

**CROSS-LANGUAGE ACOUSTIC ADAPTATION
FOR AUTOMATIC SPEECH RECOGNITION**

by

Christoph Nieuwoudt

Submitted in partial fulfilment of the requirements for the degree

Philosophiae Doctor

in the

Faculty of Engineering

UNIVERSITY OF PRETORIA

April 2000

Summary

Speech recognition systems have been developed for the major languages of the world, yet for the majority of languages there are currently no large vocabulary continuous speech recognition (LVCSR) systems. The development of an LVCSR system for a new language is very costly, mainly because a large speech database has to be compiled to robustly capture the acoustic characteristics of the new language.

This thesis investigates techniques that enable the re-use of acoustic information from a source language, in which a large amount of data is available, in implementing a system for a new target language. The assumption is that too little data is available in the target language to train a robust speech recognition system on that data alone, and that use of acoustic information from a source language can improve the performance of a target language recognition system.

Strategies for cross-language use of acoustic information are proposed, including training on pooled source and target language data, adaptation of source language models using target language data, adapting multilingual models using target language data and transforming source language data to augment target language data for model training. These strategies are allied with Bayesian and transformation-based techniques, usually used for speaker adaptation, as well as with discriminative learning techniques, to present a framework for cross-language re-use of acoustic information. Extensions to current adaptation techniques are proposed to improve the performance of these techniques specifically for cross-language adaptation. A new technique for transformation-based adaptation of variance parameters and a cost-based extension of the minimum classification error (MCE) approach are proposed.

Experiments are performed for a large number of approaches from the proposed framework for cross-language re-use of acoustic information. Relatively large amounts of English speech data are used in conjunction with smaller amounts of Afrikaans speech data to improve the performance of an Afrikaans speech recogniser. Results indicate that a significant

reduction in word error rate (between 26% and 50%, depending on the amount of Afrikaans data available) is possible when English acoustic data is used in addition to Afrikaans speech data from the same database (i.e. both sets of data were recorded under the same conditions and the same labelling process was used). For same-database experiments, best results are achieved for approaches that train models on pooled source and target language data and then perform further adaptation of the models using Bayesian or discriminative techniques on target language data only. Experiments are also performed to evaluate the use of English data from a different database than the Afrikaans data. Peak reductions in word error rate of between 16% and 35% are delivered, depending on the amount of Afrikaans data available. Best results are achieved for an approach that performs a simple transformation of source model parameters using target language data, and then performs Bayesian adaptation of the transformed model on target language data.

Keywords: multilingual speech recognition, cross-language acoustic adaptation, Bayesian adaptation, transformation-based adaptation, minimum classification error adaptation

Opsomming

Spraakherkenningstelsels is reeds ontwikkel vir die groot tale van die wêreld, maar vir die meerderheid van tale bestaan daar tans geen groot-woordeskat kontinuespraakherkenningstelsels nie. Die ontwikkeling van 'n groot-woordeskat kontinuespraakherkenningstelsel vir 'n nuwe taal is baie duur, hoofsaaklik omdat 'n groot databasis opgestel moet word om die akoestiek van 'n nuwe taal op robuuste wyse te vervat.

Die tesis ondersoek tegnieke wat die hergebruik van akoestiese inligting van 'n brontaal, waarvoor 'n groot hoeveelheid data beskikbaar is, toe te laat in die implementering van 'n stelsel vir 'n nuwe teikentaal. Die aanname word gemaak dat te min data beskikbaar is vir die teikentaal om 'n robuuste spraakherkenningstelsel mee af te rig, en dat akoestiese inligting in 'n brontaal gebruik kan word om die herkenning van 'n teikentaalstelsel te verbeter.

Strategieë vir die gebruik van akoestiese inligting oor taalgrense heen word voorgestel en sluit in: afrigting op gepoelede brontaal- en teikentaaldata, aanpassing van brontaalmodelle met teikentaaldata, aanpassing van multitaalmodelle met teikentaaldata en transformasie van brontaaldata om teikentaaldata aan te vul vir afrigting van modelle. Hierdie strategieë word met Bayes en transformasie tegnieke, wat gewoonlik vir sprekeraanpassing gebruik word, en diskriminerende afrigtingstegnieke gebruik om 'n raamwerk vir die gebruik van akoestiese inligting oor taalgrense daar te stel. Uitbreidings van bestaande tegnieke word voorgestel om die herkenning van die tegnieke te verbeter vir kruis-taal aanpassing. 'n Nuwe tegniek vir transformasie van variansieparameters en 'n kostegebaseerde uitbreiding van die minimum klassifikasiefout tegniek word voorgestel.

Eksperimente word uitgevoer vir 'n groot aantal benaderings uit die voorgestelde raamwerk vir kruis-taal hergebruik van akoestiese inligting. Relatief groot hoeveelhede Engelse spraakdata word gebruik tesame met kleiner hoeveelhede Afrikaanse spraakdata om die werkverrigting van 'n Afrikaanse herkenningstelsel te verbeter. Die resultate dui aan dat 'n beduidende vermindering in woordfouttempo (tussen 26% en 50%, afhangende van die hoe-

veelheid Afrikaanse data wat beskikbaar is) moontlik is wanneer Engelse data tesame met Afrikaanse data van dieselfde databasis gebruik word (dit wil sê beide datastelle is onder dieselfde toestande opgeneem en dieselfde etiketteringsproses is gebruik). Vir dieselfde databasis eksperimente word die beste resultate bereik vir benaderings wat modelle afrig op gepoelde brontaal- en teikentaaldata, en wat dan verdere afrigting van modelle volgens Bayes of diskriminasiegebaseerde tegnieke uitvoer met slegs teikentaaldata. Eksperimente word ook uitgevoer om die gebruik van Engelse spraakdata van 'n verskillende databasis as die Afrikaanse data te evalueer. Piek verminderings in fouttempo tussen 16% en 35% word gelewer, afhangende van die hoeveelheid Afrikaanse data wat beskikbaar is. Beste resultate word bereik vir 'n benadering wat 'n eenvoudige transformasie van bronmodelparameters uitvoer met gebruik van teikentaaldata, en dan Bayes aanpassing van die getransformeerde model uitvoer met teikentaaldata.

Sleutelwoorde: multitaalspraakherkenning, kruis-taal akoestiese aanpassing, Bayes aanpassing, parameter transformasie, minimum klassifikasiefout aanpassing

Contents

1	Introduction	1
1.1	Speech recognition fundamentals	2
1.2	Multilingual speech recognition systems	5
1.3	Speaker adaptation techniques	6
1.4	Cross-language re-use of acoustic information	9
1.5	Organisation of thesis	12
1.6	Contributions of thesis	13
2	Background	16
2.1	Hidden Markov modelling framework	17
2.1.1	Feature extraction	17
2.1.2	Continuous density hidden Markov models	18
2.1.3	Duration modelling	20
2.1.4	Hidden Markov model training	21
2.1.5	Pattern matching	27
2.2	Multilingual speech recognition	30
2.2.1	Bootstrapping of new target language recognisers	31

2.2.2	Explicitly multilingual systems	33
2.2.3	Cross-language use of acoustic data for new target languages	34
3	Speaker adaptation theory	36
3.1	Background on speaker adaptation	36
3.1.1	Speaker variation	37
3.1.2	Speaker normalisation	38
3.1.3	Modes of applying speaker adaptation	39
3.1.4	Categories of speaker adaptation	40
3.2	Bayesian adaptation	42
3.2.1	Bayes estimators	43
3.2.2	Gaussian density parameter distributions	46
3.2.3	Mixture density HMM parameter distributions	57
3.2.4	Estimation algorithm	65
3.2.5	Prior density estimation	68
3.3	Transformation-based adaptation	72
3.3.1	Linear transformation of the Gaussian mean	73
3.3.2	Variance transformation	79
3.3.3	Non-linear transformation adaptation	85
3.3.4	Transformation for normalisation before training	86
3.4	Combined Bayesian and transformation-based adaptation	88
3.4.1	Linear transformation-MAP	89
3.4.2	MAP-MLLR	91

3.4.3	Comparison of MLLR-MAP and MAP-MLLR	93
3.5	Discussion	94
4	Discriminative learning theory	95
4.1	Discriminative optimisation criteria	96
4.1.1	Corrective training	97
4.1.2	Maximum mutual information (MMI)	98
4.1.3	Minimum error rate	100
4.2	Minimum classification error approach	101
4.2.1	Optimisation criterion	101
4.2.2	Gradient descent optimisation	104
4.2.3	HMM parameter update	104
4.2.4	MCE training for HMMs	109
4.2.5	Applications	111
4.3	Discriminative optimisation of duration modelling parameters	112
4.4	Discriminative optimisation of linear model transformations	114
4.5	Cost-based MCE	117
4.5.1	String-level MCE	118
4.5.2	Incorporating cost into the loss function	119
4.5.3	Estimating cost based on word error	121
4.5.4	Modifying the misclassification measure	125
4.6	Discussion	132

5	Cross-language acoustic adaptation issues	134
5.1	Language and database issues	135
5.1.1	Phonetic inventories and context	136
5.1.2	Labelling conventions	137
5.1.3	Phonetic mapping	138
5.1.4	Database issues	141
5.2	Strategies for using multilingual data sources	142
5.2.1	Data pooling	142
5.2.2	Model combination	143
5.2.3	Model adaptation	143
5.2.4	Combined pooling and adaptation	144
5.2.5	Data augmentation	145
5.2.6	Combined augmentation and adaptation	147
5.3	Cross-language model adaptation issues	147
5.3.1	Bayesian adaptation	148
5.3.2	Transformation-based adaptation	149
5.3.3	Discriminative adaptation using MCE	152
5.4	Discussion	156
6	Cross-language recognition on SUN Speech	157
6.1	The SUN Speech database	158
6.2	Experimental protocol	159
6.2.1	General system setup	159

6.2.2	Phoneme recognition experiments	160
6.2.3	Word recognition experiments	161
6.3	Initial phoneme recognition experiments	162
6.3.1	Overall phoneme recognition performance	162
6.3.2	Individual phoneme recognition performance	163
6.4	Multilingual data pooling	168
6.5	Bayesian adaptation	169
6.5.1	Cross-language model adaptation	170
6.5.2	Cross-language adaptation of variance	172
6.5.3	Data pooling followed by adaptation	173
6.5.4	Pooling-variance parameter adaptation	174
6.5.5	MAP versus MSE estimation	174
6.6	Transformation-based adaptation	177
6.6.1	Cross-language model adaptation	178
6.6.2	Data pooling followed by adaptation	180
6.7	Combined transformation-Bayesian adaptation	182
6.7.1	MLLR-MAP	182
6.7.2	MAP-MLLR	183
6.8	Discriminative adaptation	184
6.8.1	Data pooling followed by adaptation	186
6.8.2	Improving best performing models	188
6.9	Discussion of results	190

7	Cross-language TIMIT - SUN Speech recognition	193
7.1	TIMIT - SUN Speech phonetic mapping	194
7.2	Multilingual data pooling	196
7.3	Bayesian adaptation	197
7.3.1	Adaptation performance	198
7.3.2	Variance parameter adaptation	199
7.3.3	Pooling-adaptation performance	200
7.3.4	Pooling-variance parameter adaptation	201
7.4	Transformation-based adaptation	202
7.5	Combined transformation-Bayesian adaptation	205
7.5.1	MLLR-MAP	205
7.5.2	MAP-MLLR	207
7.6	Discriminative adaptation	208
7.6.1	Data pooling followed by adaptation	209
7.6.2	Improving best performing models	211
7.7	Data augmentation	214
7.8	Augmentation followed by adaptation	216
7.9	Discussion of results	217
8	Conclusion	220
8.1	Future research	224
A	SUN Speech database	225

A.1	Description	225
A.2	Subdivision into training and test sets	226
A.3	Phonetic content and labelling	227
B	TIMIT - SUN Speech phonetic mapping	231
C	MCE update derivations	235
C.1	Mixture weight derivative	235
C.2	Transition probability derivative	236

List of Abbreviations

ANN	Artificial neural network
CBLF	Cost-based loss function
CBMM	Cost-based misclassification measure
CDHMM	Continuous density hidden Markov model
CDR	Connected digit recognition
CMS	Cepstral mean subtraction
CRBMM	Cost and reward-based misclassification measure
DCT	Discrete cosine transform
DTW	Dynamic time warping
EM	Expectation maximisation
ESHMM	Expanded state hidden Markov model
FFT	Fast Fourier transform
GPD	Generalised probabilistic descent
HMM	Hidden Markov model
LDA	Linear discriminant analysis
LVCSR	Large vocabulary continuous speech recognition
MAP	Maximum <i>a posteriori</i>
MAPLR	Maximum <i>a posteriori</i> linear regression
MCE	Minimum classification error
MFCC	Mel-scaled cepstral coefficient
ML	Maximum likelihood
MLLR	Maximum likelihood linear regression
MLP	Multi-layer perceptron
MMI	Maximum mutual information
MSE	Minimum square error
PCA	Principal component analysis
SA	Speaker adaptive
SD	Speaker dependent
SI	Speaker independent
VTLN	Vocal tract length normalisation
WER	Word error rate

List of Symbols

A	a state transition probability matrix
D	the feature dimension
K	the number of mixtures in a state
M	the number of classes/HMMs
N	the number of states in an HMM
R	a Gaussian precision matrix
S	a sample variance
T	the number of time frames in an observation sequence
W	a transformation matrix
X	an observation sequence
a	an HMM transition probability
c	a mixture weight
q	a state sequence
m	the mean vector of a Gaussian prior distribution
r	a Gaussian precision value
v	a parameter of the mixture weight Dirichlet prior distribution
v	the target variance parameter
w	the relative variance of the mean in the prior
x	a feature vector
Λ	the parameters of a set of HMMs
Σ	a Gaussian covariance matrix
Υ	the precision of the covariance Wishart prior distribution
α	a parameter of a gamma distribution
β	a parameter of a gamma distribution
γ	the state/mixture occupancy variable
	the slope of the sigmoid in the MCE loss function
ϵ	the MCE update parameter
ζ	the cost associated with a phoneme misclassification
η	a parameter of the transition probability Dirichlet prior distribution
	the scaling factor in the MCE misclassification measure
θ	the offset in the MCE loss function
κ	the reward associated with a phoneme classification decision
λ	the parameters of an HMM
μ	a Gaussian mean vector
ξ	the transition count variable
σ	a Gaussian variance vector (for diagonal covariance)
v	the extended mean vector
τ	the variance of the mean in the prior (univariate)
ϖ	a learning rate parameter