

9-2006

## Multilingual Phoneme Models for Rapid Speech Processing System Development

Eric G. Hansen

Follow this and additional works at: <https://scholar.afit.edu/etd>

 Part of the Computational Linguistics Commons, and the Language Interpretation and Translation Commons

---

### Recommended Citation

Hansen, Eric G., "Multilingual Phoneme Models for Rapid Speech Processing System Development" (2006). *Theses and Dissertations*. 3515.  
<https://scholar.afit.edu/etd/3515>

---

This Thesis is brought to you for free and open access by the Student Graduate Works at AFIT Scholar. It has been accepted for inclusion in Theses and Dissertations by an authorized administrator of AFIT Scholar. For more information, please contact [AFIT.ENWL.Repository@us.af.mil](mailto:AFIT.ENWL.Repository@us.af.mil).



MULTILINGUAL PHONEME MODELS FOR RAPID  
SPEECH PROCESSING SYSTEM DEVELOPMENT

THESIS

Eric G. Hansen

AFIT/GE/ENG/06-62

DEPARTMENT OF THE AIR FORCE  
AIR UNIVERSITY

**AIR FORCE INSTITUTE OF TECHNOLOGY**

Wright-Patterson Air Force Base, Ohio

APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED.

The views expressed in this thesis are those of the author and do not reflect the official policy or position of the United States Air Force, Department of Defense, or the United States Government.

AFIT/GE/ENG/06-62

MULTILINGUAL PHONEME MODELS FOR RAPID  
SPEECH PROCESSING SYSTEM DEVELOPMENT

THESIS

Presented to the Faculty  
Department of Electrical and Computer Engineering  
Graduate School of Engineering and Management  
Air Force Institute of Technology  
Air University  
Air Education and Training Command  
In Partial Fulfillment of the Requirements for the  
Degree of Master of Science in Electrical Engineering

Eric G. Hansen, B.S.E.E. University of Dayton

September 2006

APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED.

MULTILINGUAL PHONEME MODELS FOR RAPID  
SPEECH PROCESSING SYSTEM DEVELOPMENT

Eric G. Hansen

Approved:

/signed/

21 August 2006

---

Dr. Steven C. Gustafson (Advisor)

date

/signed/

21 August 2006

---

Dr. John M. Colombi (Member)

date

/signed/

21 August 2006

---

Dr. Timothy R. Anderson (Member)

date

/signed/

21 August 2006

---

Dr. Raymond E. Slyh (Member)

date

*Abstract*

Current speech recognition systems tend to be developed only for commercially viable languages. The resources needed for a typical speech recognition system include hundreds of hours of transcribed speech for acoustic models and 10 to 100 million words of text for language models; both of these requirements can be costly in time and money. The goal of this research is to facilitate rapid development of speech systems to new languages by using multilingual phoneme models to alleviate requirements for large amounts of transcribed speech. The GlobalPhone database, which contains transcribed speech from 15 languages, is used as source data to derive multilingual phoneme models. Various bootstrapping processes are used to develop an Arabic speech recognition system starting from monolingual English models, International Phonetic Association (IPA) based multilingual models, and data-driven multilingual models. The Kullback-Leibler distortion measure is used to derive data-driven phoneme clusters. It was found that multilingual bootstrapping methods outperform monolingual English bootstrapping methods on the Arabic evaluation data initially, and after three iterations of bootstrapping all systems show similar performance levels. Applications of this research are in speech recognition, word spotting, information retrieval, and speech-to-speech translation.

### *Acknowledgements*

I would first like to thank Drs. Slyh and Anderson for the many hours spent discussing ideas and implementation issues to allow this research to progress to this point and for helping brainstorm research paths which lead to this work in the first place. Thank you also to Drs. Gustafson and Colombi for timely and valuable feedback during the experimentation time and for review of this thesis.

I would like to thank Dr. Bryan Pellom, formerly of the University of Colorado at Boulder, for his time developing the SONIC speech recognition system and for his time spent discussing the inner workings of SONIC with me and how best to use it to achieve the goals of this research.

I would like to thank Dr. Grant McMillan, my supervisor, for allowing flexible work hours to allow these graduate studies to take place.

I would also like to thank Mr. Brian Ore for work and discussions on the formatting of the GlobalPhone database. Thanks go to Mr. Dave Hoeferlin for keeping the computer systems up and running during these experiments.

Thanks go to Dr. Tanja Schultz for collecting the GlobalPhone database and allowing for its use by the speech community.

Special thanks to my wife who has supported me and been at my side during this entire process.

Eric G. Hansen

## *Table of Contents*

	Page
Abstract . . . . .	iv
Acknowledgements . . . . .	v
List of Figures . . . . .	viii
List of Tables . . . . .	ix
List of Abbreviations . . . . .	xi
I. Introduction . . . . .	1
II. Background . . . . .	3
2.1 Speech Recognition . . . . .	3
2.1.1 Continuous Speech . . . . .	3
2.1.2 Large Vocabulary ASR . . . . .	5
2.1.3 Speaker Independence . . . . .	5
2.2 Units of Speech . . . . .	5
2.2.1 Phonemes in Context . . . . .	6
2.3 Multilingual Research . . . . .	7
2.3.1 International Phonetic Alphabet . . . . .	7
2.3.2 Multilingual vs. Monolingual Speech Recognition	7
2.3.3 Data-Driven Approaches . . . . .	9
2.4 Automatic Speech Recognition . . . . .	9
2.4.1 General Signal Processing . . . . .	9
2.4.2 Mel Frequency Cepstral Coefficients . . . . .	11
2.4.3 Pronunciation Lexicon or Dictionary . . . . .	12
2.4.4 Language Model . . . . .	12
2.5 Hidden Markov Models . . . . .	13
2.5.1 HMM Training . . . . .	14
2.5.2 Baum-Welch Re-estimation . . . . .	15
2.5.3 Viterbi Algorithm and Decoding . . . . .	16
2.6 The SONIC Speech Recognition System . . . . .	16
2.6.1 Speech Detection and Feature Representation .	17
2.6.2 Acoustic Model . . . . .	17
2.6.3 Monophone Acoustic Models . . . . .	20
2.6.4 Triphone Acoustic Models . . . . .	20
2.6.5 Model Adaptation . . . . .	21
2.6.6 Porting SONIC to Other Languages . . . . .	22

	Page
III. Experimental Results . . . . .	24
3.1 Language Inventory . . . . .	24
3.1.1 Details of Languages used in this Research . . . . .	24
3.1.2 Description of the Arabic partition of GlobalPhone	24
3.2 Multilingual Phoneme Set . . . . .	25
3.3 Performance Metrics . . . . .	25
3.3.1 Phoneme Error Rate – Equally Likely Phonemes	28
3.3.2 Word Error Rate . . . . .	28
3.3.3 Phoneme Error Rate – Word Language Model .	29
3.3.4 Phoneme Confusion Matrix . . . . .	29
3.4 Bootstrapping from English Results . . . . .	30
IV. Porting SONIC to Arabic . . . . .	32
4.1 Bootstrapping from English . . . . .	32
4.2 Bootstrapping from IPA-Based Multilingual Phonemes .	32
4.2.1 Building ML-IPA Acoustic Models . . . . .	33
4.3 Bootstrapping from Data-Driven Multilingual Phonemes	33
4.3.1 Kullback-Leibler Distortion . . . . .	34
4.3.2 Data-Driven Phoneme Clusters . . . . .	35
4.3.3 Building the ML-DD Acoustic Models . . . . .	35
4.4 Bootstrapping Results . . . . .	37
4.5 Adapting Multilingual AMs to Arabic . . . . .	40
4.6 Supplementing with IPA and Data-Driven Multilingual Data . . . . .	40
V. Conclusions . . . . .	44
5.1 Review . . . . .	44
5.2 Future Work . . . . .	44
Appendix A. Phoneme Confusion Matrices . . . . .	46
Appendix B. Additional Information . . . . .	66
Bibliography . . . . .	69

## *List of Figures*

Figure	Page
2.1. IPA Chart . . . . .	8
2.2. Block diagram of an automatic speech recognition system. . . . .	10
2.3. A typical left-to-right HMM . . . . .	14
2.4. Process of calculating Mel Frequency Cepstral Coefficients . . . . .	18
2.5. Example HMM sequence for the word “one” . . . . .	18
2.6. Example decision tree for the base phoneme /AA/ . . . . .	21
2.7. Block diagram of bootstrapping an ASR system from an initial language to a target language . . . . .	23
4.1. Monophone AM results for Bootstrapping Experiments . . . . .	39
4.2. Triphone AM results for Bootstrapping Experiments . . . . .	39
4.3. Monophone AM Adaptation PER-ELP Results . . . . .	41
4.4. Triphone AM Adaptation PER-ELP Results . . . . .	41
4.5. Monophone AM Supplementation PER-ELP Results . . . . .	43
4.6. Triphone AM Supplementation PER-ELP Results . . . . .	43

## *List of Tables*

Table	Page
3.1. Multilingual phonemes as listed in the GlobalPhone dictionaries	26
3.2. Continuation of the Multilingual phonemes as listed in the GlobalPhone dictionaries . . . . .	27
3.3. Error rates on four languages in the GlobalPhone database using the bootstrapping from English method and monophone AMs .	31
3.4. Error rates on four languages in the GlobalPhone database using the bootstrapping from English method and triphone AMs . .	31
3.5. Error rates on three languages in the GlobalPhone database as stated in Schultz . . . . .	31
4.1. Data-driven phoneme groupings based on the KL-distance metric for various thresholds . . . . .	36
A.1. Phoneme Confusion Matrix for monophone AM bootstrapped from English AM with zero iterations . . . . .	50
A.2. Phoneme Confusion Matrix comparing differences in monophone AM bootstrapped from English AM with zero iterations to monophone AM bootstrapped from ML-IPA. . . . .	51
A.3. Phoneme Confusion Matrix comparing differences in monophone AM bootstrapped from English AM with zero iterations to monophone AM bootstrapped from ML-DD10. . . . .	52
A.4. Phoneme Confusion Matrix for monophone AM bootstrapped from English AM with three iterations . . . . .	53
A.5. Phoneme Confusion Matrix comparing differences in monophone AM bootstrapped from English AM with three iterations to monophone AM bootstrapped from ML-IPA. . . . .	54
A.6. Phoneme Confusion Matrix comparing differences in monophone AM bootstrapped from English AM with three iterations to monophone AM bootstrapped from ML-DD10. . . . .	55
A.7. Phoneme Confusion Matrix for triphone AM bootstrapped from English AM with zero iterations . . . . .	56

Table	Page
A.8. Phoneme Confusion Matrix comparing differences in triphone AM bootstrapped from English AM with zero iterations to tri-phone AM bootstrapped from ML-IPA. . . . .	57
A.9. Phoneme Confusion Matrix comparing differences in triphone AM bootstrapped from English AM with zero iterations to tri-phone AM bootstrapped from ML-DD10. . . . .	58
A.10. Phoneme Confusion Matrix for triphone AM bootstrapped from English AM with three iterations . . . . .	59
A.11. Phoneme Confusion Matrix comparing differences in triphone AM bootstrapped from English AM with three iterations to tri-phone AM bootstrapped from ML-IPA. . . . .	60
A.12. Phoneme Confusion Matrix comparing differences in triphone AM bootstrapped from English AM with three iterations to tri-phone AM bootstrapped from ML-DD10. . . . .	61
A.13. Phoneme Confusion Matrix comparing differences in monophone AM bootstrapped from English AM with three iterations to this AM supplemented with monophone data from IPA labels. . . .	62
A.14. Phoneme Confusion Matrix comparing differences in monophone AM bootstrapped from English AM with three iterations to this AM supplemented with monophone data from ML-DD5b labels.	63
A.15. Phoneme Confusion Matrix comparing differences in triphone AM bootstrapped from English AM with three iterations to this AM supplemented with triphone data from IPA labels. . . . .	64
A.16. Phoneme Confusion Matrix comparing differences in triphone AM bootstrapped from English AM with three iterations to this AM supplemented with triphone data from ML-DD5b labels. . . .	65
B.1. American English phoneme set used by SONIC . . . . .	66
B.2. Number of words in each language's dictionary . . . . .	66
B.3. Count and average duration of each multilingual phoneme in the GlobalPhone test subset . . . . .	67
B.4. Continuation of the count and average duration of each multilingual phoneme in the GlobalPhone test subset . . . . .	68

*List of Abbreviations*

Abbreviation		Page
ASR	Automatic Speech Recognition . . . . .	3
HMMs	Hidden Markov Models . . . . .	3
AM	Acoustic Model . . . . .	4
LM	Language Model . . . . .	4
IPA	International Phonetic Alphabet . . . . .	7
DD	Data-Driven . . . . .	9
MFCCs	Mel Frequency Cepstral Coefficients . . . . .	11
DCT	Discrete Cosine Transform . . . . .	11
CMS	Cepstral Mean Subtraction . . . . .	12
GMMs	Gaussian Mixture Models . . . . .	13
SAD	Speech Activity Detection . . . . .	17
EM	Expectation Maximization . . . . .	19
SMAPLR	Structural Maximum a Posteriori Linear Regression . . . . .	22
TL	Target Language . . . . .	22
MLLR	Maximum Likelihood Linear Regression . . . . .	22
PER	Phoneme Error Rate . . . . .	23
WER	Word Error Rate . . . . .	23
OOV	Out of Vocabulary . . . . .	25
ELP	Equally-Likely Phonemes . . . . .	28
WLM	Word Language Model . . . . .	29
ML	Multilingual . . . . .	33
KL	Kullback-Leibler . . . . .	34
KLDM	Kullback-Leibler Distance Measure . . . . .	34

# MULTILINGUAL PHONEME MODELS FOR RAPID SPEECH PROCESSING SYSTEM DEVELOPMENT

## I. Introduction

The concept of a machine that is able to understand human speech has been around since before computers were developed. The thought of using a machine to translate from one language to another has been around nearly as long. The computing power now available in portable devices may one day lead to the universal language translators that are currently science fiction.

A key component in automatic speech translation systems first recognizes the foreign speech, and this recognition is the focus of this research. While much money and research has gone into building speech recognition systems in English and other commercially viable languages, there is a need to be able to rapidly build speech recognition systems in other languages as economic and political landscapes change. A recent example of this requirement is the tsunami that occurred in Indonesia in 2004. The U.S. sent troops for aid relief, and the call was put out for tools to help the troops communicate with the local population. At the time, there were no Indonesian language speech recognition systems available, nor would there be any time soon due to the long development time required. The typical process for building a speech recognition system for a given language is to first collect hundreds of hours of speech data, transcribe this speech, and then to collect 10 to 100 million words of text for a language model. There is also a requirement for a pronunciation dictionary, and time is required for algorithm development and for building the speech recognition system. To allow speech recognition systems to be developed rapidly, especially for speech translation purposes, new processes must be investigated that can access the current pool of large speech resources in a few languages, and then share this information by porting to many languages, thereby minimizing the data collection stage of building a speech recognition system.

One recent idea is to use multiple languages for which there are data to support the development of new systems in other languages. Two approaches include *multilingual* speech recognition systems that recognize multiple languages with one set of models and *monolingual* speech recognition systems that draw from a multilingual training space. The second approach is the focus of this research and it considers the various ways that multilingual data can be used to build phoneme models in other languages.

This thesis is organized as follows. First, a background in speech recognition, units of speech, and past multilingual research is given. An overview follows of automatic speech recognition processes, including a description of Hidden Markov Models (used to model the phonemes) and specific details on the SONIC speech recognition system, which is the software used for all modeling experiments. Chapter three discusses the experimental results using the GlobalPhone database, outlines the performance metrics and discusses the baseline experiments. Chapter four discusses various approaches to integrating the multilingual training data to build an Arabic speech recognition system and discusses the results of each experiment. Chapter five contains a summary of the work presented and highlights ideas for future investigation. The Appendix contains the Phoneme Confusion Matrices discussed in Chapter four. These matrices help provide a more thorough analysis of the results than average phoneme error rates provide.

## II. Background

This chapter addresses the main areas of statistical automatic speech recognition (ASR). First, an overview of terms used to describe speech recognition systems and basic concepts is given. Next, “units of speech,” specifically phonemes, are described. Then, previous multilingual research as related to this work is discussed. Finally, the different components of an ASR system are described in detail, including the inner workings of the SONIC speech recognition system, which is used for the experiments in this research. Special attention is given to Hidden Markov Models (HMMs) as they are the most common statistical models used for ASR and are used within the SONIC speech recognition system.

### 2.1 *Speech Recognition*

Recognizing speech is a difficult problem under even the best circumstances. Every person has a different style of speaking (accent and dialect), which of course can vary based on health reasons, emotional reasons, and the meaning of what the speech is supposed to convey. In addition, when background noise, or additional speakers are included, the challenge of speech understanding becomes greater. Finally the phrase *context is everything* can play a big role in speech recognition. The two phrases below contain similar phonetic information, but with different voicings convey two different meanings.

*Recognize Speech*

*Wreck a nice beach*

Humans can easily distinguish between the two phrases, but we rely on years of learning and have adapted into excellent pattern recognizers. A computer must be trained to learn such differences.

*2.1.1 Continuous Speech.* Continuous speech can be either read or spontaneous, but it is much more difficult for automatic speech recognition systems than

isolated word recognition because of coarticulation effects as our articulators (tongue, teeth, lips, mouth, etc.) move from word to word. An audio signal is sampled at some set rate, and a sequence of feature vectors is extracted for a defined window of time. This feature vector sequence is the observation space,

$$O = \{o_1, o_2, \dots, o_N\},$$

where  $N$  is the total number of observation frames. From this sequence, the goal is to determine the most likely word sequence that could have produced this observation sequence. Let  $W$  represent a sequence of hypothesized words uttered by the speaker,

$$W = \{w_1, w_2, \dots, w_m\}.$$

Let  $\hat{W}$  be the most likely sequence of words given the sequence of feature vectors, then

$$\hat{W} = \operatorname{argmax}_w P(W|O).$$

From Bayes' Rule,

$$P(W|O) = \frac{P(O|W)P(W)}{P(O)}.$$

Given  $O$ ,  $P(O)$  is constant for all possible word sequences, so

$$\hat{W} = \operatorname{argmax}_w \frac{P(O|W)P(W)}{P(O)} = \operatorname{argmax}_w P(O|W)P(W).$$

The term  $P(O|W)$  is provided by the Acoustic Model (AM), which is estimated from transcribed speech training data.  $P(W)$  is derived from the Language Model (LM), which characterizes the probability of observing a sequence of words based on prior knowledge and is estimated from textual training data. This model is the basis for modern statistically-based speech recognizers, not rule-based systems. The focus of this research is to investigate different strategies to derive  $P(O|W)$ .

*2.1.2 Large Vocabulary ASR.* Large vocabulary is a term used to describe modern continuous speech ASR systems and is in contrast to “command and control” recognition applications that have a very limited vocabulary. Ten years ago, 1000 words was a large vocabulary. Today, a 65,000 or even 100,000 word vocabulary is common. Careful search strategies are required to traverse this increased space, and trade-offs are ever present between accuracy, vocabulary size, and algorithm speed.

*2.1.3 Speaker Independence.* A *speaker-independent* speech recognition system is capable of recognizing speech from any speaker, including speakers outside the training space. To achieve speaker independence, a wide variety of training data is used to provide a good representation of the “speaker space” including accent and dialect. In building a speaker independent model, accuracy is lost for any one speaker because of model generalization. Speaker-dependent systems are tuned to a specific group of speakers (gender-specific or speaker-specific for example) and can improve speech recognition performance when the models are applied to the proper group. Speaker adaptation techniques can be employed, and given a small amount of data can readily adapt the speaker-independent models to speaker-dependent models for improved performance. A trade-off exists in flexibility and performance.

## **2.2 Units of Speech**

Words are the most natural units of speech on which to build statistical models for speech recognizers, and for small vocabulary tasks whole-word models are sometimes used. However, due to their lack of generality, word models are not ideal for large vocabulary tasks, as they require large amounts of training data. Syllables are a smaller level unit of speech, but syllables also suffer from a lack of generality. A still smaller unit of speech is the phoneme. It is standard procedure to build large vocabulary speech recognizers using statistical models of phonemes, especially of phonemes in various left and right phonetic contexts.

A phoneme is the smallest meaningful contrastive unit in the phonology of a language [27]. The sounds associated with each phoneme usually have some articulatory gesture(s) in common. Each language has its own set of phonemes. For example, English has roughly 50 phonemes, while Turkish has approximately 30 phonemes. Phonemes can be grouped together into two primary categories, consonants and vowels. Consonants can be further broken down into nasals, plosives, fricatives, approximants, trills, and flaps. Vowels are broken down by where the tongue is positioned within the mouth cavity and the shape of the lips that produces a particular vowel.

*2.2.1 Phonemes in Context.* The definition of each phoneme can be described by specific positions of the articulators (tongue, lips, teeth, etc). However, it takes time to move the articulators from one position to the next, and often the proper position for a given phoneme is never reached as the articulators finish one phoneme and are already moving on to produce the next phoneme. In saying the word “happy”, while the /h/ is being spoken it will be influenced by /a/, (known as a *anticipatory coarticulation*), which in turn, is modified as the articulators position themselves for /p/. The /a/ in “happy” will have different characteristics than the /a/ in “cat” even though they are the same phoneme. The degree of coarticulation between two sounds is dependent, but is not restricted to, the interval between adjacent sounds. The phoneme /k/ has a substantial amount of lip rounding when the next sound is round as in “coo”, but if there is a separation between the /k/ and the rounding, as in “clue”, the /k/ is less rounded.

Two types of models are investigated, monophone and triphone models. Monophone models treat each phoneme as an independent sound, and all instances of a given phoneme are treated equally (i.e., the phonemes are context-independent). Triphone models are groups of models for the same base phoneme that differ based on the phonemes that occur before and after the base phoneme, thus taking into account the immediately preceding and following phonetic contexts (i.e., the phonemes are context-dependent).

## 2.3 Multilingual Research

*2.3.1 International Phonetic Alphabet.* The International Phonetic Alphabet (IPA) was developed in the late 1800’s to create a separate symbol for each contrastive sound occurring in the human language. Figure 2.1 shows the latest revision of the IPA [17]. The symbols on the chart are modified Greek and Latin letters. Between the main symbols and the diacritic marks (see the bottom of the chart), all known sounds of the languages of the world can be represented. Looking at the top most table of consonants, one can see the “manner of articulation” in the rows and the “place of articulation” in the columns. The grayed-out sections are judged to be impossible to humanly produce. A closeup of the vowel section of the chart shows the rows defining the position of the tongue at the roof of the mouth, while the columns represent the position of the tongue at either the front or the back inside the mouth cavity. The rest of the symbols on the chart are used to modify the base symbols and hence create a huge range of sounds. If a Russian word is transcribed with the phoneme /n/ and a French word is also transcribed with the phoneme /n/ it can be said that the two words share the same basic sound. However, in reality there can be variations among the same phoneme across languages, so the IPA is really a categorical simplification of the phonetic context of languages which in truth can be more continuous in nature.

*2.3.2 Multilingual vs. Monolingual Speech Recognition.* Much research has been conducted on the topic of multilingual speech recognition [1,3–5,8,14,43–45,47]. This past research utilizes the similarity of sounds across languages and efficiently builds acoustic models that could recognize multiple languages. In all cases there is a trade-off in performance for the flexibility of the multilingual acoustic models, but there is also success in finding sounds across languages that have many characteristics in common. In [26] cross-language approaches augment a new target language with existing source language acoustic information focusing on the adaptation and transformation from source language to target language and resulting in word error

## THE INTERNATIONAL PHONETIC ALPHABET (2005)

### CONSONANTS (PULMONIC)

	Bilabial	Labio-dental	Dental	Alveolar	Post-alveolar	Retroflex	Palatal	Velar	Uvular	Pharyngeal	Epi-glottal	Glottal
Nasal	m	ɱ		n		ɳ	ɲ	ŋ		N		
Plosive	p b	ɸ ɖ		t d		t̪ d̪	c ɟ	k ɣ	q ɣ		ʔ	ʔ
Fricative	ɸ β	f v	θ ð	s z	ʃ ʒ	ʂ ʐ	ç ɿ	x ɣ	x̪ ɣ̪	h ɺ	h̪ ɻ̪	
Approximant		v		ɹ		ɻ	ɻ̪					
Trill	B			r						R		
Tap, Flap		v		t̪		t̪̪						
Lateral fricative				ɬ ɭ		ɬ̪ ɭ̪						
Lateral approximant				l		l̪	ɺ	ɺ̪	ɺ̪̪			
Lateral flap				ɺ		ɺ̪						

Where symbols appear in pairs, the one to the right represents a modally voiced consonant, except for murmured f.

Shaded areas denote articulations judged to be impossible. Light grey letters are unofficial extensions of the IPA.

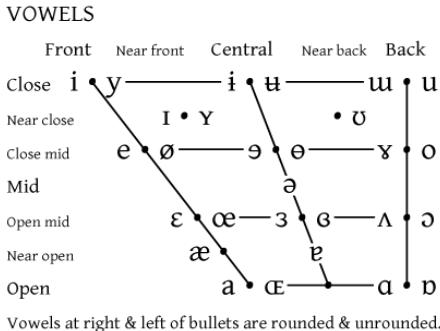
### CONSONANTS (NON-PULMONIC)

Anterior click releases (require posterior stops)	Voiced implosives	Ejectives
ʘ Bilabial fricated	ɓ Bilabial	' Examples:
Laminal alveolar fricated ("dental")	ɗ Dental or alveolar	p' Bilabial
! Apical (post)alveolar abrupt ("retroflex")	ʃ Palatal	t' Dental or alveolar
ǂ Laminal postalveolar abrupt ("palatal")	g Velar	k' Velar
ǁ Lateral alveolar fricated ("lateral")	ɠ Uvular	s' Alveolar fricative

### CONSONANTS (CO-ARTICULATED)

ʍ	Voiceless labialized velar approximant
w	Voiced labialized velar approximant
ɥ	Voiced labialized palatal approximant
χ	Voiceless palatalized postalveolar (alveolo-palatal) fricative
ʐ	Voiced palatalized postalveolar (alveolo-palatal) fricative
ɧ	Simultaneous x and ʃ (disputed)

kp ts Affricates and double articulations may be joined by a tie bar



### SUPRASEGMENTALS

' Primary stress	" Extra stress	Level tones	Contour-tone examples:
, Secondary stress	[fɔuənə'tifən]	é ˥ Top	é ˥ Rising
εː Long	é ˧ Half-long	é ˦ High	é ˨ Falling
ε Short	é ˧ Extra-short	é ˧ Mid	é ˧ High rising
. Syllable break	˨ Linking (no break)	è ˩ Low	è ˩ Low rising
INTONATION		è ˨ Bottom	è ˨ High falling
	˥ Minor (foot) break	Tone terracing	è ˧ Low falling
	˨ Major (intonation) break	↑ Upstep	è ˧ Peaking
	˥ Global rise	˥ Global fall	è ˧ Dipping

**DIACRITICS** Diacritics may be placed above a symbol with a descender, as ɻ̪. Other IPA symbols may appear as diacritics to represent phonetic detail: t̪ (fricative release), b̪ (breathy voice), ՚ (glottal onset), ՚ (epenthetic schwa), ՚ (diphthongization).

SYLLABICITY & RELEASES	PHONATION	PRIMARY ARTICULATION	SECONDARY ARTICULATION						
ɳ ɻ̪	Syllabic	ɳ ɖ	Voiceless or Slack voice	t̪ b	Dental	t̪w d̪w	Labialized	ڙ ڦ	More rounded
ɛ ڦ	Non-syllabic	ڦ ڏ	Modal voice or Stiff voice	t̪ d̪	Apical	t̪̪ d̪̪	Palatalized	ڙ ڦʷ	Less rounded
t̪ʰ h̪ t̪	(Pre)aspirated	ɳ ڦ	Breathy voice	t̪ ڏ	Laminal	t̪v d̪v	Velarized	ڦ ڦ̪	Nasalized
d̪n	Nasal release	ɳ ڦ	Creaky voice	ڦ t̪	Advanced	t̪v d̪v	Pharyngealized	ڦ ڦ̪	Rhoticity
d̪l	Lateral release	ɳ ڦ	Strident	i̪ t̪	Retracted	t̪ z	Velarized or pharyngealized	ڦ ڦ̪	Advanced tongue root
t̪̪	No audible release	ɳ ڏ	Linguolabial	ä̪ ɻ̪	Centralized	ڦ	Mid-centralized	ڦ ڦ̪	Retracted tongue root
ɛ ڦ̪	Lowered (ڦ̪ is a bilabial approximant)	ɛ ɻ̪	Raised (ɻ̪ is a voiced alveolar non-sibilant fricative)						

Figure 2.1: IPA Chart [17].

rate reductions. The approaches investigated in this research are drawn from this background of incorporating multilingual speech data in various approaches to derive new language speech recognition systems. Instead of using the similarity of the multiple languages' phonemes to build a multilingual speech recognizer, the multilingual data is used to help build a monolingual recognizer in a new and (theoretically,) data sparse language.

*2.3.3 Data-Driven Approaches.* To address issues of using categorical based multilingual labeling standards (such as the IPA), data-driven (DD) approaches are investigated so phoneme models across languages can be grouped based on their acoustical properties regardless of phonetic category. In [21] a distance measure between two phoneme models based on a relative entropy-based distance metric is proposed. In [23] phoneme models are grouped based on IPA labels and a log-likelihood approximated distance measure, and experiments are conducted on the SpeechDat database that contains speech from multiple languages. Finally, in [42] a data-driven approach to generate phonetic broad classes is taken using the phoneme confusion matrix to derive the phonetic classes. Results from these papers lead to further investigation into data-driven approaches discussed in this research.

## **2.4 Automatic Speech Recognition**

Figure 2.2 shows an overview of a typical automatic speech recognition system. This section covers each block in the diagram and how each component works to create an ASR system.

*2.4.1 General Signal Processing.* Speech digitally sampled at sixteen kHz results in a usable frequency bandwidth of eight kHz, which covers the majority of acoustic information carried by human speech. However, 16,000 samples per second results in a large number of samples, so the signal is parameterized with a much lower information rate. The signal processing front end extracts important information contained in the speech signal and ideally is designed to show consistent patterns for

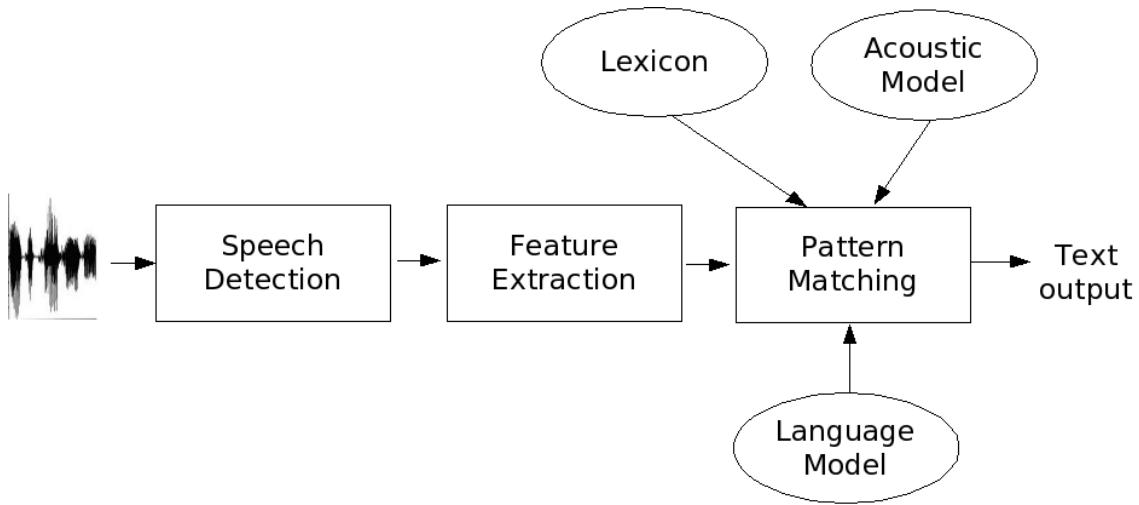


Figure 2.2: Block diagram of an automatic speech recognition system.

the same phonemes across speaker gender, age, accent, and dialect, and is channel independent.

The first step in most speech processing systems *pre-emphasizes* the signal by applying a first order difference equation

$$\tilde{s}(n) = s(n) - as(n - 1)$$

to the samples  $s(n); n = 1, \dots, N$  in each window of samples. In the above equation  $a$  is the pre-emphasis coefficient, which normally is between  $0.9 \leq a \leq 1.0$ . This filtering process increases the energies of the high frequency spectrum to compensate for the approximately -6dB/octave spectral slope of the speech signal, which is mostly attributable to the glottal source (i.e., the airflow through the vocal chords during voiced speech) [16].

The second step windows the data. To avoid end-point problems due to signal truncation for a frame of speech, a weighting window is applied to smooth the end-points. Typical window smoothing options are *Hamming*, *Hanning*, or *Raised Cosine* windows.

**2.4.2 Mel Frequency Cepstral Coefficients.** Mel Frequency Cepstral Coefficients (MFCCs) are the most common feature representation for speech processing. MFCCs are a combination of filter-bank analysis and cepstral analysis. Filter-bank analysis represents the signal spectrum by the log-energies at the output of a filter-bank, where the filters are overlapping band-pass filters spread along the frequency axis. This representation gives a rough representation of the signal spectral shape. The center frequencies of the filters are spread evenly on the Mel scale, which takes into account the relationship between frequency and “perceived” pitch and is related to how the human ear operates. The Mel warping is approximated by

$$f_{Mel} = 2595 \log_{10}\left(1 + \frac{f_{Linear}}{700}\right).$$

The log-energy filter outputs with  $P$  filters are

$$e[j] = \log\left(\sum_{k=0}^{N-1} w_j[k] |S_{Mel}[k]|^2\right) \text{ for } j = 1, \dots, P,$$

where  $w[j]$  represents the  $j^{th}$  filter to the  $k^{th}$  discrete frequency of the sampled signal  $s(n)$  and  $|S_{Mel}[k]|$  represents the DFT magnitude spectrum of  $s(n)$  warped onto the Mel frequency scale.

The “cepstrum” is defined as the inverse Fourier transform of the logarithm of the Fourier transform. Cepstral coefficients represent the spectral envelope of the speech. In practice,  $M$  cepstral coefficients are obtained by decorrelating the filter-bank energies via a Discrete Cosine Transform (DCT)

$$c_t[i] = \sqrt{\frac{2}{P}} \sum_{j=1}^P \left( e[j] \cos \frac{\pi i}{P} (j - 0.5) \right) \text{ for } i = 1, \dots, M.$$

By decorrelating the features, one can easily use diagonal covariance HMMs, which in turn reduces the computational requirements of the models.

The final step to computing the MFCCs is to subtract the mean from each coefficient to account for amplitude differences; this process is referred to as Cepstral Mean Subtraction (CMS).

*2.4.3 Pronunciation Lexicon or Dictionary.* A pronunciation lexicon, sometimes referred to as a dictionary, contains pronunciations for the words in the recognition vocabulary. A pronunciation is a sequence of phonemes used to pronounce the word. The vocabulary of the recognizer is defined by the LM, which is discussed in the next section. The pronunciation lexicon may contain multiple pronunciations for the same word. An example of a word with two pronunciations is

**ACCIDENTAL:** AE K S AX D EH N AX L

**ACCIDENTAL(2):** AE K S AX D EH N T AX L

*2.4.4 Language Model.* A statistical N-gram LM is used to predict word selection for the recognizer output,

$$P(w_1, w_2, \dots, w_k) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)\dots P(w_k|w_1, w_2, \dots, w_{k-1}).$$

For the experiments the LM is a trigram model with back-offs built using the CMU-Cambridge Language Modeling Toolkit [7]. “Trigram” means that word probabilities are derived from the base word and the two-word preceding context:

$$P(W) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2).$$

The term “back-off” refers to the language model backing off to either bigrams:

$$P(W) = P(w_1)P(w_2|w_1)$$

or a unigram:

$$P(W) = P(w_1).$$

An example of a trigram is the word “recognition” given the word “speech” and given the word “automatic” preceding. If a trigram does not occur, the LM reverts (or backs down) to bigrams or unigrams. An example of a bigram is given the word “recognition”, the probability of the word “speech” preceding it. An example of an unigram is the probability of the word “recognition” in the training text [19].

## 2.5 *Hidden Markov Models*

A brief description of Markov processes is presented in this section, followed by descriptions of the training and decoding algorithms for Hidden Markov Models (HMMs). Further details are in [20, 25, 31, 32].

A discrete Markov process has a finite number of states,  $N$ , which form a Markov chain. At discrete time intervals the system may undergo state changes according to a set of transition probabilities, where a transition into the same state is allowed. An  $n^{th}$  order Markov process depends on the current state at time  $t$  as well as all its  $n$  previous states.

Each state of the Markov chain has an associated random output function. These random functions are known as observation distributions or emission distributions and typically are chosen to be Gaussian Mixture Models (GMMs) for ASR. At a discrete time instance,  $t$ , the Markov process is assumed to be in state  $s_t$ , and an observation  $o_t$  is generated according to the emission probability associated with the state. The system may generate a state change or it may stay in the current state at the next time instance. Thus, a Markov model generates a set of observations according to a set of transition probabilities and output distributions. If the state sequence of the underlying Markov chain is unknown or hidden from the observer, this is called an HMM. In speech recognition the state sequence is not directly observed, so HMMs are appropriate. Part of the task of speech recognition is estimation of the underlying state sequence.

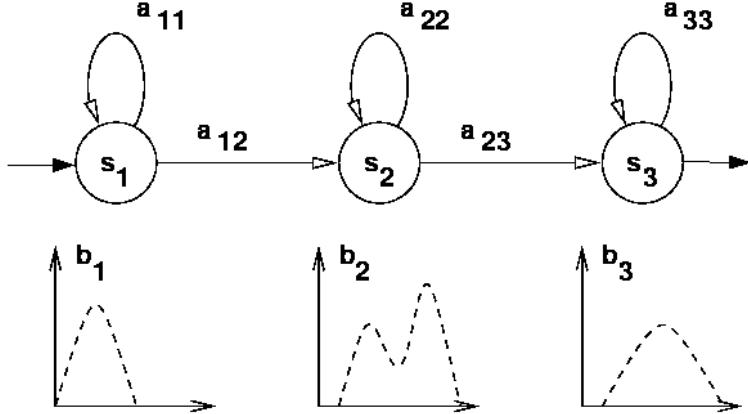


Figure 2.3: Shows a typical left-to-right HMM. The model is entered from the left with an initial transition probability of entering the model in state 1 ( $s_1$ ). The transition probabilities are shown as  $a_{ij}$  where  $i$  is the current state and  $j$  is the next state. The output distributions  $b$  are modeled as Gaussians and would output an observation  $o_j$  when state  $j$  is entered.

There is no limit to the order of the Markov chain, but for speech recognition restriction to a first-order, left-to-right Markov process typically is made. In first-order processes the current state depends only on the immediately preceding state and no other history. This assumption is invalid for speech recognition, but it drastically reduces the complexity of the model and has been shown to give useful results.

Figure 2.3 shows a typical three-state left-to-right HMM. The model is entered from the left with an initial transition probability of entering the model in state 1 ( $s_1$ ). The transition probabilities are shown as  $a_{ij}$  where  $i$  is the current state and  $j$  is the next state. The transition probabilities of the HMMs represent the statistical duration information of the phonemes. The output distributions  $b_i$  are modeled as GMMs and output observation  $o_i$  when state  $i$  is entered.

**2.5.1 HMM Training.** For generative models the model  $\lambda$  that most likely generated an observed sequence of observations  $O = o_1, \dots, o_T$  must be determined. An iterative process known as the forward-backward algorithm is used. The forward

variable probability  $\alpha_t(i)$  is

$$\alpha_t(i) = p(o_1, \dots, o_t, s_t = i | \lambda). \quad (2.1)$$

Equation 2.1 gives the probability of observing the partial sequence of observations  $o_1, \dots, o_t$ , up to time  $t$  and being in state  $i$  at time  $t$  given model  $\lambda$ . For an  $N$ -state model, initially

$$\alpha_1(i) = \pi_i b_i(o_1), \quad 1 \leq i \leq N,$$

where  $\pi_i$  is the probability of starting in state  $i$  at time  $t$ . For the left-to-right model,  $\pi_1 = 1$ . By iterating and summing, a trellis of forward probabilities is generated for state  $j$  at time  $t+1$ :

$$\alpha_{t+1}(j) = \left[ \sum_{i=1}^N \alpha_t(i) a_{ij} \right] b_j(o_{t+1}),$$

where  $1 \leq t \leq T-1$  and  $1 \leq j \leq N$ . The forward probability that model  $\lambda$  generated the full observation sequence  $O$  is

$$P(O|\lambda) = \sum_{i=1}^N a_T(i).$$

A backward variable  $\beta_t(i)$  is defined in a similar fashion,

$$\beta_t(i) = p(o_{t+1}, \dots, o_T, s_t = i | \lambda)$$

or

$$\beta_t(i) = \sum_{j=1}^N a_{ij} b_j(o_{t+1}) \beta_{t+1}(j),$$

*2.5.2 Baum-Welch Re-estimation.* A procedure is needed for adjusting the parameters of a model given the training data. An iterative method known as the Baum-Welch algorithm is commonly used for this purpose. The re-estimation formulas of the Baum-Welch algorithm provide a method of recomputing the transition and emission probabilities of a HMM using the forward and backward probability equa-

tions shown in the previous section. Every iteration of the Baum-Welch re-estimation is guaranteed to increase the total likelihood of the model, (if using Gaussian Mixture output distributions), generating the observation sequence,  $p(O|\lambda)$ , unless a maximum is reached, at which point the likelihood remains constant. The proof of this property is in [25].

**2.5.3 Viterbi Algorithm and Decoding.** For decoding purposes, given an observed acoustic sequence it is necessary to find the maximum likelihood path (i.e., the best state sequence) through a composite model, which is a concatenation of phoneme models (HMMs) into a larger model (HMM network) constrained by the lexicon (dictionary) and the LM. Due to the high computational load of the forward probability calculations that find the best state sequence, a dynamic programming algorithm called the Viterbi algorithm [11] is typically used. The likelihood is computed using the forward probabilities, except that the summation is replaced by an *argmax* operation,

$$\phi_t(j) = \operatorname{argmax}_i \{\phi_{t-1}(i)a_{ij}\}b_j(o_t),$$

where  $\phi_t(j)$  is the highest score along a path at time  $t$ , which accounts for the first  $t$  observations and ends in state  $j$ . The maximum likelihood approximation  $\hat{p}(O|\lambda)$  is then given by  $\operatorname{argmax}_i \{\phi_T(i)\}$  for an  $N$ -state model. A trace back through the trellis created by the Viterbi algorithm, selecting state  $j$  with the highest  $\phi$  at time  $t$ , yields the best state sequence (in reverse order).

## 2.6 The SONIC Speech Recognition System

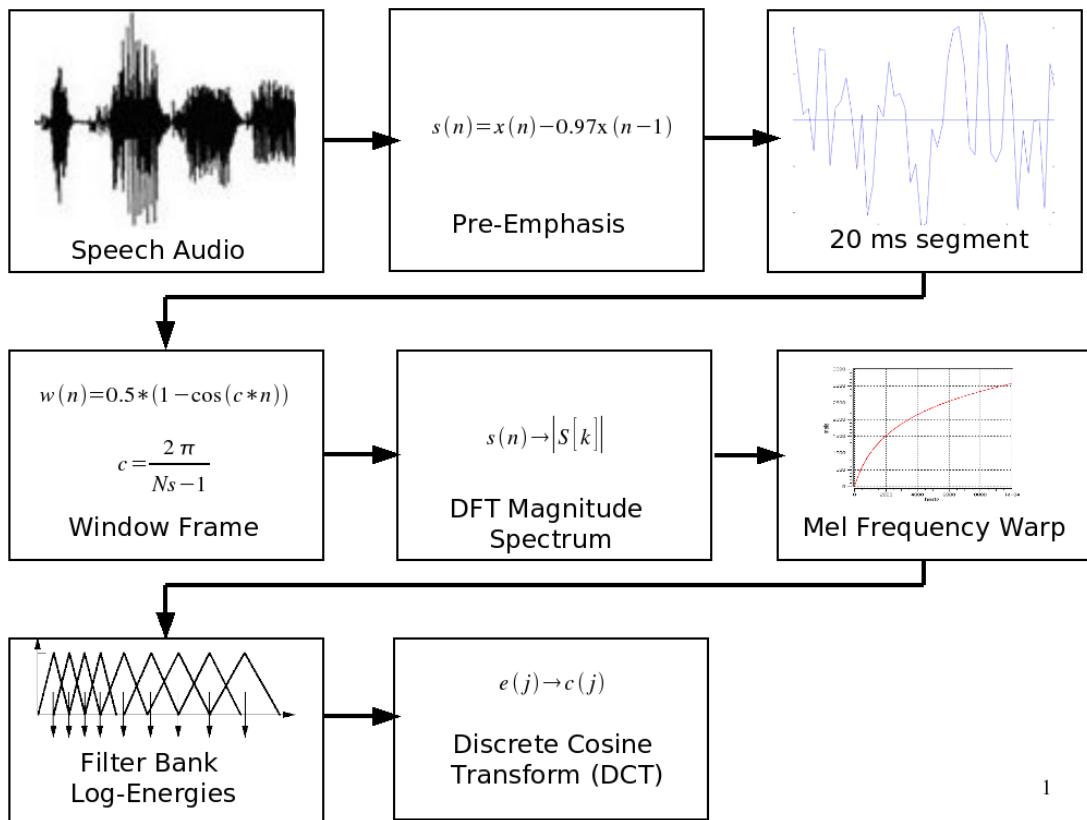
SONIC is the University of Colorado’s Continuous Speech Recognizer. It is designed for research and development of new algorithms for continuous speech recognition. The rest of this section describes in detail each of the components in SONIC. All experiments discussed here were run using SONIC 2.0 beta5 [28].

*2.6.1 Speech Detection and Feature Representation.* Speech detection is sometimes referred to as Speech Activity Detection (SAD) and is a process that marks regions of a speech utterance according to whether they are speech or non-speech (silence, background noise, cough, etc.). Non-speech segments are ignored for recognition. Within SONIC, the SAD is built around a two-state HMM with one state representing speech and the other state representing non-speech. These HMMs are pre-trained and are held constant through all experiments. For both SAD and ASR, SONIC parameterizes the speech into a feature vector.

SONIC computes a 39-dimensional feature vector consisting of 12 MFCCs and the normalized frame energy along with the first- and second-order derivatives of the features for both modeling and decoding stages. The feature vector is calculated every 10 ms from a sliding window of 20 ms of audio. A block diagram of the feature extraction process is shown in Figure 2.4. In general, the process involves:

- Pre-emphasize the signal by a factor of 0.97
- Window 20ms of speech with a raised cosine window
- Compute the FFT
- Warp the frequencies to the Mel scale
- Perform filter-bank analysis (Log of output filter energies)
- Compute the Discrete Cosine Transform (decorrelate features)
- Subtract the Cepstral mean (reduces channel mismatch)
- Compute signal energy of the 20ms speech window
- Compute the first and second order derivatives of all features

*2.6.2 Acoustic Model.* The AM contains all information related to phonetics and channel condition (telephone, microphone, etc.). Each phoneme is represented by a three state continuous-density HMM. The models are continuous because of the underlying GMMs for each state. The three states of each HMM represent the



1

Figure 2.4: Process of calculating Mel Frequency Cepstral Coefficients (MFCCs).

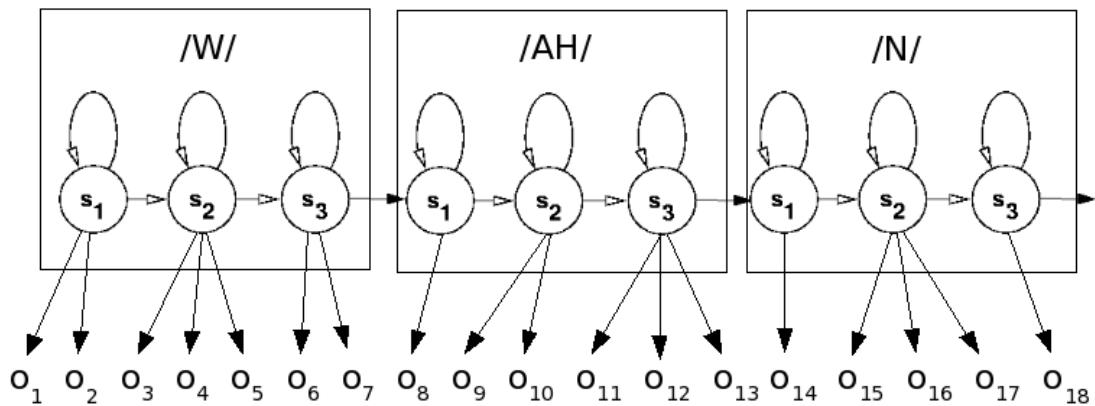


Figure 2.5: Example HMM sequence for the word “one” Each phoneme HMM has three states representing the beginning, middle, and end of the phoneme. Here  $O = o_1, \dots, o_{18}$  for 18 output observation vectors.

beginning, middle, and end of each phoneme. Figure 2.5 shows an example of an HMM network for the word “one.” Both male and female genders have their own set of phoneme models. The phoneme HMMs can be either monophone or triphone models. State transitions in SONIC are modeled using a two-parameter gamma distribution [30] rather than the typical HMM state transition probabilities which by default have a geometric distribution. As reported in [2, 18, 33], using explicit state duration models with HMMs improves recognition accuracy. In [30], it is shown that the gamma distribution is used to fit the measured phoneme duration distribution, and this outperforms the standard methods of using HMM transition probabilities.

Each HMM state within SONIC is represented by a mixture of  $M$  Gaussian distributions. In all experiments, the number of Gaussian distributions is fixed at 32 per HMM state. Each 39-dimensional Gaussian is represented by a weight  $w_m$ , a mean vector  $\mu_m$ , and a covariance matrix. After processing the feature vector through the DCT, a diagonal covariance matrix  $\sigma_m^2$  is assumed. The likelihood calculation is

$$p(o_t|\lambda) = \sum_{m=1}^M \frac{w_m}{(2\pi)^{\frac{D}{2}} \sqrt{\prod_{d=1}^D \sigma_m^2}} \exp \left( -\frac{1}{2} \sum_{d=1}^D \frac{(o_t - \mu_m)^2}{\sigma_m^2} \right),$$

where  $p(o_t|\lambda)$  is the likelihood of observing the  $t^{th}$  feature vector  $o_t$  given the model  $\lambda$ . As discussed,  $D$  is 39 in these experiments. Also, the weights, sum to 1:

$$\sum_{m=1}^M w_m = 1.$$

SONIC actually uses a Viterbi-based training algorithm. During model training, the frames of training data for each base phoneme are clustered to the HMM states by an Expectation-Maximization (EM) algorithm. The trainer code estimates HMM parameters in the maximum likelihood sense one base phoneme at a time using 10 EM iterations [28]. This method is not as thorough as the Forward-Backward algorithm discussed in Section 2.5.1, but Viterbi-based training is fast and efficient.

The training frames for each phoneme are defined by a Viterbi-based alignment process. Recall, the reference transcripts only exist at the utterance level, but SONIC uses the initial AM, the dictionary, and the transcripts to *align* the phoneme labels with the audio data. This alignment process is repeated at each iteration (of either bootstrapping or adaptation techniques) with the updated AMs.

**2.6.3 Monophone Acoustic Models.** Monophone AMs, also known as context-independent models, model all instances of a given phoneme. The resulting monophone model is a three-state, 32-mixture HMM trained on all available frames of data for the particular phoneme being modeled.

**2.6.4 Triphone Acoustic Models.** Triphone AMs, also known as context-dependent models, are clustered depending on the phoneme immediately to the left and right of the phoneme being modeled. The resulting triphone model consists of multiple three-state, 32-mixture HMMs, where each HMM has a list of phonemes for the left context and a list of phonemes for the right context. The training data for each HMM consists of only the training frames that fit the context.

An example of the difference between monophone and triphone AMs is as follows: Given the following three words,

**cat:** K AE T

**hat:** H AE T

**scat:** S K AE T

build a monophone and triphone AM for the phoneme /AE/. The monophone AM for /AE/ is trained on all feature vectors within the start and end times of the phoneme /AE/ for the three words, and results in a single three-state, 32-mixture HMM modeling /AE/. The triphone AM contains two three-state, 32-mixture HMMs with one HMM trained on the feature vectors of /AE/ that fall between /K/ and /T/, and another HMM trained on the feature vectors of /AE/ that fall between /H/ and /T/.

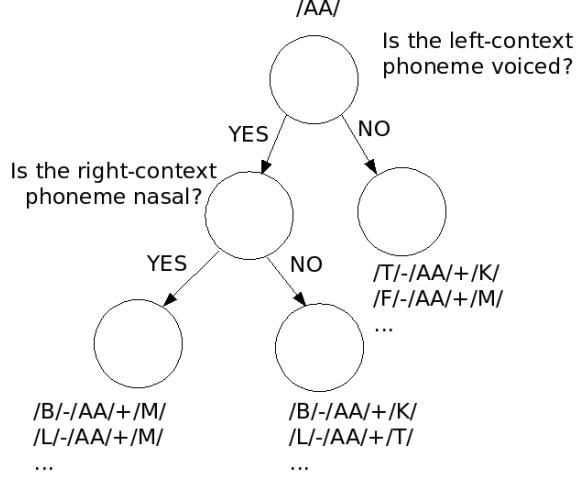


Figure 2.6: Example decision tree for the base phoneme /AA/ [28].

For a language that has 50 phonemes, there are  $50^3$  possible triphone units. However, not every context must appear, and sometimes there is not enough training data to properly model certain contexts. Typically, an automatic clustering method is used to reduce a triphone model to 5000-6000 triphone units. SONIC uses a decision tree method to cluster the contexts into similar acoustical characteristics. A binary decision tree is automatically derived based on certain rules such as which phonemes are *voiced*, which are *nasals*, etc. At each node of the tree these questions are asked, and the splitting continues until either the change in likelihood due to splitting is below a pre-determined threshold or until the number of frames assigned to the leaf node falls below a pre-determined minimum frame count. Feature vectors assigned to the leaf nodes are then used to estimate the actual HMM parameters. An example decision tree is shown in Figure 2.6. A theoretical discussion of decision tree methods for training acoustic models is in [24, 29].

**2.6.5 Model Adaptation.** Model adaptation techniques modify the model parameters based on new data without discarding previously trained models. Two reasons to adapt models are highlighted here. First, adapting speaker independent models to gender dependent models has been shown to improve recognition performance. Second, language adaptation may be realized, where for example, a multilin-

gual model is adapted to a monolingual model. Depending on how the adaptation is implemented, the AMs can adapt too quickly to the new data (similar to disregarding the previous model), or not adapt quickly enough and thus never adequately model the new data.

During the normal process of building AMs in SONIC, a gender-independent AM is built using all the alignments of the training data. Gender-dependent AMs are then adapted from this gender-independent AM by updating the weights and means of the HMMs using the gender-specific training data. If the gender of the test file is known (which can be found by manual or automatic means), the corresponding gender-specific AM is used for the decoding of the audio, and this procedure has been shown to yield improved performance compared to a gender-independent AM. Gender-dependent AMs are used for all the experiments discussed in this research.

The Structural Maximum a Posteriori Linear Regression (SMAPLR) method described by Siohan et al. [40,41] is used to adapt the initial AM to the target language (TL) AM. SMAPLR, (like Maximum Likelihood Linear Regression, MLLR) [12]), estimates a set of regression class transformations to maximize the likelihood of the adaptation data against the existing HMM model. SMAPLR and MLLR differ in that the number of regression transforms as determined by the amount of data and the phonetic content of the adaptation data for SMAPLR, where-as for MLLR the number of regression transforms are determined by the user. The regression class transforms in the implementation of SMAPLR in SONIC are used to transform the Gaussian mean parameters and to adjust the variances of the Gaussian distributions.

*2.6.6 Porting SONIC to Other Languages.* Figure 2.7 shows an example of the steps needed to bootstrap a SONIC ASR system from one language to another. The first step uses an initial AM from a *large resource* language, such as English. The second piece of information needed for bootstrapping is a *Phoneme Map File*. The Phoneme Map File is used to map TL phonemes to the nearest initial language phonemes so that the AMs of the initial language can be used to create an initial

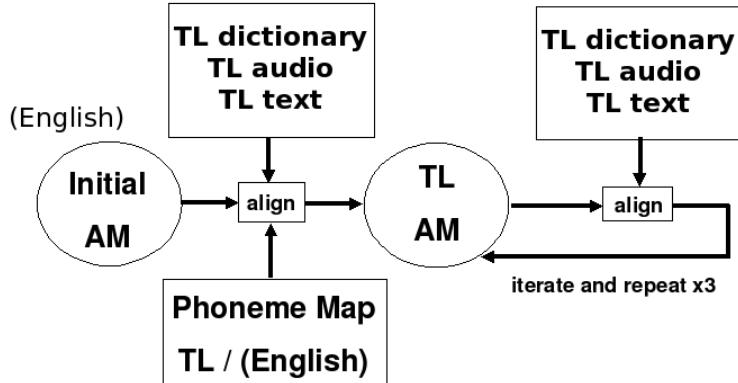


Figure 2.7: Block diagram of bootstrapping an ASR system from an initial language to a target language (TL). Here the initial AM is assumed to be English, hence the phoneme map file contains mappings from the TL to English. The TL dictionary, audio, and text are needed for each training iteration. After each iteration of training the new AM is used to re-align the TL audio to create the next (better) iteration of bootstrapping.

alignment of the TL speech. The initial alignments are used to build the initial TL AMs, which are then used to realign the text in an iterative fashion. Mappings are typically created manually based on IPA. Tables 3.1 and 3.2 show the mappings used in this research for several TLs to English. It is allowable to have a single initial language phoneme map to multiple TL phonemes (as in /a/ and /al/ in German both mapping to English /AA/). After the first alignment of the acoustic data, the TL phoneme models are all distinct. Once these two components are in place, the TL dictionary is used in conjunction with the Viterbi algorithm to time align the audio data with the transcripts. These alignments are then used to train the TL AM for iteration zero. With these new models the Viterbi alignment is repeated using the same TL dictionary, audio data, and transcript files as before. After each iteration of these steps, the previous TL AM is discarded and the new TL AM is trained with the updated alignments. After three iterations of this process, the TL AM is considered final, and this is the model used for decoding the test set data, which results in Phoneme Error Rate (PER) and Word Error Rate (WER) values.

### III. Experimental Results

#### 3.1 Language Inventory

The GlobalPhone Database consists of 15 different languages of high quality read speech data with the source text being news websites [34]. The number of speakers vary with language, and the amount of speech data varies from 16 to 33 hours per language. Because the audio data is read from on-line texts, the transcriptions are extracted from the original text. Transcripts are available at the utterance level only. Word and phoneme alignments must be derived by automatic means.

*3.1.1 Details of Languages used in this Research.* A subset of the GlobalPhone languages is chosen for this research to include Croatian, German, Japanese, Turkish, and Arabic. For each of these languages an ASR system is ported using SONIC according to the method of Section 2.6.6 starting with English AMs, but the Arabic partition is the focus for multilingual porting methods. Narrowing the experiments to Arabic allows five diverse languages to be used for training purposes — namely, Croatian, English, German, Japanese, and Turkish. This subset of languages is chosen from GlobalPhone as being somewhat diverse and yet with good (but not total) coverage of all the phonemes of Arabic. Another reason for choosing this subset of languages (with the exception of English, which is not included in GlobalPhone) is that each of these languages has a pronunciation dictionary in the IPA notation, which allows the proposed multilingual experiments.

*3.1.2 Description of the Arabic partition of GlobalPhone.* The Arabic training data consists of 14.5 hours of speech data spread across 68 different speakers (39 female, 29 male). The evaluation set has 2 hours of speech data (589 utterances) spoken by five different speakers (2 female, 3 male). The dictionary and transcripts are in a Romanized format. An example is below.

Example transcript: *KaARiThTei AL-BuWYiNGh NYuWYuWRK*

Translation: “Collapse of the Boeing in New York” (a headline phrase).

Example dictionary entries:

**KaARiThTei**: k a al r i T t i

**AL-BuWYiNGh**: al l b ul j i n rr

**NYuWYuWRK**: n j ul j ul r k

The pronunciation dictionary contains all words in both training and evaluation sets (i.e., no out-of-vocabulary OOV instances). The dictionary released with GlobalPhone does not account for all the words in the transcripts, but an automatic letter-to-sound rule system is trained according to the procedures of [28] to account for the missing words. Some pronunciations are still missing after this step, so an in-house language expert manually updates such pronunciations.

### ***3.2 Multilingual Phoneme Set***

Each language of the GlobalPhone Database has a dictionary, or pronunciation lexicon, for each of the words of the utterances covered by that section of the database. The dictionaries include an ASCII representation of each phoneme for each language. All of the languages chosen for this research have the pronunciations listed in a multilingual phoneme representation. A list of these phonemes is in Tables 3.1 and 3.2. Each row signifies a multilingual phoneme as listed in the GlobalPhone dictionaries. The last column “ENG” lists a “close” English phoneme mapping in terms of the English phoneme set used by SONIC. See Appendix B for examples of the sounds these English phonemes represent which is taken from [28]. There are some rows where the only multilingual phoneme listed falls under the Arabic column, which means that phoneme is not represented in the four non-English languages used as training data.

### ***3.3 Performance Metrics***

Various performance metrics are used to compare the various acoustic modeling approaches tested. Word and phoneme error rates are the most widely used metrics to compare ASR system performance. However, for phoneme recognition, overall

Croatian	German	Japanese	Turkish	Arabic	English
a	a			a	AA
	al			al	AA
				alal	AA
		ab	ab		AA
		abl			AA
	ae				AE
	aI			aI	AY
	atu				AX
	aU			aU	AW
b	b	b	b	b	B
cp					CH
	C			C	SH
				Cl	SH
d	d	d	d	d	D
				dd	D
dp					JH
dZ		dZ	dZ		JH
				D	DH
e	e	e	e		EH
	el	el			EH
	etu				AX
	eU				OY
f	f	f	f	f	F
g	g	g	g		G
				G	JH
	h	h	h	h	HH
				H	HH
				Hq	EH
i	i	i	i	i	IY
	il	il		il	IY
			i2		IH
j	j	j	j	j	Y
k	k	k	k	k	K
l	l	l	l	l	L
				ll	L
L					L
m	m	m	m	m	M
				ml	M

Table 3.1: Multilingual phonemes as listed in the GlobalPhone dictionaries for Croatian, German, Japanese, Turkish, and Arabic. The phonemes listed for English are from the phoneme set used in the SONIC ASR system.

Croatian	German	Japanese	Turkish	Arabic	English
n	n	n	n	n	N
				nl	N
	ng				NG
nj					N
		nq			NG
o	o	o	o		OW
	ol	ol			OW
	oe		oe		ER
	oel				ER
p	p	p	p		P
				q	K
		Q		Q	SIL
r	r		r	r	R
				rl	R
				rr	G
s	s	s	s	s	S
				sl	S
sj					S
	S	S	S	S	SH
				Sl	SH
			sft		Y
t	t	t	t	t	T
				td	T
ts	ts	ts			TS
tS		tS	tS		CH
				T	TH
u	u		u	u	UW
	ul			ul	UW
	ue		ue		UW
	uel				UW
v	v		v		V
		w		w	W
		W			UH
		Wl			UH
x	x			x	HH
z	z	z	z	z	Z
zj					ZH
			Z	Z	DH

Table 3.2: Continuation of the Multilingual phonemes as listed in the GlobalPhone dictionaries for Croatian, German, Japanese, Turkish, and Arabic. The phonemes listed for English are from the phoneme set used in the SONIC ASR system.

error rates collapse the performance of all the phonemes into one value, sometimes masking individual gains or degradations. To provide a more in-depth look into what the updated AMs contribute towards system performance, the Phoneme Confusion Matrix is also used.

*3.3.1 Phoneme Error Rate – Equally Likely Phonemes.* PER using an equally-likely phoneme “language model” (PER-ELP) means that the system performance is solely due to the AMs. As seen previously in Figure 2.2, the three components needed for the Pattern Matching Block are the LM, AM, and Dictionary. The LM for ELP has equal probability that any phoneme can occur after any other phoneme. Further, the dictionary does not contain words, but is merely a list of phonemes for that language. With these settings, the recognizer output is a stream of phonemes, where the recognizer decision is based solely on the frames of features every 10ms and how closely they match the AMs. This output stream is then compared to the reference phoneme stream (derived from aligning the reference transcripts). The PER takes into account insertion, deletion, and substitution errors. The following example highlights each of these errors:

	<b>The</b>				<b>dog</b>				<b>jumped</b>			
Reference:	DH	AX	–	D	AO	GD	–	JH	AH	M	PD	TD
ASR output:	DH		–	D	AO	G	–	JH	AH	M	PD	TD EH
Error type:	<i>deletion</i>				<i>substitution</i>				<i>insertion</i>			

*3.3.2 Word Error Rate.* WER is another metric used to rate ASR systems. To compute a WER, the ASR system being evaluated has the same AMs as are used to compute the PER, but the LM is made up of word trigrams and the dictionary has the pronunciation for each of the words in the LM. The recognizer output is based not only on the AM information, but also on the statistics of the LM. A well-built LM can mask flaws in poorly trained AMs and vice versa. A problem with WER is the need for exactness. If the reference transcript contains the word “computer” but

the recognizer outputs the word “compute,” an error is listed even though it is wrong by only one phoneme.

*3.3.3 Phoneme Error Rate – Word Language Model.* A second form of PER allows the use of a trained word LM (PER-WLM) as opposed to an equally-likely phoneme LM. The LM is the same when running the recognizer to output word strings. In other words, the recognizer determines the best word given a set of acoustical feature frames and the LM, and that word, along with the pronunciation, are output. In the previous example of the words “computer” and “compute”, one deletion error is returned for the missing phoneme “r”.

*3.3.4 Phoneme Confusion Matrix.* To fully analyze speech recognizer performance, a phoneme confusion matrix is used to display which phonemes an ASR system outputs correctly, and also which phonemes an ASR system confuses with other phonemes. A phoneme confusion matrix labels each row and each column with every valid phoneme for the ASR system employed. The rows represent the actual phonemes (truth) and the columns represent the output phonemes (recognizer output). Therefore, the diagonal of the matrix shows the number of times the recognition output matched the truth for a given phoneme. Reading across a row for a given phoneme, each other cell represents the number of times the recognizer output that wrong phoneme. High values off the diagonal represent highly confusable phonemes and show where work should be focused to improve recognition performance.

Usually, phonemes are grouped into categories such as vowels, nasals, fricatives, etc., so that trends can be analyzed in the errors. One expects a long “a” and a short “a” to show some confusability. One does not expect a long “a” and a “s” to be confused often. If this second trend occurs the phoneme confusion matrix would be the tool used to focus on these errors and narrow system tuning. See Appendix A for phoneme confusion matrices.

### ***3.4 Bootstrapping from English Results***

A preliminary set of experiments conducted using the bootstrapping procedure outlined in Section 2.6.6 starting with English AMs to build phoneme and word recognizers for different languages. The English AMs are trained on the Wall Street Journal Database [13] which consists of 73 hours of English audio data. The first step to the bootstrapping aligns the new language transcriptions to the audio data using the pronunciation dictionary for that language. In this step the English AMs and a Phoneme Map file that assigns the closest English phoneme to each phoneme of the new language (see Tables 3.1 and 3.2 for these mappings) is used for the alignment. Once this initial alignment is completed for each training file, the alignments are used to train three-state, 32-mixture HMMs for each phoneme of the new language. This AM is then used to realign the training files, and these new alignments are used to train new AMs. This process is iterated three times (see Figure 2.7), and the resulting AMs are used for decoding the test data.

Bootstrapping has been successfully used by others to build ASR systems in new languages, see [22, 35–38, 46]. Various methods of bootstrapping exist, especially in how model parameters are updated to build the TL acoustic models and to what extent model parameter sharing is used (if any) between the training language acoustic models and the TL acoustic models.

Tables 3.3 and 3.4 list the PER-ELP, the PER-WLM, and the WER for four languages out of the GlobalPhone Database for both monophone and triphone AMs. The four languages and the amount of training data used are Croatian (12 hours), German (14 hours), Japanese (24 hours), and Turkish (14 hours). (Japanese transcripts are in a Romanized format.) Error rates improve when triphone models are used instead of monophone models, which shows that there is enough training data in the languages in GlobalPhone to build proper triphone AMs.

The error rates are consistent with previous research conducted on the GlobalPhone database [38] as shown in Table 3.5 (Croatian is not evaluated in [38]). Error

	Croatian	German	Japanese	Turkish
PER-ELP	36.4	53.8	35.8	45.1
PER-WLM	19.5	17.1	16.8	3.7
WER	47.3	34.4	47.7	10.9

Table 3.3: Error rates on four languages in the GlobalPhone database using SONIC to bootstrap from English AMs to monophone AMs for that language. All results are after three iterations of alignment (see Figure 2.7).

	Croatian	German	Japanese	Turkish
PER-ELP	33.1	50.0	30.0	40.1
PER-WLM	10.8	11.5	8.8	3.2
WER	32.3	22.3	36.4	10.5

Table 3.4: Error rates on four languages in the GlobalPhone database using SONIC to bootstrap from English AMs to triphone AMs for that language. All results are after three iterations of alignment (see Figure 2.7).

rates in Table 3.5 are for triphone AMs. Differences in error rates between those experiments and those discussed in this research can be attributed to the following: different recognition systems are used, different subsets of GlobalPhone (training and evaluation sets) are used, and different word language models are used. The German PER-ELP is higher than the rest of the languages because the number of German phonemes is 41, which is a much larger pool of phonemes to choose from compared to Croatian with 30 phonemes, Japanese with 31 phonemes, and Turkish with 29 phonemes. The reason the Turkish PER-WLM has low error (3.2% error) is perhaps because of the properties of the Turkish language [6], which result in long words, allowing easier acoustical discrimination (assuming no OOV).

	German	Japanese	Turkish
PER-ELP	44.5	33.8	44.1
WER	11.8	10.0	16.9

Table 3.5: Error rates on three languages in the GlobalPhone database as stated in Schultz [38]. All results are based on triphone AMs. Differences between these results and the results reported for SONIC include: differences in ASR systems, differences in the partitions of the GlobalPhone database, and differences in WLMs.

## IV. Porting SONIC to Arabic

### 4.1 *Bootstrapping from English*

The first set of experiments to build an Arabic speech recognizer started with English AMs that were then bootstrapped into Arabic AMs with the process of Section 2.6.6. The resulting Arabic AMs were used to decode the test data and compute PER-ELP, PER-WLM, and WER values. These results are named “Boot-ENG-0” and “Boot-ENG-3” for results after zero and after three iterations of bootstrapping, respectively.

### 4.2 *Bootstrapping from IPA-Based Multilingual Phonemes*

As indicated in Tables 3.1 and 3.2, there are multiple instances of the same English phoneme mapping to different Arabic phonemes. Some of the mappings are coarse approximations because there are phonemes in Arabic that do not occur in English (Arabic has a “pharyngeal voiced fricative”). The premise behind using multiple languages to train the initial AMs is that such training would provide greater coverage of the Arabic phoneme space so that the initial Arabic alignments and the final Arabic AMs would be more accurate. As mentioned in Section 3.1.1, the languages chosen from the GlobalPhone database have a common multilingual phoneme set based on the IPA. Thus, in aligning the phoneme /f/, one could draw on /f/ AMs from Croatian, German, Japanese, and Turkish. On the other hand, English does not have a phoneme that closely corresponds to the Arabic /x/, but Croatian and German do, so their /x/ models might be useful in aligning the initial Arabic. There are still 17 Arabic phonemes that are not represented by the four chosen training languages. These 17 phonemes are “mapped” to the next closest multilingual phoneme. “Closest” is determined by manner, then voicing, then place for consonants and according to the IPA vowel chart for vowels. Also, note that due to accent and dialect issues, just because all of the languages have an /f/ does not mean that they are all equally close to the Arabic /f/.

*4.2.1 Building ML-IPA Acoustic Models.* The process used to build the initial Multilingual ML-IPA based AMs is as follows. First, all alignment information from each of the four training languages is gathered. This alignment information is from the third iteration of bootstrapping from English to that particular language. The alignment information contains phoneme labels (based on the IPA) and start and ending times that reference the audio file pertaining to the transcript.

To build monophone ML-IPA AMs, all training data across the four training languages are grouped, and a three-state, 32-mixture HMM is trained and used for the initial alignment of the Arabic acoustic data. For the 17 Arabic phonemes that do not exist in Croatian, German, Japanese, and Turkish, the base phoneme ML-IPA AM, as determined by the mapping file, is used.

The procedure used to build triphone ML-IPA AMs is more complicated. Recall that a triphone represents the base phoneme in the context of the preceding and following phoneme. When creating multilingual AMs for the purpose of Arabic speech recognition, one must be careful to only use phonemes that occur in Arabic. If a set of alignments are given to the HMM trainer that contain a non-Arabic phoneme, the resulting HMM model contains a context-dependent model accounting for that non-Arabic phoneme, and SONIC posts an error when attempting to decode with this model because that phoneme does not occur in the Arabic phoneme set. To ensure against this error, the following process is developed. For a given base phoneme, all multilingual training alignment data is grouped together. This pool is then refined by discarding any frames involving triphones that contain contexts with non-Arabic phonemes. This process is repeated for each base phoneme, and the result is a ML-IPA triphone model that exists in the phoneme space of Arabic.

### ***4.3 Bootstrapping from Data-Driven Multilingual Phonemes***

One issue with using the IPA-based multilingual approach is that, for Arabic, there are 17 phonemes that must be mapped from “close” phonemes, where the idea of closeness is categorical in nature and somewhat subjective. A second issue is that just

because several languages have the same basic phoneme according to the IPA, it does not follow that the phonemes all have exactly the same acoustic properties. There are generally differences between like phonemes across different languages. One way to address these issues is to implement a data-driven approach based on a *measure* of the distance between different phonemes from the different training languages and to allow the “close” matching phonemes to group together to form the new phoneme models for the initial Arabic AMs.

*4.3.1 Kullback-Leibler Distortion.* The distance measure chosen is based on the Kullback-Leibler (KL) distortion measure [9,39], which is designed to measure the “difference” between Gaussians. Let  $N(\mu_1, \Sigma_1)$  and  $N(\mu_2, \Sigma_2)$  denote two Gaussian distributions for vectors of length  $d$ , then the KL distortion between them is

$$KL_{dist(N_1, N_2)} = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{trace}(\Sigma_2^{-1} \Sigma_1) + (\mu_1 - \mu_2)^T \Sigma_2^{-1} (\mu_1 - \mu_2) \right].$$

This expression is implemented by training three-state, single-mixture HMMs for each phoneme for each language. By using a single mixture, there is only one 39-dimensional Gaussian for each state. (There is no known closed-form expression for the KL distortion measure for GMMs, although there has been some work on approximations [15].)

The KL distortion measure is not symmetric, (i.e.,  $KL_{dist(N_1, N_2)} \neq KL_{dist(N_2, N_1)}$ ). To form a distance measure, one can average  $KL_{dist(N_1, N_2)}$  and  $KL_{dist(N_2, N_1)}$ , we refer to this as the KL distance measure (KLDM). One other factor is handling the three states (beginning, middle, ending) for each phoneme in order to obtain a single distance measure between two phoneme models. Two approaches are investigated. The first approach equally weights the contributions of the three states, while the second approach weights the contributions of the states 0.25, 0.50, 0.25 under the assumption that the middle state, which is sometimes longer in duration, has fewer coarticulation effects and might be a better indicator of the distance between phoneme models. It is found that both weighting schemes return the same phoneme grouping trends, and

therefore all data-driven phoneme clusters discussed in this research are derived by equally weighting the state distances.

*4.3.2 Data-Driven Phoneme Clusters.* To derive data-driven Arabic phoneme clusters, each single mixture (three-state) Arabic phoneme model is compared using the KLDM to each multilingual phoneme model from each of the four training languages. All models (including Arabic) are those bootstrapped from English after three iterations. Table 4.1 displays the resulting phoneme clusters for four threshold values using the KLDM. Trends in the phoneme groupings show Turkish phonemes dominating the closest distances to Arabic, with Croatian phonemes also matching often. With narrow threshold values (5 and 10), the data-driven phoneme clusters generally match the phonetic content of the Arabic phoneme; thus, the KLDM generally yields results that are intuitive. Relaxing the threshold value to 15, it is seen that “plosives” are grouped together in the example of /b/ and /d/ for AR-/dd/ and /k/, /p/, and /t/ for AR-/k/ and /t/. Also note that the voicing matches in these two examples. Further, AR /m/ includes both nasals /m/ and /n/.

All experiments involving these data-driven multilingual phoneme clusters are referred to in the remainder of this thesis as ML-DD(threshold value) experiments, where the threshold value is that of the KLDM.

*4.3.3 Building the ML-DD Acoustic Models.* The process used to build the ML-DD AMs is similar to that used to build the ML-IPA models, except instead of the IPA labels driving the frame combinations, the phoneme clusters seen in Table 4.1 determine which training frames are grouped.

To build a monophone ML-DD AM for a given phoneme and threshold, all training data corresponding to phoneme models found to be within the threshold (or the closest match when no phoneme models were within threshold) are grouped and used to train a three-state, 32-mixture HMM, which is used for the initial alignment of the Arabic acoustic data.

Arabic phoneme	Threshold			
	5	10	15	20
a	TU-ab*	TU-ab*	TU-ab,CR-a	JA-ab,TU-e,i2
aI	TU-e*	TU-e*	TU-e*	TU-e*
al	TU-i2*	TU-i2	GE-etu,r,TU-r	
alal	TU-ab*	TU-ab*	TU-ab*	TU-ab*
aU	TU-o*	TU-o*	TU-o*	TU-o*
b	CR-b,TU-b		JA-b	CR-d,TU-d
C	TU-S*	TU-S	CR-sj,JA-S	TU-s,tS
Cl	JA-S*	JA-S*	JA-S*	JA-S,TU-S
d	CR-d	TU-d	TU-b	CR-b,g,JA-d,TU-g
D	JA-d*	JA-d*	JA-d*	CR-d,JA-b,d,TU-d
dd	CR-d*	CR-d*	CR-d,JA-b,TU-b,d	CR-b,g,JA-d
f	CR-f,TU-f		GE-f	JA-f
G	TU-Z*	TU-Z	TU-dZ	CR-dp,zj
h	TU-h*	TU-h*	GE-h,TU-h*	
H	JA-h*	JA-h*	JA-h*	JA-h*
Hq	TU-ab*	TU-ab*	TU-ab*	TU-ab*
i	TU-i	CR-i	JA-i,TU-ue	GE-etu,i,JA-W,TU-i2
il	CR-i*	CR-i*	CR-i,TU-i	JA-i,il
j	TU-j*	TU-j		
k	TU-k*	CR-k,TU-k	CR-t,JA-k,TU-p,t	CR-p,GE-g,k
l	TU-l*	GE-l,TU-l	CR-L	CR-n,r,JA-g,l,TU-n,r
ll	CR-L*	CR-L*	CR-L*	CR-L*
m	TU-m	CR-m	CR-n,GE-m,JA-m,TU-n	JA-n,TU-l
ml	CR-m*,TU-m*	CR-m*,TU-m*	CR-m*,TU-m*	CR-m*,TU-m*
n	TU-n	CR-n	CR-m,GE-n,TU-m	CR-nj,JA-n,TU-1
nl	JA-n*	JA-n*	JA-n*	JA-n
q	CR-k*	CR-k	TU-k	JA-k,TU-p
Q	CR-x*	CR-x*	CR-x*	CR-x*
r	TU-r*	CR-r,TU-r		CR-l,GE-l,r,JA-l,TU-i2
rl	CR-r*	CR-r*	CR-r*	CR-r
rr	TU-h*	TU-h*	TU-h	CR-v,JA-g,TU-l,v
s	TU-s	CR-s	GE-s,z,JA-s	TU-z
S	TU-s*	TU-s	GE-s,z,JA-s	TU-z
sl	TU-s*	TU-s		CR-s,JA-s
Sl	TU-s*	TU-s		CR-s,JA-s
t	TU-t	CR-t	GE-t,TU-k,p	CR-k,p,GE-d,JA-k,t
T	TU-f*	TU-f	CR-f	CR-x
td	TU-t*	CR-t,TU-t	CR-p,TU-p	CR-k,GE-t,JA-t,TU-k
u	TU-u	CR-u	CR-o,GE-u	JA-o,TU-i2,o
ul	GE-u*	GE-ol,u	CR-o,u,TU-o,u	GE-o,ul,JA-o
w	CR-l*,u*	CR-l*,u*	CR-l*,u*	CR-l*,u*
x	CR-x*	CR-x*	CR-x*	CR-x
z	CR-z*,TU-z*	CR-z,TU-z		GE-z,JA-z
Z	JA-g*,TU-l*	JA-g*,TU-l*	JA-g*,TU-l*	JA-g,TU-l

Table 4.1: Data-driven phoneme groupings based on the KL-distance metric for various thresholds. The first two letters signify the source language (CR-Croatian, GE-German, JA-Japanese, TU-Turkish), while the letters after the hyphen signify the phoneme label. As the threshold increases, only new phonemes are listed; all previous phoneme groupings still hold. The (\*) signifies that no phoneme is below the given threshold, and the phoneme listed is the closest in distance to the Arabic phoneme.

To illustrate this process for triphones, an example is given first for the base phoneme AR-/i/ with a KLDM threshold of 10 and then discussed in further detail.

Given that the Turkish training data contains the triphone: TU /f/-/i/-/m/, Table 4.1 shows that, TU-/f/ maps to AR-/f/ and AR-/T/ at a KLDM of 10. Also, TU-/m/ maps to AR-/m/ and AR-/ml/ at a KLDM of 10. Therefore, the training frames affiliated with the Turkish sequence /f/-/i/-/m/ are added to the multilingual training data to create the following four Arabic contexts: /f/-/i/-/m/, /f/-/i/-/ml/, /T/-/i/-/m/, and /T/-/i/-/ml/.

Thus, to build triphone ML-DD AMs, steps similar to those used to build the ML-IPA AMs are employed. For each Arabic base phoneme and for a given KLDM threshold, the set of clustered phoneme multilingual training data is input. Each triphone in the multilingual training data is then examined to see if it contains valid Arabic phonemes. If so, the triphone training frames are stored, otherwise the triphone data are discarded. If multiple mappings exist for a training phoneme, that triphone is repeated to cover all mappings.

#### 4.4 Bootstrapping Results

Figure 4.1 shows error rates for the three different metrics used to evaluate the different modeling methods. In all experiments the pronunciation dictionary and language model are held constant; only the AM varied. These error rates are computed using the program *Sclite* [10] and represent the average of the errors for the five Arabic speakers (589 utterances) in the evaluation set.

The first four bar groupings in Figures 4.1 and 4.2 show the results of decoding the Arabic speech with AMs that do not contain any Arabic data. From the PER-ELPs, it is seen that the monophone performance does not vary greatly between AMs, but for the triphones the ML-IPA AM performs 12% better (in an absolute sense) than the best ML-DD AM. For PER-WLM, the ML-DD10 monophone AM performs best, with the other ML-DD monophone AMs not much reduced in performance.

The ML-DD10 triphone AM performs the worst of the triphone models for this error metric. For the WER, the ML-DD10 monophone models performs the best (as with the PER-WLM); however, for triphones, ML-DD20 models performed the best (as for the PER-WLM). Because the monophone ML-DD10 AM returns the best WER out of the three ML-DD experiments, only this KLDM threshold is investigated further with the bootstrapping and adaptation iterations.

The next two groupings of bars in Figures 4.1 and 4.2 show the results of the different error metrics after the initial bootstrap of Arabic data (0 iter) and after three iterations of bootstrapping (3 iter). All three metrics show that the ML-DD10 AM performs best on the initial alignment of the Arabic for both monophones and triphones. After three iterations of bootstrapping the ML-DD10 is still best for monophones but not for triphones; although the difference is only 1.3% worse for the PER-ELP, only 0.2% worse for the PER-WLM, and only 0.8% worse for the WER (all percents in an absolute sense). Overall, after three iterations of bootstrapping, all three modeling procedures tend to converge to similar performance levels. Tables A.1–A.12 show the phoneme confusion matrices for each of these experiments and are discussed in further detail in the Appendix.

There are notable challenges with the version of the transcripts and dictionary used in these experiments which may account for high WER and PER values. The first is the added complexity in the dictionary of the prefix *AL* (Arabic for “the”) attached to many words. Entries exist for both the root word and the root word with the prefix attached. Also, sometimes *AL* is written as *AeL*. By adding these entries, there is a larger search space and a greater chance of error. The second issue with the transcripts is their inconsistency. With the aforementioned discussion about dictionary entries with and without prefixes, the transcripts are found to sometimes have the prefix attached and sometimes not. Also, in the transcripts certain phrases (perhaps Arabic colloquialisms) have multiple words grouped together by “underbars”. The phrases are left “as is”, and dictionary entries are created to account for these “words”.

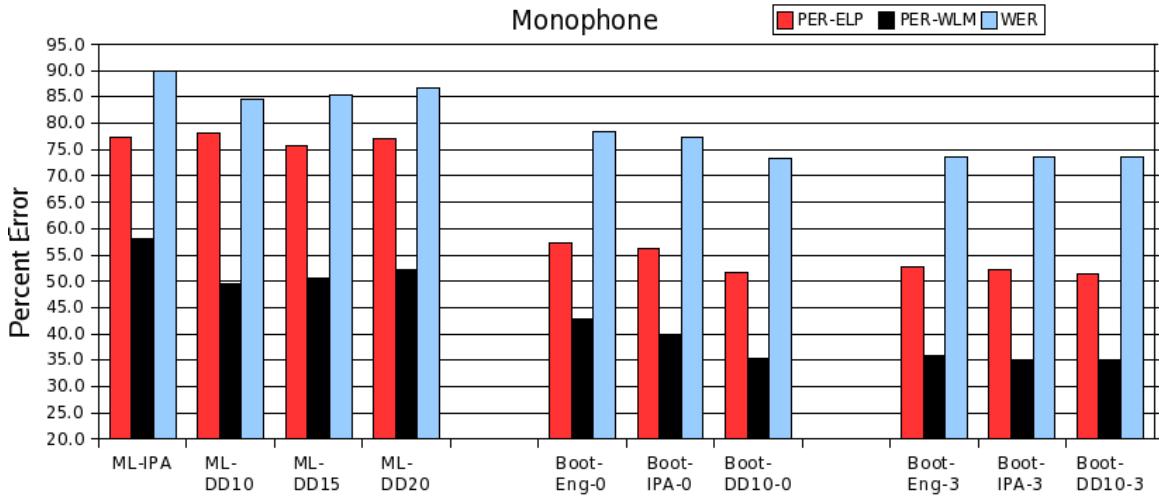


Figure 4.1: Monophone AM results on Arabic test data for Bootstrapping Experiments. The first grouping shows results of using AMs with no Arabic acoustic data, the second grouping is after an initial alignment and training stage of Arabic data, and the third grouping is after three iterations of bootstrapping.

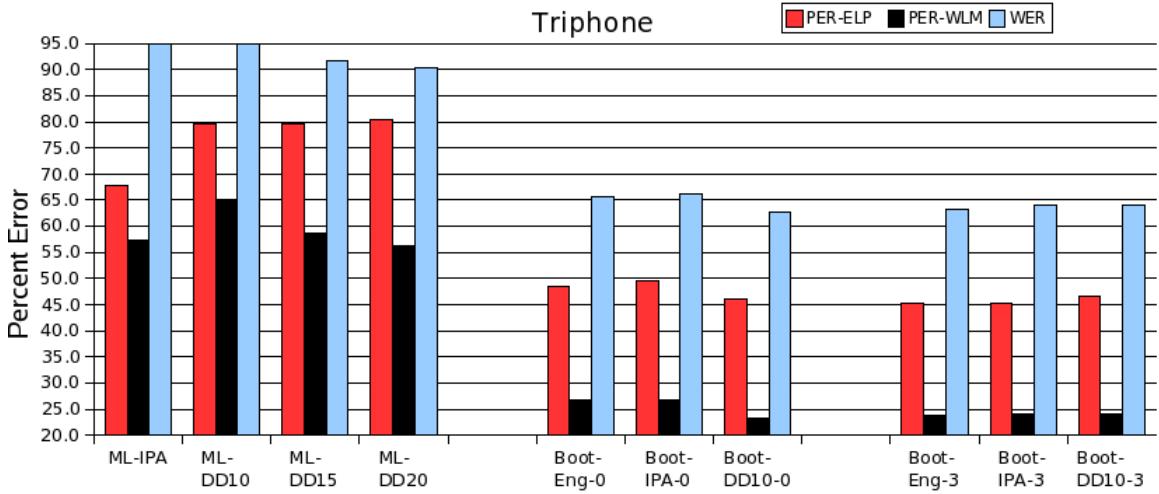


Figure 4.2: Triphone AM results on Arabic test data for Bootstrapping Experiments. The first grouping shows results of using AMs with no Arabic acoustic data, the second grouping is after an initial alignment and training stage of Arabic data, and the third grouping is after three iterations of bootstrapping.

These phrases do not occur repeatedly in the training data, and therefore affect the LM statistics.

#### ***4.5 Adapting Multilingual AMs to Arabic***

The motivation behind trying adaptation as opposed to bootstrapping is the fact that bootstrapping discards all previous AMs and only uses the most recent model trained on whatever acoustic data is available for the given target language. In the case of small resource languages, if the initial AM has a reasonable amount of coverage of the target language phoneme set, adaptation might tune the acoustic model while still allowing generalization of the initial model trained on larger amounts of data.

Using the SMAPLR adaptation scheme as discussed in Section 2.6.5, an alternative porting method to bootstrapping is evaluated. An initial multilingual AM is *adapted* with training data from the Arabic training set. Initial AMs were either ML-IPA or ML-DD10 AMs. English AMs were not used because of difficulties getting access to the exact alignments that went into the English models. The SMAPLR approach automatically creates regression class transformations and adjusts the mean and variance values of the HMM components. The phoneme duration parameters (represented by a gamma distribution) are not modified.

As indicated in Figures 4.3 and 4.4, multiple iterations of adaptation improve ASR performance; however, after three iterations the adapted AMs do not outperform the bootstrapped AMs. Since adaptation takes substantially more computational power and time than bootstrapping, no further model adaptation experiments are conducted.

#### ***4.6 Supplementing with IPA and Data-Driven Multilingual Data***

A last set of experiments investigated *supplementing* the Arabic data with multilingual data. The motivation is that if some multilingual data matches close enough to the Arabic data, it might be utilized to increase the amount of training data for the

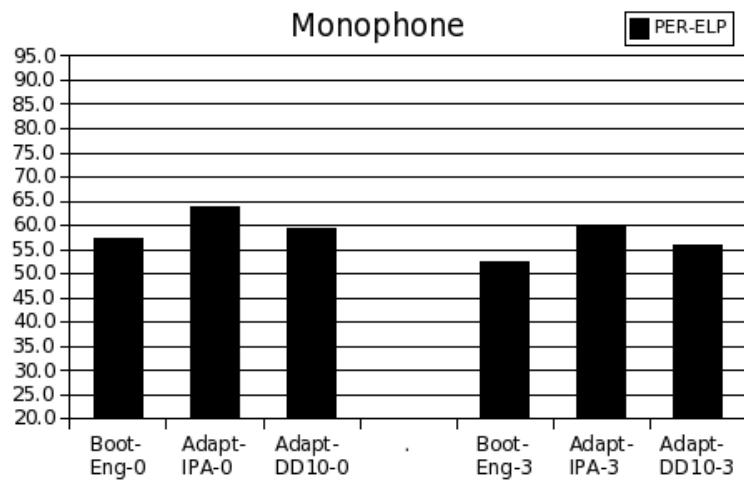


Figure 4.3: Monophone AM adaptation PER-ELP results on the Arabic test data. Recognition performance improves after three iterations of adaptation but does not outperform the baseline bootstrapped AMs.

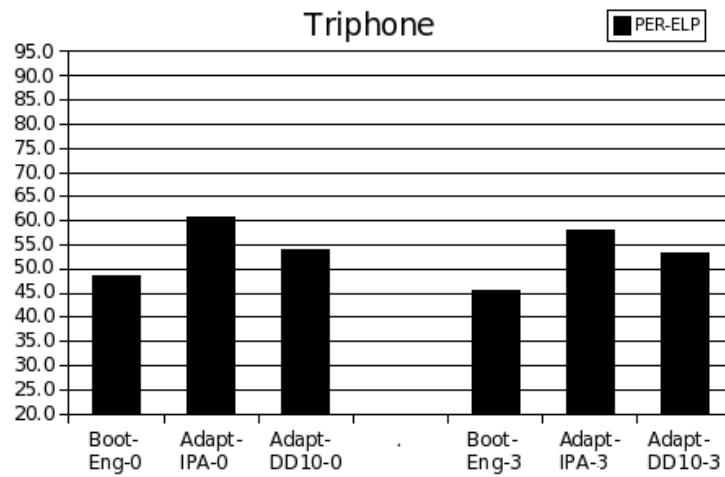


Figure 4.4: Triphone AM adaptation PER-ELP results on the Arabic test data. Recognition performance improves after three iterations of adaptation but does not outperform the baseline bootstrapped AMs.

Arabic models. The first step of the supplementation process is to bootstrap up Arabic, Croatian, German, Japanese, and Turkish AMs using the standard bootstrapping process starting with English AMs. After three iterations of bootstrapping, these final phonetic alignments are saved. The next step depends on whether the “closeness” of the multilingual data is determined by IPA or by the data-driven approach.

The first supplementation experiment uses the IPA labels to group the multilingual data with the Arabic data. For each base phoneme, multilingual acoustic data that matches the Arabic phoneme labels are added to the Arabic data while taking careful consideration not to include non-Arabic phonemes in any triphone context as discussed in Section 4.2.1.

The next group of supplementation experiments uses the data-driven clusters indicated in Table 4.1. The first set of experiments supplemented the Arabic data with the multilingual phonemes listed in the column corresponding to a given KLDM threshold, resulting in supplementing the Arabic data with multilingual data from phonemes either within threshold or with the closest matching multilingual phoneme. The second set of experiments (designated with a *b* after the threshold) requires the supplementation data to be at or below the given threshold. If no multilingual data matches close enough to the Arabic data, that Arabic phoneme model is not supplemented with any extra data. Again, the triphone models are built taking into consideration non-Arabic phoneme contexts as discussed in Section 4.3.3.

Figures 4.5 and 4.6 display the PER-ELPs for all these supplementation experiments. None of the overall PER-ELPs outperform the baseline three iteration bootstrapping from English AMs results. However, there are differences in the individual phoneme performances as seen in the confusion matrices shown in Tables A.13, A.14, A.15, and A.16.

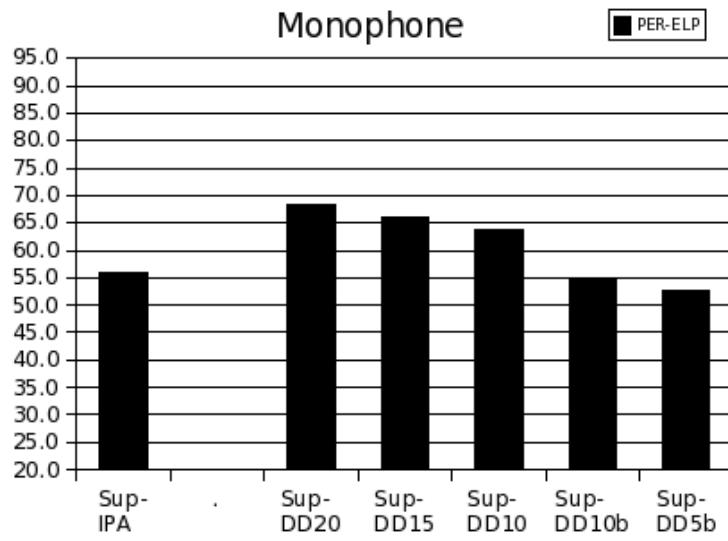


Figure 4.5: Monophone AM supplementation PER-ELP results on the Arabic test data. All AMs stated as bootstrapped from English AMs, with multilingual training data supplementing the Arabic data with different levels of supplementation.

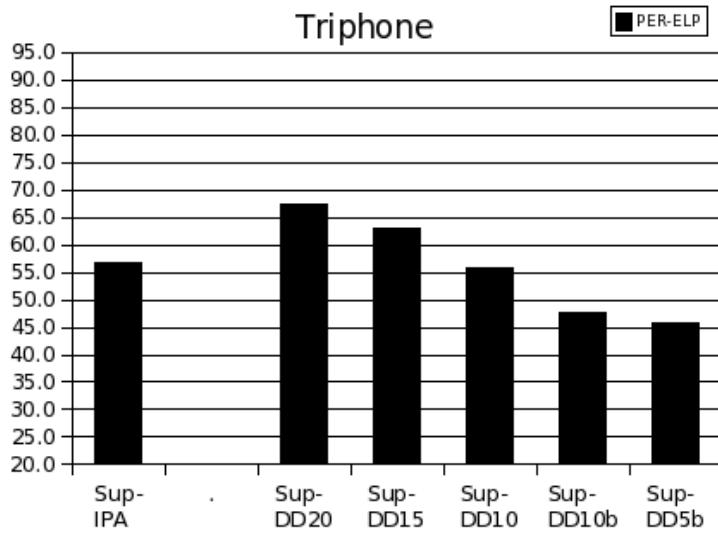


Figure 4.6: Triphone AM supplementation PER-ELP results on the Arabic test data. All AMs stated as bootstrapped from English AMs, with multilingual training data supplementing the Arabic data with different levels of supplementation.

## V. Conclusions

### 5.1 *Review*

In this work multilingual phoneme models are investigated for porting to a new language. Croatian, English, German, Japanese, and Turkish training data are used in various approaches to build an Arabic ASR system. The first approach used standard bootstrapping methods from English AMs. Next IPA-based labels are used to create multilingual AMs that are then bootstrapped to Arabic AMs. Finally, data-driven phoneme clusters based on the KLDM are used to build multilingual AMs that are then bootstrapped to Arabic AMs. The initial bootstrapped data-driven AMs return the lowest WER on the Arabic evaluation set at 73.4% for monophones and 62.7% for triphones. After three iterations of bootstrapping the data-driven phoneme clusters return the lowest PER-ELP for monophones at 51.3% and the IPA-based AMs returned the lowest PER-ELP for triphones at 45.4%.

SMAPLR adaptation is employed to adapt IPA-based multilingual AMs and data-driven multilingual AMs to Arabic data. While multiple iterations of the adaptation improve the AMs, the rate of improvement is not as rapid as for the bootstrapping methods. After three iterations of SMAPLR, no adapted AM outperformed any bootstrapped AM.

Using IPA labels and the KLDM, multilingual training data are used to supplement Arabic training data and to build new Arabic AMs. The phoneme confusion matrices from these experiments show that certain phonemes increase in performance from the supplementation of the multilingual data, but the overall PER-ELP does not show any improvement from the standard method of bootstrapping.

### 5.2 *Future Work*

The results of the phoneme clusters from the KLDM provide for some interesting possibilities. If the KLDM is run on the Arabic AMs, close matching Arabic phonemes may point to an area of high confusability and the need for more focused training data for the particular phonemes. This practice could be applied to other languages

to evaluate the AM space of an ASR and to point to problems with the AMs before testing on an evaluation set is complete. Another possibility is to compute KLDMs for triphones. Also, regarding the KLDM, the addition of phoneme duration information could be added as parameter(s) to the distance measure. Extending the KLDM to multiple mixtures with an approximation is another extension to this work.

These experiments could also be repeated using the HMM Toolkit (HTK) software package instead of SONIC to confirm results and to compare the forward-backward training methodology, as opposed to the align-train-realign Viterbi-based methodology of SONIC.

Also, modifications to the Arabic dictionary and transcripts should be made to remove inconsistencies and the appended word colloquialisms. Changes to the dictionary would have drastic effects on automatic word alignment, which directly affects the way the final AMs are trained. By reducing the error rates in general, different trends using the multilingual modeling approaches may be found.

Other experiments could entail constraining the amount of TL training data even further (1, 5, or 10 hours) and to compare how the bootstrapping and adaptation methods perform.

Combining the results from each AM experiment may result in overall lower error rates. It can be seen that each different AM approach produced different errors across the phoneme set, so combining system results could build upon the benefits of each of the approaches. Another way of combining AM approaches is a merger of the IPA and DD methods: use the IPA mapping for all “known” phonemes in the TL, and use the DD methods for all “unknown” phonemes in the TL.

Finally, changing the number and set of languages used for training and testing could be used to verify the consistency of these results and to measure the robustness of these techniques. What if more languages were added to the language training set? What if the number was kept the same but languages were changed? How does the number of “unknown” phonemes in the TL affect the outcome of the AM approaches?

## Appendix A. Phoneme Confusion Matrices

The following tables display phoneme confusion matrices (as described in Section 3.3.4), for various experiments run on the Arabic evaluation data of the GlobalPhone database. All percentages represent results based on the PER-ELP metric.

For tables that show recognition results in percent, high values along the diagonal and low values in all other cells are signs of a well-performing recognition system. For tables that compare recognition results between two systems with percent differences, positive values along the diagonal and negative values in off-diagonal cells show where the system improved compared to the baseline system.

Table A.1 (page 50) shows the phoneme confusion matrix for the results of decoding the Arabic evaluation data with the initial monophone AM bootstrapped from English AMs with zero iterations. The phoneme labels are arranged to group broad phonetic classes together. Phonemes that are difficult to recognize ( $\leq 20\%$  correct) include /ml/, /Z/, /sl/, /Sl/, and /Cl/. Phonemes that are easier to recognize ( $\geq 70\%$  correct) include /aU/, /f/, /C/, /x/, /G/, /b/, and /k/.

Table A.2 (page 51) compares the PER-ELP-based performance differences using the monophone ML-IPA bootstrapped AM relative to using the monophone English bootstrapped AM. All values are in percent, and positive numbers along the diagonal show where the value increases using the ML-IPA AM versus the English bootstrapped AM, and negative numbers along the diagonal show where the value decreases using the ML-IPA AM versus the English bootstrapped AM. The phonemes /al/, /il/, /u/, /l/, /n/, /f/, and /dd/ increase the most (by over 5% each), while /i/, /r/, and /sl/ decrease the most (by over 5% each).

Table A.3 (page 52) compares performance differences using the monophone ML-DD10 bootstrapped AM relative to using the monophone English bootstrapped AM. This table shows even greater differences than ML-IPA AM. The phonemes /a/, /u/, /l/, /r/, /n/, /Z/, /Q/, and /Hq/ show the greatest improvement, while the phonemes /al/, /rl/, /nl/, and /sl/ show the greatest decline in performance. In general, all *long* phonemes (designated with a *l*) show some degree of degradation

using the ML-DD10 bootstrapped AM, but their corresponding short versions have performance gains.

Tables A.4, A.5, and A.6 show the same order of experiments as just discussed, but after three iterations of bootstrapping the AMs. After three iterations the AMs bootstrapped from English have the same difficulty in recognizing the phonemes: /ml/, /Z/, /sl/, /Sl/, and /Cl/ (see Table A.4, page 53). The phonemes that are easier to recognize are the same as for the initial bootstrapped AMs: /aU/, /f/, /C/, /x/, /G/, /b/, and /k/, with the addition of /m/ and /Hq/.

Table A.5 (page 54) compares the performance difference between three iterations of ML-IPA bootstrapped monophones and three iterations of English bootstrapped monophones. The overall difference in PER-ELP is a 0.5% improvement for the ML-IPA AM, and no phonemes improve by at least 5%, but /i/ and /rl/ decrease by at least 5%.

Table A.6 (page 55) compares the performance difference between three iterations of ML-DD10 bootstrapped monophones and three iterations of English bootstrapped monophones. The PER-ELP improves by 1.3% absolute, and the phoneme confusion matrix shows that /alal/ and /l/ improve and /al/ and /rl/ degrade the most.

Table A.7 (page 56) shows the phoneme confusion matrix for decoding the Arabic evaluation data with the initial triphone AMs bootstrapped from English AMs. The phonemes that are difficult to recognize include /ml/, /Z/, /sl/, /Sl/, and /Cl/. Phonemes that are easier to recognize include /aU/, /f/, /C/, /x/, /G/, /b/, /k/, /a/, and /m/.

Table A.8 (page 57) compares the performance differences of using the triphone ML-IPA bootstrapped AMs and the triphone English bootstrapped AMs. The PER-ELP for the initial ML-IPA AMs is 1.2% worse than the English bootstrapped AMs, but the following phonemes improve performance by at least 5%: /alal/, /al/, /il/,

/ul/, /l/, /rr/, and /dd/. Five phonemes are degraded in performance by at least 5%: /i/, /u/, /x/, /td/, and /Hq/.

Table A.9 (page 58) compares the performance differences of using the triphone English bootstrapped AMs to the triphone ML-DD10 bootstrapped AMs. The PER-ELP for the initial ML-DD10 AMs is 2.3% better than the English bootstrapped AMs. Five phonemes, /u/, /l/, /r/, /rr/ and /Hq/ improve recognition by at least 5% and one phoneme, /sl/, is degraded in performance by more than 5%.

Table A.10 (page 59) shows the phoneme confusion matrix for decoding the Arabic evaluation data after three iterations of bootstrapping starting from the triphone English AMs. After three bootstrapping iterations, the phonemes that are difficult to recognize remain constant; however, /r/ and /Hq/ are added to the list of phonemes that are easier to recognize, while /aU/ downgrades to 67.9% correct.

Table A.11 (page 60) compares the performance differences of using the triphone ML-IPA bootstrapped AMs after three iterations and the triphone English bootstrapped AMs after three iterations. The PER-ELP between the two experiments is identical at 45.4%, but /ul/ is recognized 7.7% better with the ML-IPA AM and /ll/, /al/, and /u/ are recognized 7.6%, 5.5%, and 5.2% worse with the ML-IPA AM respectively. All other phonemes have a change of less than  $\pm 5\%$ .

Table A.12 (page 61) compares the performance differences of using the triphone ML-DD10 bootstrapped AMs after three iterations and the triphone English bootstrapped AMs after three iterations. The PER-ELP with the ML-DD10 AM increases to 46.7%, due in part to the fact that /a/, /ll/ and /dd/ are recognized 11.7%, 6.6%, and 8.9% worse than the English bootstrapped AM, respectively. However, the phoneme /ul/ is recognized 8.9% better than the English bootstrapped AM.

Table A.13 (page 62) compares the performance difference between three iterations of English bootstrapped monophones to this AM supplemented with monophone data based on IPA labels . The overall difference is 3.4% worse in PER-ELP for the ML-IPA supplemented AM, but from the phoneme confusion matrix it can be seen

that the individual phoneme performances varies drastically. Six phonemes are improved in correct recognition by 5.0% to 17.7%, while 13 phonemes degrade in correct recognition by -6.5% to -43.2%.

Table A.14 (page 63) shows the results of reducing the amount of supplementation data. With the threshold set to  $5b$ , only the following Arabic phonemes are supplemented: {/i/, /u/, /m/, /n/, /f/, /s/, /b/, /t/, /d/}, and these phonemes are only supplemented by multilingual data that match with a KLDM of “five” or less. Out of this set of phonemes, all increase in recognition performance except /m/ (-1.0%) and /s/ (-0.9%). Note that, /nl/ is recognized 9.1% better with these AMs.

Table A.15 (page 64) compares the performance difference between three iterations of English bootstrapped triphones to this AM supplemented with triphone data based on IPA labels. The PER-ELP is degraded by 11.4%, and 21 phonemes decrease in recognition performance by -6.4% to -38.5%. The phonemes /ml/, /nl/, /Z/, /sl/, and /Sl/ increase in recognition performance by 7.1% to 10.5%.

Table A.16 (page 65) shows the results of supplementing the English bootstrapped triphone AM with multilingual data that match within a KLDM of “five” or less. The phoneme set is the same as mentioned before, {/i/, /u/, /m/, /n/, /f/, /s/, /b/, /t/, /d/}, and each of these phonemes increase in recognition performance except /m/ (-0.2%) and /d/ (-1.9%).

































## Appendix B. Additional Information

Phoneme	Example	Phoneme	Example
AA	father	DD	had <u>_</u>
AE	mad	KD	talk <u>_</u>
AH	but	JH	Jerry
AO	for	K	kitten
AW	frown	L	listen
AX	alone	M	manager
AXR	butter <u>_</u>	N	nancy
AY	hire	NG	fishing <u>_</u>
B	bob	OW	cone
CH	church	OY	boy
D	don't	P	pop
PD	top <u>_</u>	R	red
TD	lot <u>_</u>	S	sonic
DX	butter	TS	bits
DH	them	GD	mug <u>_</u>
EH	bed	SH	show
ER	bird	T	tot
EY	state	TH	thread
F	friend	UH	hood
G	grown	UW	moon
HH	had	V	very
IH	bitter	W	weather
IX	roses <u>_</u>	Y	yellow
IY	beat	Z	bees <u>_</u>
BD	tab	ZH	measure

Table B.1: American English phoneme set used by SONIC [28].

Language	Words in Dictionary
Arabic	47,688
Croatian	24,186
German	40,706
Japanese	32,543
Turkish	31,944

Table B.2: Number of words in each language's dictionary.

Phoneme	Croatian	German	Japanese	Turkish	Arabic
a	(9749) 70.0	(4009) 78.7			(21845) 47.0
al		(1645) 112.9			(10010) 38.3
alal					(60) 106.8
ab			(12548) 68.6	(8860) 74.7	
abl			(185) 125.4		
ae		(271) 106.0			
aI		(1275) 113.2			(420) 86.2
atu		(1444) 60.0			
aU		(482) 135.5			(444) 97.5
b	(1283) 78.4	(1776) 68.4	(953) 68.9	(1858) 74.0	(2376) 71.9
cp	(633) 121.0				
C		(1101) 96.4			(531) 111.1
Cl					(143) 104.1
d	(2965) 63.0	(4203) 50.8	(2088) 57.2	(3316) 64.6	(2168) 77.6
dd					(502) 64.9
dp	(191) 91.6				
dZ	(11) 57.3		(1318) 90.3	(970) 92.8	
D					(417) 63.4
e	(7670) 56.2	(2464) 64.1	(4864) 70.5	(8596) 72.2	
el		(2405) 69.3	(1694) 124.8		
etu		(7252) 45.2			
eU		(335) 115.7			
f	(189) 100.0	(2545) 107.5	(319) 76.4	(346) 104.7	(1662) 90.8
g	(1451) 63.8	(1975) 65.5	(2218) 55.3	(1013) 84.5	
G					(994) 90.4
h		(819) 69.4	(1808) 69.8	(762) 64.5	(1417) 67.0
H					(1103) 108.8
Hq					(2274) 66.9
i	(8047) 58.2	(4340) 51.4	(9975) 54.5	(7312) 54.6	(9076) 44.5
il		(2408) 73.2	(160) 129.3		(3774) 77.7
i2				(3939) 45.0	
j	(3293) 50.3	(347) 66.5	(2149) 58.7	(2538) 69.2	(2069) 62.5
k	(3353) 85.1	(1826) 96.2	(7625) 76.4	(3559) 91.6	(1017) 98.6
l	(2190) 55.1	(3118) 56.1	(3501) 48.5	(5548) 49.7	(7955) 44.4
ll					(399) 66.0
L	(421) 57.8				
m	(2450) 80.1	(2252) 80.7	(2311) 81.1	(2719) 72.5	(4262) 64.6
ml					(180) 76.5

Table B.3: Count and average duration of each multilingual phoneme in the GlobalPhone test subset. The number of times the phoneme occurs is in “()”, followed by the average duration (in ms) as determined by the Viterbi-based alignment of the reference transcripts by using the third iteration of English bootstrapped AMs for that language.

Phoneme	Croatian	German	Japanese	Turkish	Arabic
n	(4906) 60.0	(9179) 66.5	(5179) 59.8	(6273) 62.4	(3805) 65.8
nl					(461) 72.9
ng		(833) 88.4			
nj	(617) 82.2				
nq			(3981) 81.3		
o	(7985) 61.5	(1396) 72.3	(7732) 69.8	(1826) 96.2	
ol		(1187) 92.1	(3043) 126.5		
oe		(126) 77.1		(854) 89.7	
oel		(125) 91.1			
p	(2752) 94.3	(1041) 94.7	(264) 73.8	(745) 102.3	
q					(1489) 105.6
Q			(941) 93.6		(4002) 62.1
r	(5011) 46.8	(6131) 48.8		(5430) 47.1	(3043) 51.1
rl					(237) 62.1
rr					(215) 62.6
s	(4334) 101.0	(2990) 107.9	(3057) 96.9	(2752) 119.8	(1990) 103.6
sl					(294) 101.9
sj	(565) 114.4				
S		(1521) 109.2	(2613) 104.2	(1112) 117.5	(489) 99.1
Sl					(129) 100.5
sft				(740) 53.7	
t	(3898) 71.7	(5991) 70.9	(4454) 69.2	(3228) 88.7	(4557) 78.2
td					(848) 89.5
ts	(1334) 103.7	(1762) 119.6	(987) 94.5		
tS	(811) 109.2		(849) 104.4	(679) 108.8	
T					(511) 89.0
u	(3865) 72.5	(2157) 65.7		(2496) 55.2	(3201) 52.3
ul		(752) 76.6			(1128) 81.1
ue		(348) 53.9		(1411) 53.0	
uel		(375) 76.5			
v	(3218) 50.8	(1628) 66.0		(838) 62.9	
w			(878) 83.1		(2392) 59.0
W			(6301) 46.6		
Wl			(1099) 101.1		
x	(969) 76.8	(695) 93.8			(543) 107.7
z	(1624) 89.9	(1530) 92.8	(635) 74.6	(1052) 91.0	(481) 94.8
zj	(390) 94.1				
Z				(61) 119.8	(187) 55.2

Table B.4: Continuation of the count and average duration of each multilingual phoneme in the GlobalPhone test subset. The number of times the phoneme occurs is in “()”, followed by the average duration (in ms) as determined by the Viterbi-based alignment of the reference transcripts by using the third iteration of English bootstrapped AMs for that language.

## Bibliography

1. J. Billa, K. Ma, J. McDonough, G. Zavaliagkos, D. Miller, K. Ross, and A. El-Jaroudi. Multilingual speech recognition: The 1996 Byblos callhome system. In *Proceedings of Eurospeech 1997*, 1997.
2. A. Bonafonte, X. Ros, and J. Marino. An efficient algorithm to find the best state sequence in HSMM. In *Proceedings of Eurospeech 1993*, 1993.
3. P. Bonaventura, F. Gallochchio, and G. Micca. Multilingual speech recognition for flexible vocabularies. In *Proceedings of Eurospeech 1997*, 1997.
4. U. Bub, J. Koehler, and B. Imperl. In-service adaptation of multilingual hidden markov models. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, 1997.
5. B. Byrne and et. all. Towards language independent acoustic modeling. In *The 1999 Johns Hopkins University Language Engineering Workshop*, 1999.
6. K. Çarkı, P. Geutner, and T. Schultz. Turkish LVCSR: Towards better speech recognition for agglutinative languages. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, June 2000.
7. P. Clarkson and R. Rosenfeld. Statistical language modeling using the CMU-Cambridge toolkit. In *Proceedings of Eurospeech 1997*, September 1997.
8. C. Corredor-Ardoy, L. Lamel, M. Adda-Decker, and J. Gauvain. Multilingual phone recognition of spontaneous telephone speech. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, May 1998.
9. R. Duda, P. Hart, and D. Stork. *Pattern Classification*. John Wiley and Sons, Inc., New York, NY, 2001.
10. J. Fiscus. Sclite overview. Online Tutorial, 1996. Available at <http://www.icsi.berkeley.edu/Speech/docs/sctk-1.2/sclite.htm>.
11. G. Forney. The Viterbi algorithm. *Proceedings of the IEEE*, 61(3), March 1973.
12. M. Gales and P. Woodland. Mean and variance adaptation within the MLLR framework. *Computer Speech and Language*, 10:249–264, 1996.
13. J. Garofalo, D. Graff, D. Paul, and D. Pallet. CSR-1 (WSJ0) complete. Linguistic Data Consortium Website, 1993. Available at <http://www.ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC93S6A>.
14. S. Gokcen and J. Gokcen. A multilingual phoneme and model set: Toward a universal base for automatic speech recognition. In *IEEE Workshop on Automatic Speech Recognition and Understanding*, December 1997.

15. J. Golberger and H. Aronowitz. A distance measure between GMMs based on the unscented transform and its application to speaker recognition. In *Proceedings of Interspeech 2005*, September 2005.
16. B. Gold and N. Morgan. *Speech and Audio Signal Processing*. John Wiley and Sons, Inc., New York, NY, 2000.
17. IPA. The international phonetic association - IPA chart. *Journal of the International Phonetic Association* 23, 1993.
18. B. Juang, L. Rabiner, S. Levinson, and M. Sondhi. Recent developments in the application of Hidden Markov Models to speaker independent isolated word recognition. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, 1985.
19. D. Jurafsky and J. Martin. *Speech and Language Processing*. Prentice Hall, New York, NY, 2000.
20. D. J. Kershaw. *Phonetic Context-Dependency In a Hybrid ANN/HMM Speech Recognition System*. Ph.D. dissertation, St. John's College, University of Cambridge, January 1997.
21. J. Koehler. Multi-lingual phoneme recognition exploiting acoustic-phonetic similarities of sounds. In *Proceedings of the 1996 International Conference on Spoken Language Processing*, October 1996.
22. J. Koehler. Language adaptation of multilingual phone models for vocabulary independent speech recognition tasks. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, May 1998.
23. J. Koehler. Multilingual phone models for vocabulary-independent speech recognition tasks. *Speech Communication*, 35:21–30, 2001.
24. C. Legetter and P. Woodland. Maximum likelihood linear regression for speaker adaptation of continuous density Hidden Markov Models. *Computer Speech and Language*, 9, 1995.
25. S. Levinson, L. Rabiner, and M. Sondhi. An introduction to the application of the theory of probabilistic functions of a Markov process to automatic speech recognition. *The Bell System Technical Journal*, 62(4), April 1983.
26. C. Nieuwoudt and E. Botha. Cross-language use of acoustic information for automatic speech recognition. *Speech Communication*, 38:101–113, 2002.
27. D. O'Shaughnessy. *Speech Communications*. IEEE Press, New York, NY, 2000.
28. B. Pellom and K. Hacıoğlu. SONIC: The University of Colorado continuous speech recognizer. Technical Report TR-CSLR-2001-01, University of Colorado, Boulder, May 2005.

29. D. Pye and P. Woodland. Experiments in speaker normalization and adaptation for large vocabulary speech recognition. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, April 1997.
30. J. Pylkkonen and M. Kurimo. Duration modeling techniques for continuous speech recognition. In *Proceedings of the 2004 International Conference on Spoken Language Processing*, May 2004.
31. L. Rabiner. A tutorial on Hidden Markov Models and selected applications in speech recognition. *IEEE Transactions in Acoustic Speech and Signal Processing*, 77(2), February 1989.
32. L. Rabiner and B. Juang. *Fundamentals of Speech Recognition*. Prentice Hall, New York, NY, 1993.
33. M. Russell and A. Cook. Experimental evaluation of duration modeling techniques for automatic speech recognition. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, 1987.
34. T. Schultz. Globalphone: A multilingual speech and text database developed at Karlsruhe University. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, September 2002.
35. T. Schultz and A. Waibel. Fast bootstrapping of LVCSR systems with multilingual phoneme sets. In *Proceedings of Eurospeech 1997*, September 1997.
36. T. Schultz and A. Waibel. Multilingual and crosslingual speech recognition. In *Proceedings of the DARPA Broadcast News Transcription and Understanding Workshop*, February 1998.
37. T. Schultz and A. Waibel. Polyphone decision tree specialization for language adaptation. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, June 2000.
38. T. Schultz and A. Waibel. Language independent and language adaptive large vocabulary speech recognition. In *Speech Communication*, February 2001.
39. J. Silva and S. Narayanan. An upper bound for the Kullback-Leibler divergence for left-to-right transient Hidden Markov Models. In *IEEE Transactions on Information Theory*, 2005.
40. O. Siohan, C. Chesta, and C.-H. Lee. Joint maximum a posteriori adaptation of transformation and HMM parameters. *IEEE Transactions on Speech and Audio Processing*, 9(4):417–428, 2001.
41. O. Siohan, T. Myrvoll, and C.-H. Lee. Structural maximum a posteriori linear regression for fast HMM adaptation. *Computer Speech and Language*, 16, 2002.
42. A. Žgank, B. Horvat, and Z. Kačić. Data-driven generation of phonetic broad classes, based on phoneme confusion matrix similarity. *Speech Communication*, 47:379–393, 2005.

43. A. Waibel, P. Geutner, L. Mayfield-Tomokiyo, T. Schultz, and M. Woszczyna. Multilinguality in speech and spoken language systems. In *Proceedings of the IEEE*, August 2000.
44. B. Walker, B. Lackey, J. Muller, and P. Schone. Language-reconfigurable universal phone recognition. In *Proceedings of Eurospeech 2003*, September 2003.
45. F. Weng, L. H. Bratt, Neumeyer, and A. Stolcke. A study of multilingual speech recognition. In *Proceedings of Eurospeech 1997*, 1997.
46. B. Wheatley, K. Kondo, W. Anderson, and Y. Muthusamy. An evaluation of cross-language adaptation for rapid HMM development in a new language. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, April 1994.
47. E. Wong and S. Sridharan. Three approaches to multilingual phone recognition. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, April 2003.

REPORT DOCUMENTATION PAGE				Form Approved OMB No. 0704-0188	
<p>The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing the burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.</p> <p><b>PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.</b></p>					
1. REPORT DATE (DD-MM-YYYY) 14-09-2006	2. REPORT TYPE Master's Thesis	3. DATES COVERED (From - To) June 2005 - September 2006			
4. TITLE AND SUBTITLE  Multilingual Phoneme Models for Rapid Speech Processing System Development			5a. CONTRACT NUMBER		
			5b. GRANT NUMBER		
			5c. PROGRAM ELEMENT NUMBER		
6. AUTHOR(S)  Hansen, Eric, G.			5d. PROJECT NUMBER		
			5e. TASK NUMBER		
			5f. WORK UNIT NUMBER		
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Air Force Institute of Technology Graduate School of Engineering and Management (AFIT/EN) 2950 Hobson Way Wright-Patterson Air Force Base, OH 45433			8. PERFORMING ORGANIZATION REPORT NUMBER AFIT/GE/ENG/06-62		
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) AFRL/HECP 2255 H St. Bldg 248 Wright-Patterson Air Force Base, OH 45433 C/O Dr. Timothy Anderson Tim.Anderson@wpafb.af.mil (937)255-8817			10. SPONSOR/MONITOR'S ACRONYM(S) AFRL/HECP		
			11. SPONSOR/MONITOR'S REPORT NUMBER(S)		
12. DISTRIBUTION/AVAILABILITY STATEMENT  Approved for Public Release; Distribution Unlimited.					
13. SUPPLEMENTARY NOTES					
<b>14. ABSTRACT</b> Current speech recognition systems tend to be developed only for commercially viable languages. The resources needed for a typical speech recognition system include hundreds of hours of transcribed speech for acoustic models and 10 to 100 million words of text for language models; both of these requirements can be costly in time and money. The goal of this research is to facilitate rapid development of speech systems to new languages by using multilingual phoneme models to alleviate requirements for large amounts of transcribed speech. The GlobalPhone database, which contains transcribed speech from 15 languages, is used as source data to derive multilingual phoneme models. Various bootstrapping processes are used to develop an Arabic speech recognition system starting from monolingual English models, International Phonetic Association (IPA) based multilingual models, and data-driven multilingual models. The Kullback-Leibler distortion measure is used to derive data-driven phoneme clusters. It was found that multilingual bootstrapping methods outperform monolingual English bootstrapping methods on the Arabic evaluation data initially, and after three iterations of bootstrapping all systems show similar performance levels.					
<b>15. SUBJECT TERMS</b> Speech recognition, Multilingual, Cross-language, Hidden Markov Models.					
16. SECURITY CLASSIFICATION OF: a. REPORT U		17. LIMITATION OF ABSTRACT U	18. NUMBER OF PAGES 84	19a. NAME OF RESPONSIBLE PERSON Dr. Steven C. Gustafson	
				19b. TELEPHONE NUMBER (Include area code) (937)255-3636 x4598 Steven.Gustafson@afit.edu	