joseph Pentheroudakis, Tandy Trower, and Charles Wayne. They provided us with many wonderful comments to refine this book. Tim Moore, Russ Hall, and Jane Bonnell at Prentice Hall helped us finish this book in a finite amount of time.

Finally, writing this book was a marathon that could not have been finished without the support of our spouses, Yingzhi, Donna, and Phen, during the many evenings and weekends we spent on this project.

*Xuedong Huang
Alex Acero
Hsiao-Wuen Hon
Redmond, WA*