

# Detection of Foreign Words and Names in Written Text

By

Bashir U Ahmed

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 2005

UMI Number: 3172339

Copyright 2005 by  
Ahmed, Bashir U.

All rights reserved.

PREVIEW  
UMI<sup>®</sup>

---

UMI Microform 3172339

Copyright 2005 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  
P.O. Box 1346  
Ann Arbor, MI 48106-1346

## Dissertation Signature (Approval) Page

We hereby certify that this dissertation, submitted by Bashir U. Ahmed, satisfies the dissertation requirements for the degree of Doctor of Professional Studies and has been approved.

---

Dr. Charles Tappert

Chairperson of Dissertation Committee

10 May 2005

---

Dr. Fred Grossman

Dissertation Committee Member

10 May 2005

---

Dr. Sung-Hyuk Cha

Dissertation Committee Member

10 May 2005

School of Computer Science and Information Systems  
Pace University 2005

# **Abstract**

## **Detection of Foreign Words and Names in Written Text**

By

Bashir U Ahmed

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 2005

Tremendous research effort has gone into the field of natural language processing and understanding during the last half of the 20<sup>th</sup> century, yet it remains as one of the most challenging problems. Nevertheless, many companies have commercialized language processing applications such as Text-to-Speech (TTS) systems, domain-specific Automatic Speech Recognition (ASR) systems, Machine Translation (MT) systems, and limited-domain, speaker-dependent Speech-to-Speech (STS) prototype systems. While the quality of all natural language processing applications has improved steadily, they remain far from being perfect.

Due to globalization and the widespread use of the Internet, occurrences of foreign entities – foreign words, names, locations, and events – in written text are becoming commonplace. This further complicates the automatic processing of natural language text. Identification of such foreign entities is important for several reasons. For example, TTS systems will usually fail to pronounce such entities properly in their native language as they try to sound them using the lexicon of the base language. These foreign entities also cause problems for MT programs where translation is initially done word by word.

This dissertation simplifies the Naïve Bayesian, N-gram based text-classification algorithm by taking its logarithm, naming it the Cumulative Frequency Addition (CFA) algorithm, and applies it to three text processing tasks: language identification, name-to-nationality identification, and detection of foreign words and names. In the language identification task CFA yields 100% accuracy on string sizes greater than 150 characters. In the name-to-nationality task, it yields 86% accuracy on a 14 country database and 96% on a 7 country database within the top three choices. Identifying the nationality, or at least the language group, of a person from his/her name, can be important not only for proper pronunciation but also for purposes of forensics studies and national security. Finally, in the task of detecting foreign words we obtain 66.9% accuracy. This is the first study to apply natural language processing techniques to the latter two tasks.

## Acknowledgments

I am a dreamer and I dreamt really big for my thesis topic wanting to develop a Universal Language Translator. I must thank my advisors for bringing me back to reality. First, I would like to thank Dr. Charles Tappert, my thesis committee chairman, who has diligently invited scholars and industry experts in the latest technological fields, which has further sparked my interest in natural language processing. I also thank Dr. Fred 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 Dr. Sung-Hyuk Cha for helping me shape up the thesis idea and also providing the initial set of training data.

As it took great advice from several great advisors to finish the intellectual work, it took similar amount of support and encouragement at the home front to get through the wee hours at night. I am thankful to my seven years old daughter, Serina, for being so understanding while working on my thesis chapters and always reminding me that I was the best Dad in the “whole wide world”. My three years old daughter, Samina, provided ultimate companionship 12 past midnight rolling over my back while I was working on my chapters lying down on the bed. I surely will miss those moments and cherish them the rest of my life. The arrival of my son, Aynis, during the third year of my doctoral research could have been a show breaker, but it did not and it became a joy because of one person – my wife. I could not thank my wife, Maria, enough for being so understanding and always wanting to see me finish the thesis. Maria’s utmost confidence, “You can do this” and unselfish support helped me making my dream a reality.

My dream of going for a doctoral program began a long time ago – my grandfather seeded the dream when he promised my mother that he would finish his Ph. D. with my elder brother. My grandfather passed away without fulfilling his dream. I dared to take this challenge because I knew I would have all the needed support. I am thankful to God for giving me so much in life – my wife’s determination to see me finish the program, and my brother (my toughest competition inspired by sibling rivalry), Nasir, whose unwavering mental support boosted my confidence. I am ever so grateful to my mother, Mazeda Akter, for raising me so well and believing in me so much. I am thankful to my uncle, Mohammad Z. Haque, for sponsoring us to come to America and for being so supportive in every step of the way. I would like to thank my mother & father-in-law for believing and giving me the confidence I needed to finish. My gratitude goes to all of my brothers & sisters and their husbands and wives for always giving encouragement and helping me get through this program. I could not have finished this work if I did not have the friendship and mental support from my fellow classmates, especially Melanie Johnson and Barbara Edington for being such a good listener, Steven Golikov and Borming Chiang for providing technical support when needed.

Finally, I thank Pace University for developing this DPS program and making my dream a reality.

## Table of Contents

|   |     |
|---|-----|
| <i>Abstract</i> .....   | iii |
| <i>Acknowledgments</i> .....  | iv  |
| <i>Table of Contents</i> .....  | v   |
| <i>List of Tables</i> .....   | x   |
| <i>List of Figures</i> .....  | xii |
| CHAPTER 1 .....   | 1   |
| INTRODUCTION .....  | 1   |
| 1.1. OVERVIEW .....   | 1   |
| 1.1.1. <i>Implications of Automatic Processing of Natural Language Text</i> ..... | 3   |
| 1.1.2. <i>Difference Among Languages</i> .....                                    | 4   |
| 1.2. THE PROBLEM.....   | 5   |
| 1.2.1. <i>Precise Problem Definition</i> .....                                    | 6   |
| 1.2.2. <i>Precise Problem Definition of the Extension Work</i> .....              | 8   |
| 1.3. STATE OF THE ART OF TEXT-TO-SPEECH .....                                     | 8   |
| 1.3.1. <i>Text-To-Speech: What For?</i> .....                                     | 8   |
| 1.3.2. <i>How Current Text-To-Speech Read Text?</i> .....                         | 9   |
| 1.3.3. <i>Components of Text-To-Speech System</i> .....                           | 12  |
| 1.3.4. <i>Multilingual Text-To-Speech Systems</i> .....                           | 14  |
| 1.3.5. <i>Polyglot Text-To-Speech Systems</i> .....                               | 14  |
| 1.4. THE CHALLENGE –.....   | 14  |
| 1.4.1. <i>Multilingual Text</i> .....   | 15  |
| 1.4.1.1. <i>Sample Text Containing two Languages</i> .....                        | 15  |
| 1.4.2. <i>Type of Inclusions</i> .....  | 16  |
| 1.4.2.1. <i>Foreign Words Inclusion</i> .....                                     | 16  |
| 1.4.2.2. <i>Other Types of Inclusions</i> .....                                   | 18  |
| 1.4.2.3. <i>Names, Locations, Events Etc</i> .....                                | 20  |
| 1.5. AUTOMATION CHALLENGES .....  | 21  |
| 1.5.1. <i>Text-To-Speech Challenges</i> .....                                     | 22  |
| 1.5.1.1. <i>Special Case – Name Pronunciation</i> .....                           | 22  |
| 1.5.2. <i>Machine Translation Challenges</i> .....                                | 24  |
| 1.5.3. <i>Information Retrieval Challenges</i> .....                              | 26  |
| 1.5.4. <i>Problem Simplification and its Justification</i> .....                  | 26  |
| 1.6. APPROACH TO THE SOLUTION .....   | 27  |
| 1.6.1. <i>Three Steps</i> –.....  | 28  |
| 1.6.1.1. <i>Language Identification Module</i> .....                              | 28  |
| 1.6.1.2. <i>Identification of Nationality From Names</i> .....                    | 29  |
| 1.6.1.3. <i>Detection of Foreign Words and Names</i> .....                        | 29  |
| 1.7. OUTLINE OF THE DISSERTATION .....  | 30  |
| CHAPTER 2 .....   | 31  |
| RELATED STUDIES AND EXPERIMENTS.....  | 31  |
| 2.1. LANGUAGE IDENTIFICATION .....  | 31  |

|  |    |
|--|----|
| 2.1.1. The Data.....   | 31 |
| 2.1.1.1. Input.....  | 31 |
| 2.1.1.2. Output.....   | 32 |
| 2.1.2. Preprocessing.....  | 32 |
| 2.2. EXISTING LITERATURE ON LANGUAGE IDENTIFICATION AND FOREIGN WORD<br>DETECTION..... | 32 |
| 2.2.1. Language Identification Approaches.....   | 32 |
| 2.2.1.1. Most Frequently Occurring Unique Letter Combinations.....                     | 33 |
| 2.2.1.2. Most Frequently Occurring Short Word Method.....                              | 34 |
| 2.2.1.3. Grammatical Word Method.....  | 36 |
| 2.2.1.4. Using the Alphabet.....   | 36 |
| 2.2.1.5. Vector Quantization Using Character ASCII Values.....                         | 37 |
| 2.2.1.6. N-Gram Based Method.....  | 37 |
| 2.2.2. Foreign Words and Names Detection.....  | 38 |
| 2.3. CLASSIFICATION ALGORITHMS.....  | 39 |
| 2.3.1. N-Gram Rank Order Statistics.....   | 39 |
| 2.3.2. K-Nearest Neighbor.....   | 41 |
| 2.3.3. Naïve Bayesian.....   | 43 |
| 2.3.4. Clustering.....   | 45 |
| 2.3.5. Neural Networks.....  | 46 |
| 2.4. PERFORMANCE MEASURES.....   | 47 |
| 2.4.1. Multiple Binary Classification Tasks.....                                       | 48 |
| 2.4.1.1. Precision.....  | 49 |
| 2.4.1.2. Recall.....   | 49 |
| 2.4.1.3. Fallout/False Positives.....  | 50 |
| 2.4.1.4. Accuracy.....   | 50 |
| 2.4.1.5. Error.....  | 50 |
| 2.4.1.6. Breakeven Point.....  | 51 |
| 2.4.1.7. F-Measure.....  | 52 |
| 2.4.2. Multi-Class And Multi-Label Classification.....                                 | 53 |
| 2.5. ALGORITHM AND FEATURE SELECTION.....  | 54 |
| CHAPTER 3.....   | 55 |
| SYSTEM DESIGN AND IMPLEMENTATION.....  | 55 |
| 3.1. INTRODUCTION.....   | 55 |
| 3.2. N-GRAM BASED APPROACH.....  | 55 |
| 3.2.1. Why N-Grams?.....   | 56 |
| 3.2.1.1. Advantages Of N-Gram Based Methods.....                                       | 56 |
| 3.2.1.2. Disadvantages Of N-Gram Based Methods.....                                    | 57 |
| 3.3. WHY NOT DICTIONARY?.....  | 57 |
| 3.4. CAN DICTIONARY HELP?.....   | 58 |
| 3.5. DATABASE BASED IMPLEMENTATION.....  | 58 |
| 3.5.1. Database Design.....  | 59 |
| 3.5.2. Database Description.....   | 60 |
| 3.5.3. The Language Identification Module.....   | 60 |
| 3.5.3.1. Training N-Gram Table.....  | 60 |

|                              |   |     |
|------------------------------|---|-----|
| 3.5.3.2.                     | Test Staging Table .....  | 62  |
| 3.5.3.3.                     | Analysis Header .....   | 62  |
| 3.5.3.4.                     | Analysis Detail .....   | 63  |
| 3.5.3.5.                     | URL_Files Table .....   | 63  |
| 3.5.4.                       | <i>The Identification of Nationality From Names Module</i> .....        | 64  |
| 3.5.4.1.                     | Athens2004olympic Table .....   | 64  |
| 3.5.4.2.                     | Coaches2004Athens Table .....   | 66  |
| 3.5.5.                       | <i>Foreign Words Detection From Native Text Module</i> .....            | 67  |
| 3.5.5.1.                     | Dictionaries Table .....  | 67  |
| 3.5.5.2.                     | Mixed_File_Source Table .....   | 68  |
| 3.6.                         | THE USER INTERFACE .....  | 70  |
| 3.6.1.                       | <i>Behind The Scene</i> .....   | 71  |
| 3.7.                         | CLASSIFICATION .....  | 74  |
| CHAPTER 4                    | .....   | 75  |
| LANGUAGE CLASSIFICATION TASK | .....   | 75  |
| 4.1.                         | INTRODUCTION .....  | 75  |
| 4.2.                         | LANGUAGE IDENTIFICATION .....   | 76  |
| 4.2.1.                       | <i>Methodology</i> .....  | 76  |
| 4.2.2.                       | <i>Collection Of Text Samples And Creation Of N-Gram Profiles</i> ..... | 77  |
| 4.2.3.                       | <i>N-Gram Frequency Calculation</i> .....                               | 79  |
| 4.2.4.                       | <i>N-Gram Rank Ordering</i> .....                                       | 81  |
| 4.3.                         | TESTING PROCEDURES .....  | 82  |
| 4.3.1.                       | <i>Classification by Rank-Order Statistics</i> .....                    | 82  |
| 4.3.2.                       | <i>Classification Using Cumulative Frequency Addition</i> .....         | 87  |
| 4.3.3.                       | <i>Classification Using Naïve Bayesian Classifier</i> .....             | 90  |
| 4.4.                         | RESULTS .....   | 91  |
| 4.5.                         | PITFALLS OF N-GRAM BASED CLASSIFIERS – FALSE POSITIVES .....            | 94  |
| 4.5.1.                       | <i>Eliminating False Positives</i> .....                                | 94  |
| 4.5.2.                       | <i>Language Similarity Groups</i> .....                                 | 96  |
| 4.6.                         | CONCLUSION .....  | 97  |
| CHAPTER 5                    | .....   | 99  |
| NAME IDENTIFICATION TASK     | .....   | 99  |
| 5.1.                         | WHY NAME RECOGNITION MATTERS? .....                                     | 99  |
| 5.2.                         | WHAT’S IN A NAME? .....   | 99  |
| 5.3.                         | EXTRACTING NAMES FROM TEXT .....  | 100 |
| 5.3.1.                       | <i>Commercial Products</i> .....  | 100 |
| 5.3.1.1.                     | Nominator .....   | 100 |
| 5.3.1.2.                     | Namefinder .....  | 100 |
| 5.3.1.3.                     | Nametag .....   | 101 |
| 5.3.2.                       | <i>Academic Systems</i> .....   | 101 |
| 5.3.2.1.                     | Proper Name Analysis .....  | 101 |
| 5.3.2.2.                     | Contextual Evidence in Name Analysis .....                              | 103 |
| 5.3.2.3.                     | Processing Names Without a Name Database .....                          | 104 |
| 5.3.2.4.                     | Proper Nouns in Information Retrieval .....                             | 104 |



|   |   |     |
|---|---|-----|
| 5.3.2.5.  | Survey of the Need for Personal Name-Matching Algorithms [Borgman et al.] | 104 |
| 5.4.  | LITERATURE SUMMARY AND THE VALIDITY OF THIS WORK                          | 105 |
| 5.5.  | NATIONALITY IDENTIFICATION FROM NAME                                      | 105 |
| 5.5.1.  | <i>Methodology</i>  | 105 |
| 5.5.1.1.  | Collection of Name Samples and Creation of N-Gram Profiles                | 106 |
| 5.5.1.2.  | Pre-Processing of Names   | 107 |
| 5.5.1.3.  | Creating the Training Databases   | 108 |
| 5.5.1.4.  | N-Gram Frequency Calculation  | 110 |
| 5.5.2.  | <i>Testing Procedures</i>   | 111 |
| 5.5.2.1.  | Name Classification by Cumulative Frequency Addition                      | 111 |
| 5.5.2.2.  | Name Classification Using Naïve Bayesian Classifier                       | 112 |
| 5.5.3.  | <i>Results</i>  | 114 |
| 5.5.3.1.  | Fourteen Country Task   | 115 |
| 5.5.3.2.  | Seven Country Task  | 116 |
| 5.5.3.3.  | Name Analysis Confusion Matrix  | 117 |
| 5.6.  | CONCLUSION  | 118 |
| 5.7.  | WHAT'S NEXT?  | 119 |
| CHAPTER 6   |   | 121 |
| DETECTION OF FOREIGN WORDS AND NAMES IN BASE TEXT |   | 121 |
| 6.1.  | INTRODUCTION  | 121 |
| 6.2.  | PREPARING THE TEST SAMPLE   | 122 |
| 6.3.  | PSEUDO CODE FOR THE ARTIFICIAL CREATION OF MIXED TEXTS                    | 123 |
| 6.4.  | THE BASIC ALGORITHM   | 126 |
| 6.4.1.  | <i>Identification of the Base Language</i>                                | 128 |
| 6.4.2.  | <i>Verification by Base Dictionary</i>                                    | 129 |
| 6.4.3.  | <i>Processing Words not Found in the Base Dictionary</i>                  | 130 |
| 6.4.3.1.  | Is it a Name?   | 130 |
| 6.4.3.2.  | What is the Nationality?  | 131 |
| 6.4.3.3.  | Is it a Misspelled Word or a Foreign Word?                                | 131 |
| 6.4.3.4.  | What Language is it?  | 131 |
| 6.4.3.5.  | An Example Walk Through   | 132 |
| 6.5.  | RESULTS   | 137 |
| 6.5.1.  | <i>The Process of Evaluation</i>  | 137 |
| 6.5.1.1.  | Recall  | 139 |
| 6.5.1.2.  | Precision   | 140 |
| 6.5.1.3.  | Loss  | 140 |
| 6.5.1.4.  | False Positives   | 143 |
| 6.5.2.  | <i>Limitation of the Approach</i>   | 143 |
| 6.6.  | SUMMARY AND DISCUSSION  | 145 |
| CHAPTER 7   |   | 147 |
| CONCLUSION AND FUTURE WORK                        |   | 147 |
| 7.1.  | SUMMARY OF THE FINDINGS   | 147 |
| 7.2.  | FUTURE WORK AND CHALLENGES  | 148 |

|  |     |
|--|-----|
| 7.2.1. <i>Identification of Short Sentence Fragments</i> .....                                   | 149 |
| 7.2.2. <i>NLP Application with Other Existing Biometric Measures for Forensic Analysis</i> ..... | 150 |
| 7.2.3. <i>Compilation and Publication of Mixed Lingual Lexicon for Further Research</i><br>150   |     |
| 7.2.4. <i>Further Improvement on Accuracy</i> .....  | 150 |
| APPENDIX.....  | 152 |
| GLOSSARY OF TERMS AND ACRONYMS.....  | 152 |
| REFERENCES .....   | 155 |

PREVIEW

## List of Tables

|   |     |
|---|-----|
| Table 1.1 Multilingual E-mail thread.....   | 16  |
| Table 2.1 Unique letter combinations proposed by various authors.....   | 34  |
| Table 2.2 Most Common Short Words.....  | 35  |
| Table 3.1 Training database for language identification .....   | 61  |
| Table 3.2 Header table for storing results from language identification .....                                     | 62  |
| Table 3.3 Detail table for storing results from language identification .....                                     | 63  |
| Table 3.4 URL location of the test samples for the language identification.....                                   | 64  |
| Table 3.5 Table for storing training names for the identification of nationality.....                             | 65  |
| Table 3.6 Table for storing test names for the identification of nationality .....                                | 67  |
| Table 3.7 Table for storing dictionary words.....   | 67  |
| Table 3.8 Table for storing the names for multilingual text samples .....   | 69  |
| Code Listing 3.1 Sample code for classification of a test string.....   | 74  |
| Table 4.1 Distance calculation using rank-order statistics, reproduced from [10]. .....                           | 76  |
| Table 4.2 Training sample size and the corresponding N-gram statistics. ....                                      | 79  |
| Table 4.3 A sample of how internal and overall frequencies are calculated. ....                                   | 80  |
| Table 4.4 Sample data with calculated internal and overall rank orders, and internal and overall frequencies..... | 81  |
| Table 4.5 List of test N-grams with ranking (List extracted from the test string “Bon Jour”)                      | 83  |
| Table 4.6 Candidate N-grams from the string “Bon Jour” with their rank-order statistics.                          | 86  |
| Table 4.7 Classification of the string “Bon Jour” as French using the rank-order statistics method.....           | 86  |
| Table 4.8 Candidate N-grams from the string “Bon Jour” with their internal and overall frequency statistics.....  | 89  |
| Table 4.9. Classification of the string the string “Bon Jour” as French using the cumulative frequency sum. ....  | 90  |
| Table 4.10 Classification of the string the string “Bon Jour” as French using the Naïve Bayesian Classifier.....  | 90  |
| Table 4.11 Summary of Results.....  | 92  |
| Table 4.12 Rank-Order Statistical Results – Internal rank order vs. overall rank order...                         | 92  |
| Table 4.13 Speed of classification for cumulative frequency addition versus rank-order statistics.....            | 93  |
| Table 4.15 Language confusion matrix .....  | 96  |
| Table 4.16 Language similarity groups .....   | 96  |
| Table 5.1 Parsing names into first, last, and middle names.....   | 108 |
| Table 5.2 Training sample size and the corresponding N-gram statistics. ....                                      | 110 |
| Table 5.4 Test N-grams extracted from the first name “charles”. ....  | 111 |
| Table 5.5 Candidate N-grams from the first name “Charles” with their normalized frequency statistics.....         | 112 |
| Table 5.6 Candidate N-grams matrix from the first name “Charles” with their normalized frequency statistics.....  | 113 |
| Table 5.7 Classification of name sample data showing the top 3 choices by the CFA method.....                     | 113 |

|  |     |
|--|-----|
| Table 5.8 Classification of name sample data showing the top 3 choices by the NBC method.....                    | 114 |
| Table 5.9 Number of test names per country .....   | 115 |
| Table 5.10 Confusion matrix in % top-choice identification (names from Great Britain are included in USA). ..... | 118 |
| Table 6.1 Sample code for detecting suspect words .....  | 126 |
| Table 6.2 Ranked language list for the sample string.....  | 133 |
| Table 6.3 Identification of “Weltherztag” .....  | 134 |
| Table 6.4 Identification of “WHO” .....  | 134 |
| Table 6.5 Identification of “World Heart Federation” .....   | 136 |
| Table 6.6 Summary of results – Combination of Dictionary and Iterative CFA analysis .....                        | 138 |
| Table 6.7 Summary of results –Dictionary only approach .....   | 139 |
| Table 6.8 Sample list of common words listed in multiple language dictionaries .....                             | 143 |

## List of Figures

|  |     |
|--|-----|
| Figure 1.1 A simple functional diagram of Text-to-speech system.....                     | 12  |
| Figure 1.2 The NLP module of a general Text-To-Speech conversion system.....             | 13  |
| Figure 1.3 E-mail containing words from 4 different languages .....                      | 18  |
| Figure 2.1 Dataflow For N-Gram-Based Text Categorization .....                           | 41  |
| Figure 2.2 K-nearest neighbor example.....   | 42  |
| Figure 2.3 Mathematical definitions of performance metrics.....                          | 51  |
| Figure 3.1 ER diagram of the database.....   | 59  |
| Figure 3.2 Sample of URL_File Content.....   | 64  |
| Figure 3.3 Sample of Dictionary Content.....   | 68  |
| Figure 3.4 Sample of Mixed_file_Source Content .....                                     | 69  |
| Figure 3.5 User interface for Interactive Testing.....                                   | 71  |
| Figure 3.6 Typical Results with performance statistics.....                              | 71  |
| Figure 4.1 Percent accuracy of classification of NBC, CFA, and rank-order statistics.... | 93  |
| Figure 4.2 Linear relationship between CFA weights and the length of the test string.... | 95  |
| Figure 5.1 Percent accuracy of classification by CFA – 14 Country Task.....              | 115 |
| Figure 5.2 Percent accuracy of classification by NBC – 14 Country Task .....             | 116 |
| Figure 5.3 Percent accuracy of classification by CFA – 7 Country Task.....               | 116 |
| Figure 5.4 Percent accuracy of classification by NBC – 7 Country Task .....              | 117 |
| Figure 6.1 Basic algorithm used to detect foreign words .....                            | 128 |
| Figure 6.2 Sample of inserted and recovered foreign words .....                          | 138 |

## Chapter 1

### Introduction

#### 1.1. Overview

Due to the vast variability in characters, pronunciations, accent and letter-to-sound rules, natural language processing (NLP) task is viewed as incredibly difficult [5]. There are 82 distinct languages listed as the world's major spoken languages, each with its own set of linguistic rules and each spoken by more than 10 million primary or alternative speakers [63]. In recent years, much progress has been made in NLP research largely due to the rapid growth in computing power. Research effort in NLP in both academic and commercial settings has increased greatly due to globalization and the need to communicate internationally. This resulted in many commercial speech recognition, language identification, language translation, and text-to-speech systems. This trend has created awareness of language processing (Speech recognition, language translation and text-to-speech) as a viable commercial tool.

Globalization is manifested in e-mails, corporate documents, and newspaper articles composed in mixed languages. Investigation of newspapers articles by Pfister et.al. [50] and our own, found numerous inclusion of English text in written German, Tagalog and Swedish newspapers. Our own investigation of internal e-mails from a German company

with sister location in the United States, found numerous German inclusions in English e-mails. This is an unavoidable reality because when a native German speaker composes an e-mail in English, some German words get included without the author's conscious knowledge. All our findings confirm that the phenomenon of mixed-linguality occurs more frequently than one would normally realize. This places a great difficulty on the text-to-speech, Machine Translation and Information Retrieval systems.

Natural sounding text-to-speech systems require accurate phonological, prosodic, morphological and syntactic knowledge, which are language specific [70]. World languages differ widely from each other on these properties. Thus, to be able to read a mixed lingual text, a text-to-speech system must know the identity and context of each word in the text before pronouncing them accurately. Language identification, which can be done using either voice input or text input, is the pre-requisite of any such system. In this study, a language identification system was implemented and applied to solve the foreign word detection problem.

Corporate America is constantly looking to increase its bottom line by finding ways to increase productivity, and many corporations are trying to leverage Automatic Speech Recognition and text-to-speech technology for customer service. Another area where text-to-speech is being used is accessibility applications – many government organizations, such as IRS in the USA and Australian Government, started to use ASR/text-to-speech to train employees who are visually impaired or have typing disabilities. Many e-mail and voice mail providers, such as Verizon and SBC

communications, provide voice enabled e-mail options where users can check their e-mail via telephone or check their voice mail from their computer. Both the voice mail and e-mail contents are read by text-to-speech software. However, as text-to-speech reading is highly language dependent, a monolingual text-to-speech system fails to provide natural sounding reading of user's voice mail or e-mail that are embedded with foreign words and names. In most cases, embedded foreign words are read in garbled manner [51]. An automatic detection and tagging of foreign words in written text can play an important role in the quality of a text-to-speech system in which it acts as the communicator to the end user. For example, when a text-to-speech system is turned on in English mode and it encounters a German word, if the system can detect that the next word is German, it can then automatically switch to German lexica and pronounce the word naturally as a normal human reader would do. This could increase the user acceptance of text-to-speech systems and make commercial application more appealing. The goal of this research is to provide a mechanism to tag the foreign words, so that a text-to-speech system will have the option to switch language when appropriate. In addition to text-to-speech, Machine Translation and Information Retrieval applications will also benefit from this research.

#### **1.1.1. Implications of Automatic Processing of Natural Language Text**

Automatic processing of natural language text has huge implications in commercial applications. For example, a document service provider may encounter numerous electronic requests in varieties of languages and sorting those requests manually would be extremely inefficient or impossible. Using automatic language identification, the service provider can build different request queues sorted by languages and then forward these



queues to appropriate sub-contractors such as a sister company in another country where people understand the language. This arrangement would allow a document service provider to become a truly global company. Automatic processing of natural language text is even more important when the input needs to go through some sort of transformation along the process such as text to speech and vice-versa. In order to implement a Voice Enabled Universal Language Translator, which does not exist yet, given any text or voice sample, the translator should be able to speak it in any other human language of the world. One approach for such a system could be the translation of the source language text to the target language text, and then read the output using text-to-speech. If the source text is in pure form, this process may work very well. However, if the source text is embedded with words from a third language, then the translation will be incorrect if those foreign words are not recognized and addressed properly. In this thesis, an algorithm to tag foreign words and names in native text was described. This was done first by identifying the base language from unknown text using CFA, and then applying a new algorithm to detect any foreign words that do not belong to the base language. Findings from this research would enable easier implementation of true Polyglot text-to-speech and Machine Translation systems where any free form source text can be read or translated accurately without the pre-existing knowledge of the source text.

### **1.1.2. Difference Among Languages**

While all the Latin character based languages share a great common part – most of the 26 letters of the English alphabets are found in all western European languages- they vary to some extent with additional characters such as the accented characters and also in text

patterns. However, the differences become very pronounced when one talks about the spoken form of these languages. In spoken form, 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 [39]. 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. Each language has its own letter-to-sound rules [68]. 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. Prosodic features, rhythm, and intonation, also vary among languages [15]. 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.

## **1.2. The Problem**

The problem this thesis addresses falls into the area of natural language processing and understanding. The title of the thesis “Detection of foreign words and names in written text” summarizes the problem accurately. Existing Text-to-speech and machine translation program performs well when processing monolingual text (text written in only one language). However, these systems perform poorly when the input text contains

words and names from foreign languages that cannot be found in the base dictionary of the main text. Therefore, this thesis investigates the problem of foreign words detection in written text.

### 1.2.1. Precise Problem Definition

**Input:** a text string – a sequence of words from one or more languages.

**Output:** the input text string with the language of each word identified.

**Assumptions:** There is a clearly dominant or main language of the input text string, and we call it the base language. Although this does not cover all possibilities, this is the usual case for this problem. Words not belonging to the base language are considered “Foreign” and identified.

**Example:** Given a test string containing 9 English words, 4 German words, and 3 French words. According to our definition, English would be considered the dominant language while German and French would be considered foreign languages. Given the input described above, the target output we tried to achieve is a tagged string as follows

Input:

“Sentence artificially constructed for der Zweck der Demonstration, la traduction pas be accurate, please pay attention”

Output:

<Eng> Sentence artificially constructed for <Ger> der Zweck der Demonstration </Ger>  
<Fre> la traduction pas. </Fre> be accurate, please pay attention </Eng>

The idea is to tag the overall sentence as English but containing German and French words as foreign words. This means there will be a beginning <Eng> tag and an ending </Eng> tag plus the beginning and ending tags for other languages inside the body of the output. Any word or sentence fragment inside the body of the output that is not explicitly tagged would be assumed to be English. This would allow automatic natural language processing application like Text-to-speech or machine translation programs to start processing in English mode and then switch to the German or French only when needed.

With this approach, any common words – words that can be found in multiple languages – would default to the base language.

The problem described above is a sub-problem of a more general one - detection of language shift in written text where no notion of native or foreign is considered. Under this general scenario a text can contain any number of words from any number of languages and the challenge is to tag each word or sentence fragments containing contiguous set of words that belong to each language. A simplified version of this general problem, the one that is most common, was implemented by identifying the base language of the majority of the text, and then identifying the foreign words relative to that language.

In addition to the foreign words, an extension of this work involves the identification of the nationality or language group of proper names. This is relevant because written text often contains proper names, which require special processing. Knowing the nationality or language of a proper name can be helpful for TTS and machine translation programs.

### 1.2.2. Precise Problem Definition of the Extension Work

**Input:** A proper name such as person, organization, event, etc., in written form.

**Output:** The proper name identified as belonging to a country and subsequently associating the country to a language.

**Assumption:** The name has been tagged as a name, and then CFA was used to identify the nationality or at least the language group of the name. Note that the focus here was not on the name delimitation task, but on the identification of the nationality of the names. In the absence of a name recognition algorithm, however all sequential words that start with capital letters were assumed to be names, except for the initial words of each sentence if not immediately followed by words with capital letters.

**Example:** Given a sentence such as “My advisor is Sung-Hyuk Cha, and he is very technical,” according to the problem definition, Sung-Hyuk Cha will be identified as a name and will be tagged as Korean. The word “My” would not be considered as a name.

## 1.3. State-of-the-Art of Text-To-Speech

### 1.3.1. Text-To-Speech: What For?

**Text-to-speech** is the creation of audible speech from computer readable text. The qualities of some commercial Mono/Multi Lingual Text-to-speech systems improved a great deal and they are being used in different critical applications. Red Planet, the Warner Bros. movie released on 11/09/2000, used AT&T’s Text-to-speech as the voices of the space ship (female) and the astronauts' space suits (male). The movie “I am Sam”

starring Sean Penn and Michelle Pfeiffer, used AT&T's "Crystal" as the voice of Pfeiffer's voice-activated cell phone. In addition to its use in the movie industry, Text-to-speech systems are being investigated in the following areas [24]: Telecommunication Services, Language Education, Aid to handicapped persons, Talking books and toys, Vocal Monitoring, Multimedia, man-machine communication, and Fundamental and applied research.

### **1.3.2. How Current Text-To-Speech Read Text?**

Synthesized speech is both a triumph of technology and the fruition of a very old dream. The first "acoustic-mechanical speech machine" was introduced in 1791 by the Viennese researcher Wolfgang von Kempelen. The machine simulated the major consonant and vowel sounds with an array of vibrating reeds, like a musical instrument. But not until the advent of electronics did machines truly begin to mimic human voices [8]. In the 1950s, researchers labored to model the acoustics of the human vocal tract and the resonant frequencies, or *formants*, it generates. This approach eventually led to workable but robotic results—certainly nothing a public-relations person would call customer ready. Stephen Hawking's voice synthesizer is the most famous example. Such a voice is not good enough for normal usage.

In the 1970s, researchers at what was then Bell Labs turned to a "concatenative" approach: Instead of trying to generate a human voice from scratch, they would start with an existing voice—several hours' worth of standard English sentences spoken by a clear-voiced person—and design a computer program to splice and re-splice it to say whatever words they wanted said.

The computer program first parsed the prerecorded sentences into consonant and vowel sounds, called phonemes—perhaps 50 or 60 in the early iterations. Then the phonemes were reassembled to form new words. The recorded word *cat*, for instance, could be deconstructed into the phonemes *k*, *ae*, and *t*, which could then be rearranged to form *tack*. It worked, and it was a definite improvement over robot-speak, but it wasn't any closer to sound of Peter Jennings or whatever Journalist you like. Fifty-odd phonemes simply couldn't capture the subtle intonations of spoken language.

In the mid-1990s, armed with a new generation of supercomputers, AT&T researchers began amassing a vast digital "voice warehouse" of phonemes. Instead of one *t* sound for the computer program to choose from, there might be 10,000. By having so many sounds, it offers a little more spontaneity. By parsing phonemes into "half-phones" it is possible to offer subtler possibilities for recombination. Voice synthesis now entails properly labeling the half-phones—10,000 versions of the " $t_1$ " sound, 10,000 versions of the " $t_2$ " sound, and so on—then creating a computer algorithm to smoothly string them into words and sentences. By assembling a simple word like *cat* from its half-phones—(" $k_1$ ,  $k_2$ ,  $a_1$ ,  $a_2$ ,  $t_1$ ,  $t_2$ ")—involves billions of combinatorial decisions and presents a massive computer-processing problem [8].

Allistar Conkie of AT&T is generally credited with devising a workable solution, now known as unit-selection synthesis [8]. His solution was to assign "costs" to the innumerable choices and combinations of half-phones. Charting the "least expensive"

path through the chorus of half-phones became simply a math problem for the computer to work out. Most costs crop up where two half-phones meet and attempt to join. The computer can measure the pitch, loudness, and duration (in milliseconds) of each one and compare them. If the total energies of each are vastly different, linking them would produce a disagreeable click or pop, so the link is rated as "expensive," and the computer avoids it. Some linkages are far less likely to occur than others, as in real spoken English, certain " $k_2$ " sounds are almost never followed by certain " $a_1$ " sounds. Those links could be deemed costly, too, and the computer could avoid them altogether. The word *cat* could theoretically call upon 10,000 ways of linking the " $k_2$ " and " $a_1$ " sounds. In practice, though, fewer than 100—a manageable number of choices for the computer to handle—can pass as reasonable facsimiles of human sounds.

There are many other niggling problems to deal with in Text-to-speech, such as how to teach the speaking computer to distinguish between written words like *bow* (as in "bow and arrow") and *bow* (as in the bow of a ship), or to recognize that minus signs aren't the same as hyphens [8]. This problem is complicated even further by the embedding of foreign words and names that require language specific letter-to-sounds rules for proper pronunciation.

Despite all the difficulties, speech synthesis and Text-to-speech technology has improved a lot. AT&T's Crystal, IBM's Via Voice, Dragon Naturally Speaking are some of the well known Text-to-speech software available in the commercial arena. Voices in Text-to-speech systems do not yet sound entirely natural. In short phrases ("I'd like to buy a ticket to Stockholm"), they can pass for a human, albeit an officious one. But longer