

**An Efficient First Pass of a Two-Stage Approach for Automatic  
Language Identification of Telephone Speech**

By

Jonathan K. Law

Submitted in partial fulfillment  
of the requirements for the degree of  
Doctor of Professional Studies  
in Computing

at

School of Computer Science and Information Systems

Pace University

May 2002

UMI Number: 3118356

PREVIEW

UMI<sup>®</sup>

---

UMI Microform 3118356

Copyright 2003 by ProQuest Information and Learning Company.

All rights reserved. This microform edition is protected against  
unauthorized copying under Title 17, United States Code.

---

ProQuest Information and Learning Company  
300 North Zeeb Road  
PO Box 1346  
Ann Arbor, MI 48106-1346

## Dissertation Signature (Approval) Page

We hereby certify that this dissertation, submitted by Jonathan K. Law, satisfies the dissertation requirements for the degree of Doctor of Professional Studies and has been approved.

---

Dr. Charles Tappert  
Chairperson of Dissertation Committee

20 May 2002

---

Dr. Fred Grossman  
Dissertation Committee Member

20 May 2002

---

Dr. Zhong-hua Wang  
Dissertation Committee Member

20 May 2002

School of Computer Science and Information Systems  
Pace University 2002

## **Abstract**

### **An Efficient First Pass of a Two-Stage Approach for Automatic Language Identification of Telephone Speech**

By  
Jonathan K. Law

Submitted in partial fulfillment  
of the requirements for the degree of  
Doctor of Professional Studies  
in Computing

June 2002

Automatic language identification, recognizing a speaker's language from a speech signal, is gaining increased importance in the context of economic globalization. The most accurate systems for such identification use multiple, large vocabulary continuous speech recognizers. But their scalability is limited, as each new language added requires complete recognizers and an enormous amount of training. For increased efficiency, we propose a Two-Stage approach: an efficient clustering algorithm (least cost) first selects the top candidates at an accuracy level of over 80%. Then, speech recognizers (high cost) would be used to narrow the field from approximately three candidates to correctly identify the language.

In this research, we describe the first phase of this Two-Stage process in which we are able to narrow the list of possible languages through an extremely scalable and flexible system. We cover how test and reference patterns (acoustic feature vectors) are extracted from speech utterances, how cepstral coefficients are used, and how reference models are generated from the reference patterns using Vector Quantization (VQ) clustering algorithm. Various distance measures are also examined in the selection phase to find the best method. We show using the top-N strategy in this first stage leads to substantial improvements over existing systems in discriminating between different languages, and our experiments showed the top three choices yield 87.2% in a 5-language task, 88.6% in a 7-language task, and the top five yield 87.5% in a 10-language task. The second iteration in the 10-language task easily narrows it down to three choices with 80% probability. Further research could improve the entire process, but combining this particular methodology with current best practices has proven an extremely efficient and cost effective way to address the challenges of automatic language identification.

## Acknowledgments

It really did take a village to help me get through this program. There are many people to thank for their support and encouragement, without whose help this thesis would not have been possible.

First, I would like to thank my advisors. Dr. Tappert has diligently invited scholars and industry experts in the latest technological fields, which has further sparked my interest in speech recognition. Dr. Wang has tirelessly shown me his unquenchable curiosity and love for the subject of speech recognition technology. I must also thank Dr. Grossman for the reality check he provided, as his constant and justifiable rejections of my initially unachievable topics helped shape my work into something that could be accomplished. I must also thank the constant support and encouragement of the Computer Science department faculty, which kept me going over the last three or so years.

Next are all the friends I've made, a list is too long to mention, but whose friendship, Indian village jokes, and "Just-in-Time" motto seemed to work every time on our last minute assignments. My dear friend Patrick Wong, who introduced me to this program and also my "room mate" during these college years. My friends made those long, crazy monthly trips across the country not only enjoyable, but also memorable.

Finally, I thank my family, which has been wonderfully supportive and thoroughly believed my "midlife crisis" excuse for going back to school 20 years after my last college exam. I first thought about applying for my doctorate degree when my sister Joyce received her A.I.A. certification. Architecture was my own major in college, but with her degree she has done everything I always wanted to do. I thank her for our sibling rivalry and taking full advantage of the opportunities that our parent never had; this degree will hopefully once again put me one step ahead of my sister. My brother-in-law Oliver, who is not only a seasoned proofreader, but also works well under imposed deadlines. All errors in this paper are mine.

And the most important person behind this accomplishment is my loving wife Joan, who has put up with me all of these years and has taken great care in the raising of our five lovely college-bound children. I am indebted to her for her support and love.

## Table of Contents

<i>Abstract</i> .....	iii
<i>Acknowledgments</i> .....	iv
<i>Table of Contents</i> .....	iv
<i>List of Tables</i> .....	vii
<i>Figures List</i> .....	ix
<i>Chapter 1.</i> .....	1
<i>Introduction</i> .....	1
1.1 <i>Overview</i> .....	2
1.1.1 <i>Differences Among Languages</i> .....	3
1.1.2 <i>The Challenge</i> .....	4
1.2 <i>The Driving Force</i> .....	5
1.3 <i>Related Researches</i> .....	6
1.3.1 <i>The first twenty years</i> .....	6
1.3.2 <i>The last twenty years</i> .....	7
1.4 <i>Overview of Proposed Two-Stages Approach</i> .....	8
1.4.1 <i>First Pass - Efficient Algorithm in Selecting Top Candidates</i> .....	8
1.4.2 <i>Second Stage – Accurately Identify unknown signal</i> .....	11
1.5 <i>Outline of the Dissertation</i> .....	12
<i>Chapter 2.</i> .....	14
<i>Related Studies and Experiments</i> .....	14
2.1 <i>Overview of Speech Recognition</i> .....	15
2.1.1 <i>The Data</i> .....	15
2.1.2 <i>Preprocessing</i> .....	19
2.1.3 <i>Recognition Application</i> .....	21
2.2 <i>Related Research and Applications Approach</i> .....	25
2.2.1 <i>Gaussian Mixture Model Approach (GMM)</i> .....	26
2.2.2 <i>Phone-based Acoustic Likelihood Approach</i> .....	27
2.2.3 <i>Phone Recognition followed by Language Modeling (PRLM)</i> .....	28
2.2.4 <i>Large Vocabulary Approach</i> .....	29
2.3 <i>Pros and Cons of Related Researches</i> .....	31
2.4 <i>Why Vector Quantization for First Stage Classification?</i> .....	32
2.5 <i>Possible Language-dependent Phone Recognizer for Second Stage?</i> .....	33
<i>Chapter 3.</i> .....	35
<i>Speech Corpus and Language Task</i> .....	35

3.1	<i>Corpora from Oregon Graduate Institute (OGI)</i>	35
3.1.1	<i>Multi-Language Corpus</i>	36
3.1.2	<i>22 Language Corpus</i>	36
3.1.3	<i>Protocol</i>	36
3.2	<i>The Training Language Tasks</i>	37
3.2.1	<i>Five-Language Task</i>	38
3.2.2	<i>Seven-language task</i>	39
3.2.3	<i>Ten-Language Task</i>	40
3.3	<i>The Testing Language Tasks</i>	40
<i>Chapter 4.</i>		42
<i>System Design and Implementation</i>		42
4.1	<i>Pre-processing Module</i>	43
4.1.1	<i>Spectral Features Extraction</i>	44
4.1.2	<i>Mel Frequency Cepstral Coefficients (MFCC)</i>	45
4.2	<i>Feature Extraction Process</i>	48
4.3	<i>Statistical Pattern-Recognition module</i>	49
4.4	<i>Pros and Cons of the VQ Baseline System</i>	50
4.5	<i>Language Identification Module</i>	53
4.5.1	<i>General Framework</i>	54
4.5.2	<i>Acoustic Pre-processing</i>	56
4.5.3	<i>Language Training component</i>	56
4.5.4	<i>Language Testing Component</i>	57
4.6	<i>Baseline Benchmark Evaluation</i>	59
<i>Chapter 5.</i>		63
<i>Improvements to Baseline System</i>		63
5.1	<i>Hamming Window Sizes</i>	63
5.2	<i>Codebook Sizes on Language Modeling</i>	67
5.2.1	<i>Combining Codebook and Hamming Windows Enhancements</i>	70
5.3	<i>Delta-Delta MFCC Feature Parameterizations</i>	71
5.4	<i>Proposed experiments to Language Identification Modeling</i>	73
<i>Chapter 6.</i>		75
<i>Distance Measures Classification</i>		75
6.1	<i>Background</i>	75
6.2	<i>Distance Measures</i>	76
6.2.1	<i>Euclidean Distance (ED)</i>	78
6.2.2	<i>City Block (Manhattan) Distance (CBD)</i>	78
6.2.3	<i>Weighted Euclidean Distance (WED)</i>	80
6.2.4	<i>Segmentation of Static and Dynamic Features Distance (SSDD)</i>	83
6.2.5	<i>Results comparison on Distance Measures</i>	87
6.3	<i>Data Fusion on Distance Measures</i>	88
6.3.1	<i>City Block Distance on Segmentation of Static/Dynamic Features Distance (CBD-SSDD)</i>	89
6.3.2	<i>Weighted Euclidean Distance on Segmentation of Static/Dynamic Features Distance (WED-SSDD)</i>	91
6.3.3	<i>Alpha coefficient to Euclidean Distance (AC-ED)</i>	92

6.3.4	<i>Results comparison on Data Fusion in Distance Measures.....</i>	96
6.4	<i>Confusion Matrix.....</i>	97
6.4.1	<i>Confusion Matrix experiments on Closed-Set test .....</i>	98
6.4.2	<i>Confusion Matrix experiments on Open-Set test.....</i>	101
6.5	<i>Applying Top-N Strategy.....</i>	104
6.6	<i>Pair-wise results using English language .....</i>	108
6.7	<i>Summary and Discussion .....</i>	110
Chapter 7.	.....	112
Conclusion and Future Work.....		112
7.1	<i>Language Identification General Issues.....</i>	112
7.2	<i>Future Work Discussions and Challenges.....</i>	113
Appendix .....		116
A. 1	<i>Glossary of Terms and Acronyms.....</i>	116
References .....		120



## List of Tables

Table 3-1	Summary of the Five-Language training set.....	39
Table 3-2	Summary of the Seven-Language training set.....	39
Table 3-3	Summary of the Ten-Language training set.....	40
Table 4-1	Baseline result for 5 languages using Hamming window 256.....	60
Table 4-2	Baseline result for 7 languages using Hamming window 256.....	61
Table 4-3	Baseline result for 10-languages (Hamming window 256) .....	62
Table 5-1	Accuracy results for 5 languages test set (Hamming window 512).....	64
Table 5-2	Comparison of 5 languages test set with different Hamming window sizes	64
Table 5-3	Accuracy results for 7 languages test set(Hamming window 512).....	65
Table 5-4	Accuracy results for 10 languages test(Hamming window 512) .....	65
Table 5-5	Comparisons of Hamming Window Sizes.....	66
Table 5-6	Codebook Size Comparisons with Hamming Window Size of 512 .....	68
Table 5-7	Codebook Size Comparison with Hamming window Size of 256 .....	69
Table 5-8	Combine scoring of three tests data task.....	70
Table 5-9	Accuracy results for 5 languages test set with delta-delta expressions .....	72
Table 5-10	Accuracy results for 7 languages test set with delta-delta expressions ....	73
Table 6-1	Result on 5 languages using City Block Distance Measure (CBD).....	79
Table 6-2	Result on 7 languages using City Block Distance Measure (CBD).....	79
Table 6-3	Result on 10 languages using City Block Distance Measure (CBD).....	80
Table 6-4	Results on 5 languages using Weighted Euclidean Distance (WED).....	82
Table 6-5	Results on 7 languages using Weighted Euclidean Distance (WED).....	82
Table 6-6	Results on 10 languages using Weighted Euclidean Distance (WED).....	83
Table 6-7	Results on 5 languages using Segmentation of Static and Dynamic Euclidean Distance Measure (SSDD).....	86
Table 6-8	Results on 7 languages using Segmentation of Static and Dynamic Euclidean Distance Measure (SSDD).....	86
Table 6-9	Results on 10 languages using Segmentation of Static and Dynamic Euclidean Distance Measure (SSDD).....	87
Table 6-10	Experiment results comparison on various Distance Measures.....	87
Table 6-11	Results on 5 languages using CBD-SSDD Distance measures .....	89
Table 6-12	Results on 7 languages using CBD-SSDD Distance measures .....	90
Table 6-13	Results on 10 languages using CBD-SSDD Distance measures .....	90
Table 6-14	Results on 5 languages using WED-SSDD Distance measures.....	91
Table 6-15	Results on 7 languages using WED-SSDD Distance measures.....	91
Table 6-16	Results on 10 languages using WED-SSDD Distance measures.....	92
Table 6-17	5-languages (ED) using alpha coefficient from 0.1 to 0.9 .....	94
Table 6-18	Results comparison of with the average coefficient from 0.10 to 0.49 ....	95
Table 6-19	Results comparing distance measure models.....	96
Table 6-20	Confusion matrixes for 10-second utterances: 5language-task .....	98
Table 6-21	Confusion matrixes for 45-second utterances: 5language-task .....	99
Table 6-22	Confusion matrixes for 10-second utterances: 7language-task .....	99

Table 6-23	Confusion matrixes for 45-second utterances: 7language-task .....	100
Table 6-24	Confusion matrixes for 45-second utterances: 10-language-task .....	101
Table 6-25	Confusion matrix for 10-second utterances: 5-language-task Open-Set	102
Table 6-26	Confusion matrix for 45-second utterances: 5-language-task Open-Set	102
Table 6-27	Confusion matrix for 10-second utterances: 7-language-task Open-Set	103
Table 6-28	Confusion matrix for 45-second utterances: 7-language-task Open-Set	104
Table 6-29	Listings of Top-N strategy in languages for 5-language task .....	105
Table 6-30	Listings of Top-N strategy in combined accuracy for 5-language task ..	106
Table 6-31	Listings of Top-N strategy in combined accuracy for 7-language task ..	107
Table 6-32	Listings of Top-N strategy in combined accuracy for 10-language task	107
Table 6-33	Confusion Matrix for English test set after second iteration .....	108
Table 6-34	Results on the English-other language-pair task.....	109
Table 6-35	Results on the French-Cantonese language-pair task .....	110

PREVIEW

## Figures List

Figure 1-1	Conceptual diagram of Vector Quantization Distortion .....	10
Figure 1-2	Conceptual diagram of nearest neighbor selection .....	11
Figure 1-3	Overview of a Language-Dependent Phone Recognition system.....	12
Figure 2-1	Block diagram of a task-oriented speech-recognition system .....	15
Figure 2-2	International Phonetic Alphabet (Rev. 1993, updated 1996).....	18
Figure 2-3	Sampled wave file of spoken word “Hello” .....	23
Figure 2-4	Parallel Phone recognition and Language Modeling (PRLM)[5].....	29
Figure 4-1	Components of language Identification application .....	43
Figure 4-2	The procedure of MFCC feature extraction[20] .....	45
Figure 4-3	One frame of sampled signal for spectral analysis .....	46
Figure 4-4	Pre-emphasis processing in time domain.....	47
Figure 4-5	Frame signal converted to frequency Domain and filter bank[20] .....	48
Figure 4-6	Block diagram of pattern-recognition approach .....	49
Figure 4-7	VQ codebook Generation Algorithms .....	52
Figure 4-8	General System Frameworks .....	55
Figure 5-1	VQ Codebook Size Comparisons for Hamming Window Size of 512.....	68
Figure 5-2	Codebook Size Comparisons for Hamming Window Size of 256 .....	69
Figure 5-3	Hamming window sizing accuracy comparison .....	70
Figure 6-1	Distance Measure logic flow diagrams.....	77
Figure 6-2	Separate Static and Dynamic feature scoring process flow diagram.....	85
Figure 6-3	Bar charts showing results on Distance Measures.....	88
Figure 6-4	5-languages Distance Measures comparison .....	97

## **Chapter 1.**

### **Introduction**

Due to the rapid growth in processing power of the computer chip in recent years, hardware and software manufacturers have begun incorporating speech recognition into their products. This trend is creating awareness of speech recognition as a viable tool. Corporate America, which is constantly looking to increase its bottom line by finding ways to increase productivity, is trying to leverage speech recognition. Most major corporations in the United States have business relations with other countries around the world, and the expansion of speech recognition in understanding human speech demands that communication cross the boundaries of languages.

An automatic language identification system can play an important role in a speech recognition system in which it acts as the front-end processor, such as routing foreign telephone calls to appropriate operators with language skill. It could even save lives when split second decisions alerting the right medical personnel during an emergency are the difference between life and death. In addition, all foreign government defense departments have a keen interest in using this capability for monitoring purposes. But duration of unknown utterances and high number of languages pose challenging tasks for researchers in multilinguality and language identification. This research tries to solve the scalability and performance problems by introducing a Two-Stage approach that

would efficiently discriminate between languages and easily expands to add many languages without those complex labeling or transcription work

## 1.1 Overview

The human brain is still the most advanced and accurate language identification system today, far exceeding the data storage and processing power of a supercomputer. Humans can easily determine those languages they already know or have heard before, usually with only a short duration of utterances. But using machines to automatically and accurately identify a sample of speech by an unknown speaker is still a huge challenge in speech recognition research. The biggest problem in designing a language identification system is the scalability issue. The current best systems use multiple, large vocabulary continuous speech recognizers. As these systems include a complete word and sentence level recognizer for each language, adding a new language would involve the enormous task of labeling speech to train phone recognizers. Another popular approach exploits the phonotactic properties of the languages with recognition and language modeling occurring at the phone level. Again, the requirement of labeled speech for a large subset of the language makes it very hard to extend beyond the most common languages.

We propose using a Two-Stage approach with an efficient front-end process that would narrow the unknown utterance down to the 2 or 3 top language choices, and then use other complex systems such as phone or word recognizers for an exact match. Most top performers involve heavy use of linguistic data. Our approach instead uses the acoustic models of phones in each language and finds the highest likelihood resulting from an efficient way of clustering algorithms. This dissertation reports how we adopted

speaker identification methodology via statistical pattern recognition and applies it to language identification using Vector Quantization approach for the First Pass. We show that this method would lead to substantial improvements over existing, fast systems in discriminating between different languages, and return those top candidates with a combined accuracy of over 80 percent. A next stage would then further narrow down the identification process from here.

To fully understand this approach, we first look at the differences among languages, the driving forces and the challenges.

### **1.1.1 Differences Among Languages**

When we are listening to foreign languages, most of the time we can guess the language by the phonemes. A phoneme is the term defined by linguists to classify speech into a number of abstract categories for grouping together subsets of speech sound. For example, American English has about 40 phonemes. Even though no two speech sounds, or phones, are identical, all of the phones classified into one phoneme category are similar enough so that they convey the same meaning. But there is much overlap of the phoneme sets, and there can be differences in the way the same phoneme is interpreted in two different languages. In English, letters /l/ and /r/ are two different phonemes. Then the frequency of occurrence of phones and the phonotactic rules in languages can also differ significantly. For example, phoneme clusters /sr/ and /sp/ are quite common in Tamil and German, but not in English.[5] Prosodic features, rhythm, and intonation, also vary among languages. English is known to be stress-timed, and French is known to be

syllable-timed. All these unique features are the subject of research for automatic identification of languages by machine.

### **1.1.2 The Challenge**

To formally distinguish one language from another, these are the main categories:

Phonetics –The study of speech sounds (phones) as physical entities on a sub-language basis. There are in any language a limited number of recurrent, fairly distinctive speech units. The number of phonemes in a language ranges about 15 to 50.

Phonotactic – The study of frequency distributions and combinations of phonemes. Some combinations that occur frequently in one language are illegal in another.

Prosodic – The study of sound patterns, which can be analyzed in terms of duration, pitch and stress. For example, one syllable of most words has a heavy stress or accent that sets it off from the other syllables in stress language.

Morphology – The study of ways in which words are built up from the smallest meaningful parts.

Syntax – The study of how words can be legally strung together.

The last forty years have produced technology capable of accurately processing spoken language, which has been marketed commercially in many industries. However, only a limited vocabulary and controlled environment are being utilized. Automatic language detection encounters most of the challenges of speech recognition. They are:

- Continuous flow of speech
- Volume of speech sound wave
- Paucity of information in speech sound wave

- Variability of voice and speech patterns
- Varying noise environment

In addition, languages have characteristic sound patterns that differ in the inventory of phonological units (speech sound categories) used to produce words. The frequency of occurrence of these units differs as well, and the order in which they occur in words[4] language identification, and the ability of systems to understand multiple languages are still in their infancy.

## **1.2 The Driving Force**

Only in recent years has speech recognition achieved the reliability and flexibility to attract the interest of businesses willing to invest in the infrastructure that is needed for ongoing research in this field. External forces such as increasingly powerful microprocessors complement the significant technological advances that put speech recognition system into many small devices. The seamless integration of computer-telephony technology, and successful user acceptance of machine voice, has encouraged most retail businesses to use machine routing of calls today.

Moore's law has held true in the last twenty years of the microprocessor industry, which is the primary factor in the migration of advanced speech recognition from research lab into the business world. All computational experiments for this thesis were done on a Dell laptop computer with a 500 Mhz PIII Intel chip. Not only has the cost of equipment been decreasing, but also the size of the hardware is decreasing to the point where a digital signal-processing chip can be found in practically every consumer electronic device. In the near future, one should expect all paper user manuals or



operational instructions to be replaced by machine response in audio of a user's native language.

Global business also plays an important role in creating huge demand for telephone-based, multi-lingual speech recognition systems. It would provide a cost-effective way for businesses needing 24-hour telephone message and communication support as their business is actually non-stop amidst the world's time zones.

### **1.3 Related Researches**

Most of the approaches to language identification have adopted the successfully proven techniques used in speaker-independent speech recognition systems. This paper will briefly discuss the major breakthroughs in speech recognition in groups of twenty years to give some context to the research conducted. The first attempts to devise systems for automatic speech recognition by machine were in 1950 at various research labs in the U.S. and Great Britain[6, 11, 30]. Research in language identification was actually started twenty years after the research began in the field of speech recognition.

#### **1.3.1 The first twenty years**

In the 1970's, a highly successful group effort in large vocabulary speech recognition at IBM called the New Raleigh language[40] for simple database queries yielded the laser patent text language[17] and the Tangora[16] for dictation of simple memos. The truly speaker independent research[34] was developed at AT&T labs in 1979. Here a wide range of sophisticated clustering algorithms was used to determine the number of distinct patterns required to represent all variations of different words across a wide samplings. This research has a significant influence on the later research where

refinement was achieved over the next decades, including the framework on which this thesis is based.

A significant technology shift from template-based approaches to statistical modeling methods led to the highly successful model of hidden Markov modeling (HMM)[9]. This widespread publication and theory has been applied in virtually every speech-recognition research laboratory in the world.

### **1.3.2 *The last twenty years***

In the 1990's, the major milestone for research was the availability of a public-domain multilingual speech corpus from Oregon Graduate Institute (OGI)[27]. Prior to this, the field of spoken language ID research has suffered from the lack of a common speech corpus that could be used to evaluate different approaches to the problem. A corpus is a collection of linguistic data, either compiled as written texts or as a transcription of recorded speech. The initial corpus contained 11 languages, which was used by the National Institute of Standards and Technology (NIST)[29] to conduct an annual common evaluation of spoken language ID algorithms.

Many papers have been published during the last 20 years on this subject, and model training has involved neural networks or HMMs. Zissman[49, 50] exploited the fact that a stochastic grammar for one language can be developed based on the acoustic models of a different language. Li [21] has applied speaker recognition techniques for language identification by classifying incoming utterances based on the similarity of the speaker of that utterance with the most similar speakers of the target languages. He based this similarity measure on spectral features extracted from experimentally determined

syllabic nuclei within the utterances. Muthusamy examined pitch variation within and across broad phonetic segments in his doctoral thesis[24] in which he found other prosodic information such as duration and syllabic rate to be more useful.

## **1.4 Overview of Proposed Two-Stages Approach**

It is clear that there is really neither a defacto methodology to spoken language identification, nor a standard set number of languages to be tested. Moreover, the level of performance is not ready for commercial usage. As of today, there is no application in the market that can claim to identify even 2 languages. Since human identification performance can be very accurate with a very short duration of speech[25], it is clear that better-automated systems can be developed. We propose a simple 2-stage approach; first using an efficient clustering algorithm to select the top candidates, and then following up with speech recognizers to correctly identify the language

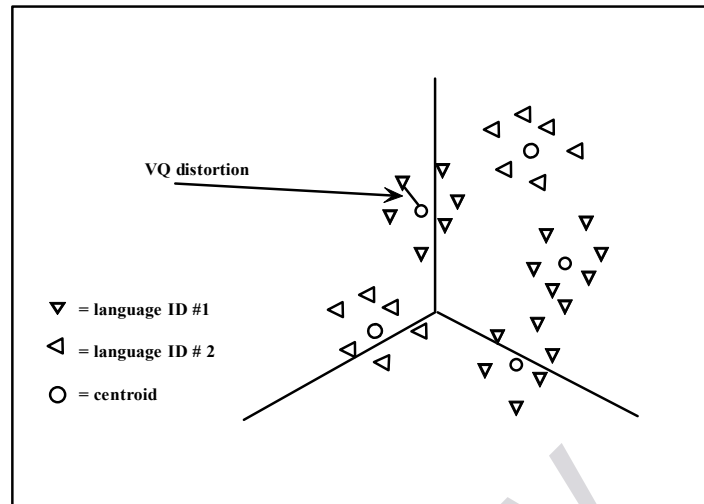
### **1.4.1 First Pass - Efficient Algorithm in Selecting Top Candidates**

In order to find the most efficient algorithm and cost effective way in selecting the top candidates among many unknown languages, we first need to understand how human beings do it. When learning to speak and understand different languages, most human beings are exposed to one particular language constantly during their “learning phase”. Over time, our data store (brain) records patterns and when the same pattern is detected, our nerve system naturally communicates with the speaker and acoustic background. Using this same basic pattern-recognition approach with machines, speech parameters are obtained using a type of spectral analysis that concentrates on the “short duration” characteristics of the speech signal. The short-time spectral measurements are performed

sequentially over time, producing a sequence of spectral feature vectors, called a speech pattern. This speech pattern is then compared with each class reference pattern and a measure of similarity between the unknown pattern and each reference pattern is calculated.

Vector Quantization (VQ) is a very efficient source-coding technique and has been widely used and successful in speech recognition research. It is a procedure that encodes a vector of input, a segment of waveform or a parameter vector that represents the segment spectrum, into an integer that is associated with an entry of a collection (codebook) of reproduction vector. The main advantage is that it requires less storage for signal analysis and reduces computation for determining similarity of speech analysis vectors. This would work well if also proven successful in language identification; as it can be easily expanded to add more languages without lots of front-end work, or transcriptions.

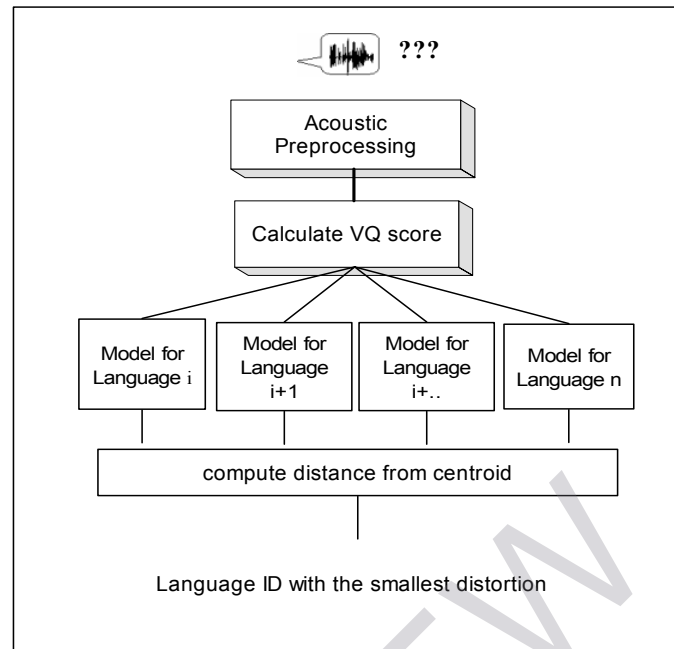
The following figure illustrates the concept with 2 languages and 3 dimensions of the acoustic space. The circles refer to the acoustic vectors from the Language ID #1 and the triangles are from #2. In the training phase, a language-specific VQ codebook is generated for each known language by clustering the training acoustic vectors. Solid circles and solid triangles for language ID #1, and #2 show the result codeword (centroids), respectively. The distance from a vector to the closest codeword of a codebook is called a VQ-distortion. During the recognition phase, an input utterance of an unknown language is “vector quantized” using each trained codebook and the total VQ distortion is computed. The language input corresponding to the VQ codebook with the smallest total distortion is identified.



**Figure 1-1**      *Conceptual diagram of Vector Quantization Distortion*

*The VQ distortion is the distance between the vector and the centroids of the same region, in this case, the same language model.*

Given an utterance of unknown language, we will use the following algorithm to calculate the score against the model for each language. The one with the smallest score is regarded as the language of this utterance. The following figure depicts the feature matching process. The unknown test utterance is matched against all language models, and the language with minimum distortion is selected as the winning language. Using the nearest-neighbor approach, we can select those top candidates with combined probability of over 80% for the next stage.



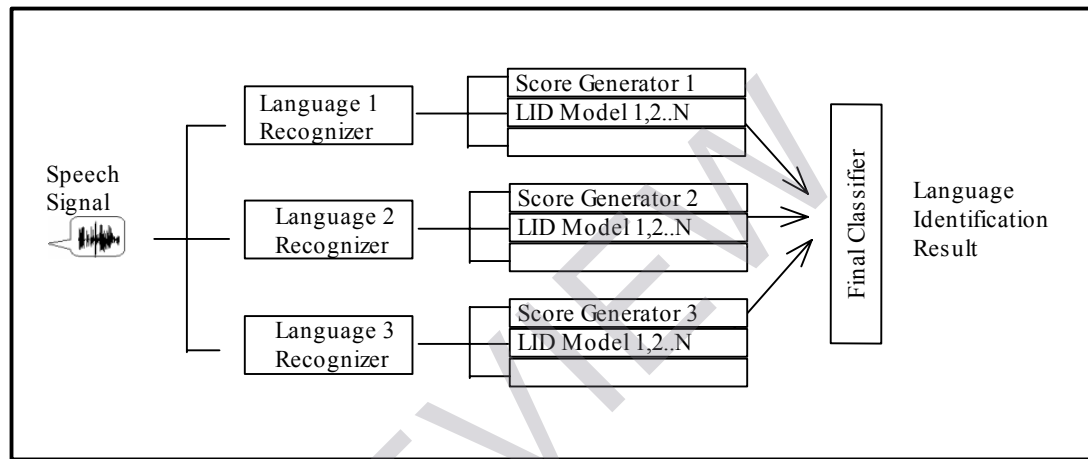
**Figure 1-2** Conceptual diagram of nearest neighbor selection

*The unknown test utterance is matched against all language models, and the language with minimum distortion is selected as the winning language*

#### **1.4.2 Second Stage – Accurately Identify unknown signal**

The second stage is the process of trying to accurately identify the unknown speech signal using various high performance approaches. Since we have narrowed down to a handle full of top candidates, there are a few options available to correctly identifying them. The most common and easiest way is to allow users to select from the option prompt generated by VoiceXML, and responded either by DTMF or speech. The multi-lingual framework by Law [20] has provided a template-based engine that would create a user-friendly native language prompt using the languages of the unknown top candidates.

Another proven successful approach, which fits nicely within our framework, in terms of feature extraction, is the parallel language dependent phone recognition by Yan [47] and Zissman et al [50]. This type of approach exploits the phonotactic properties of the languages, and also does not need to recognize words. The recognition and language modeling are done at the phone level. The following figure depicts the basic concept.



**Figure 1-3** Overview of a Language-Dependent Phone Recognition system

*The approach that exploits the phonotactic properties of the languages and not words*

This approach is able to achieve identification accuracy in excess of 80% using 10-second utterances in 6 languages [47], and over 90% on 3 or less languages. The biggest drawback is the requirement of labeled speech for a large subset of languages used. Since it is based on multiple language-specific phone recognizers, it requires labeled speech to train those recognizers.

## 1.5 Outline of the Dissertation

This chapter briefly describes an idea of how and why we use this simple Two-Stage approach for language identification. Chapter 2 will describe speech recognition

and the most common recognition applications. We will then present comparative studies, experiments, and results of those approaches we have discussed earlier. Chapter 3 provides details on the Vector Quantization approach, and the data set from the Oregon Graduate Institute and corpus structure. The technical specifications and implementation process is outlined in Chapter 4. Details of the improvement in both configuration and algorithms, and how results were generated for a 5-language set (NIST 1996), 7-language set and 10-language set as the baselines. We then performed initial improvements on various configurations, which are described in Chapter 5. The use of various distance measures and data fusion to the output scoring is detailed in Chapter 6 and revised baselines are generated. Conclusions and future work, based on our findings from this research, are discussed in chapter 7. The last section is the appendix and references.