

READABLE

An Automatic Reading Tutor for Non-reading Adults

A Special Problem
Presented to
the Faculty of the Division of Physical Sciences and Mathematics
College of Arts and Sciences
University of the Philippines Visayas
Miag-ao, Iloilo

In Partial Fulfillment
of the Requirements for the Degree of
Bachelor of Science in Computer Science by

GONZALES, Benjie Jr.
PANIZALES, John Patrick
PIORQUE, Lester

Francis Dimzon
Adviser

December 24, 2022

Abstract

Advancements in technology have allowed people to utilize machine learning techniques such as automatic speech recognition (ASR). ASR is a technology that allows computers to automatically recognize and transcribe spoken language , and it has made significant advances in recent years. Studies have been conducted and results show the potential of technology-based reading tutors in improving reading skills. However, it can still be challenging to achieve high levels of accuracy for some languages and accents, especially those that are underrepresented in ASR training data. This paper, therefore, aims to develop a system to detect the acceptability of an input reading speech relative to a reference speech in Hili-gaynon. The system will use the open source speech recognition toolkit — Kaldi, to build the acoustic model and to perform the viterbi-forced alignment process to determine the input speech’s similarity against the reference speech.

Keywords: Reading miscue detector, automated reading tutor, Hili-gaynon language, viterbi-forced alignment

Contents

1	Introduction	1
1.1	Overview of the Current State of Technology	1
1.2	Problem Statement	3
1.3	Research Objectives	3
1.3.1	General Objective	3
1.3.2	Specific Objectives	3
1.4	Scope and Limitations of the Research	4
1.5	Significance of the Research	4
2	Review of Related Literature	5
2.1	Philippine Education Crisis	5
2.2	Speech Recognition and Reading Miscue Detection	6
2.3	The Language Dilemma	6
2.4	Kaldi ASR Toolkit	7
2.5	The Hiligaynon Language	8
3	Research Methodology	9
3.1	Research Activities	9

3.1.1	Acoustic Modelling	9
3.1.2	Forced Alignment	10
3.1.3	Evaluation	10
3.2	Calendar of Activities	10
4	Preliminary Results/System Prototype	12
A	Appendix	13
B	Resource Persons	14
	References	15

List of Figures

List of Tables

2.1	Table of Hiligaynon-specific phonemes used in training the system’s acoustic model (Gavieta, et al., 2022, p. 20)	8
3.1	Timetable of Activities	11

Chapter 1

Introduction

1.1 Overview of the Current State of Technology

The ability to read is a fundamental skill that is necessary for success in many areas of life. Unfortunately, there are many adults in the Philippines who have never learned to read or who have difficulty reading due to various reasons such as illiteracy, limited education, or learning disabilities. These individuals often face significant barriers to employment, education, and social participation, leading to a cycle of poverty and marginalization.

In 2019, the Philippines achieved a literacy rate of 96.5 % for the segment of the population aged 10 and over according to the PSA's Functional Literacy, Education and Mass Media Survey (FLEMMS), as reported in an article from Business World entitled, "Literacy rate estimated at 93.8% among 5 year olds or older — PSA." Literacy was defined as the ability to read and write "with understanding of simple messages in any language or dialect." However, the same article notes that this was the same rate observed in 2013, a matter described as alarming by University of Asia and the Pacific Senior Economist Cid L. Terosa, stating that even minimal improvements should be expected especially after six years.

More recently, according to a report published by UNICEF in collaboration with UNESCO and the World Bank, the percentage of 10-year-olds in low- and middle-income countries who are unable to read is as high as 70%. This figure has likely been affected by school closures brought about by the COVID-19 pandemic. The same report also stated that only 10% of children in the Philippines were able to read simple text as of March 2022. Alarming, a separate report

published by the World Bank in 2021 found that the rate of learning poverty - defined as the inability to read simple text by age 10 - in the Philippines was at 90%. These statistics highlight the urgent need to address the education crisis in the Philippines and the need to further augment the country's current literacy situation.

To address this problem, we propose the development of an automatic reading miscue detection system called Readable, specifically targeting the Hiligaynon language. Hiligaynon, also known as Ilonggo, is an Austronesian language spoken in the Western Visayas region of the Philippines, particularly in the provinces of Iloilo, Guimaras, Negros Occidental, and Capiz. It is one of the major languages of the Philippines, spoken by millions of people as a first or second language. Our reading miscue detection system will utilize machine learning techniques including automatic speech recognition (ASR). ASR is a technology that allows computers to automatically recognize and transcribe spoken language, and it has made significant advances in recent years. However, it can still be challenging to achieve high levels of accuracy for some languages and accents, especially those that are underrepresented in ASR training data. By targeting local languages like Hiligaynon and designing our ASR system to work well for these languages, we can help ensure that our reading tutor is accessible and effective for non-reading adults in the Philippines.

While there are some similar applications like Google Read Along available for reading instruction, they may not be accessible or relevant for many non-reading adults in the Philippines due to language barriers or lack of internet connectivity. By targeting local languages like Hiligaynon and utilizing the benefits of natural language processing and machine learning techniques, our automatic reading tutor can provide personalized and effective reading instruction that is accessible and relevant for non-reading adults in the Philippines. By providing accessible, effective, and scalable reading education in Hiligaynon, we hope to improve the lives and prospects of non-reading adults in the Philippines and break the cycle of poverty and illiteracy. Children with strong literacy skills grow more consistently and confidently in their studies, and reading literacy is a crucial gateway to other learning areas such as the humanities, mathematics, and the sciences. By addressing learning poverty and promoting reading literacy, we can help ensure that children in the Philippines have the opportunity to reach their full potential and succeed in their studies.

1.2 Problem Statement

Given the current educational crisis our country is facing (UNICEF, UNESCO, & Bank, 2022) and with the aim to further improve the current state of literacy rate of our country (Hernandez, 2020), the development of automatic reading tutor systems which entails building reading miscue detection systems and other related programs becomes relevant. Furthermore, the limited resources available for Hiligaynon in the context of speech processing technologies opens a good opportunity to attempt to make a contribution for the said domain of interest.

1.3 Research Objectives

1.3.1 General Objective

The aim of this project is to develop a reading miscue detection system that would determine the acceptability of an input speech relative to a reference speech pattern.

1.3.2 Specific Objectives

Specifically, the project targets to:

1. Train and model a DNN-based acoustic model for Hiligaynon.
2. Use the developed acoustic model to derive phonemic transcriptions of the input speeches/audio.
3. Determine the acceptability of a user's speech input, given a set of predetermined words to read, in terms of its deviation from reference transcriptions via forced alignment; judging is based on a set threshold score.
4. Evaluate the model in terms of the Word Error Rate via 5-fold cross validation.

1.4 Scope and Limitations of the Research

The system is specific to the Hiligaynon language. The words used in the audio data are limited to 2-3 syllable Hiligaynon words from the book “Hiligaynon Lessons” by Cecille L. Motus. Deviations are only measured word-by-word. In terms of the toolkit, the system is limited by the features offered by Kaldi - an open source speech recognition toolkit for speech recognition and signal processing.

1.5 Significance of the Research

This project aims to take one step towards developing a solution for improving the reading skill of Filipinos, particularly non reading adults. The authors also aim to contribute to the growing efforts of including Philippine local languages, specifically, Hiligaynon in literatures related to speech processing, particularly those relevant to the development of automated reading tutors.

Chapter 2

Review of Related Literature

2.1 Philippine Education Crisis

The Philippines has faced a persistent and pervasive educational crisis, as evidenced by low literacy rates and learning outcomes, particularly among disadvantaged and marginalized groups. As mentioned earlier in the introduction, UNICEF, UNESCO, and the World Bank found that only 10% of children could read simple text as of March 2022 in the Philippines. The World Bank also found that 90% of children in the Philippines cannot read simple text by age 10. These statistics highlight the urgent need to address the education crisis in the Philippines and ensure that children have access to quality reading instruction.

Pascual and Guevara (2017) conducted a study evaluating the performance of a reading miscue detector and automated reading tutor for Filipino, a language spoken in the Philippines. The study was conducted with a group of elementary school students in the Philippines, and the results showed that the reading miscue detector and automated reading tutor were effective in improving reading skills. The students who used the technology demonstrated significant improvements in reading fluency, accuracy, and comprehension, compared to a control group. These results highlight the potential of technology-based reading tutors to improve reading skills among children in the Philippines.

2.2 Speech Recognition and Reading Miscue Detection

Reading miscue detection tasks inevitably borrow concepts from the development of speech processing technology, particularly speech recognition. Automatic Reading Tutors, such as the one developed by Pascual and Guevara (2017), are machine-aided systems designed to help its users improve their skill in reading by offering help or guidance when it detects reading miscues or disfluencies in user input reading speech. Part of the approach in their system involved deriving the phone symbol sequence corresponding to the input speech, which they force aligned with a reference speech to determine deviations based on a computed likelihood score. This process can benefit from speech recognition, where various types of acoustic and/or language models, especially machine learning models are used to make sense of acoustic signals by extracting features from the said signals. These features can then be used as inputs for analysis or to whatever tasks authors deem to be appropriate. A study by Rasmussen, Tan, Lindberg, and Jensen (2009) also used an ASR component in their system for detecting miscues in dyslexic read speech.

In their study "Listen, Attend and Spell," Chan, Jaitly, Le, and Vinyals (2015) proposed a neural network architecture for automatic speech recognition (ASR) that they dubbed the "Listen, Attend and Spell" (LAS) model. The LAS model was designed to be a more efficient and accurate ASR system, particularly for languages with limited data availability. The LAS model utilizes an attention mechanism, which allows the model to focus on specific parts of the input audio, rather than processing the entire audio signal at once. This allows the model to better handle the variability and noise present in real-world audio, and it enables the model to learn more efficiently and accurately. The LAS model was tested on several datasets and outperformed existing ASR systems, demonstrating its effectiveness for automatic speech recognition tasks.

Indeed, speech recognition is one key component of computer-assisted language learning systems.

2.3 The Language Dilemma

While there is a wealth of English-oriented ASR systems, other languages, especially lesser known languages tend to struggle in these situations. For instance, most leading tech companies tend to focus on developing speech recognition tech-

nologies for the English language such as DeepSpeech, Mozilla’s speech to text engine and OpenAI’s Whisper. In the context of the Philippines, research for developing speech processing technologies for the Filipino language is by no means a desert, however developing efficient ASR systems for the said language have yet to be seen, as noted by Dimzon and Pascual (2020) . For instance the previously mentioned authors — guided by the motivation to fill in knowledge gaps in Filipino phoneme recognition — were able to develop an “Automatic Phoneme Recognizer for Children’s Filipino Read Speech”. Additionally, Aquino, Tsang, Lucas, and de Leon (2019) was able to develop a system using a grapheme to phoneme (G2P) approach, in conjunction with selected ASR models which have been found out to be just as effective as human transcribers. Other local languages, however, are challenged by limited resources. efforts are underway. Another related study was conducted by Billones and Dadios (2014) , where they created a 5-word vocabulary speech recognition system for Hiligaynon terms used as motion commands implemented for a breast self-examination (BSE) multimedia training system.

2.4 Kaldi ASR Toolkit

Povey et al. (2011) described Kaldi as a modern toolkit for speech recognition. It is designed to be extensible and has one of the least restrictive licenses making it more accessible. Several studies have incorporated Kaldi into their implementations.

For instance, Upadhyaya, Farooq, Abidi, and Varshney (2017) , developed a continuous Hindi speech recognition model using Kaldi, citing the toolkit for its ability to create high quality lattices and sufficient speed for real time recognition. It also said that the mentioned toolkit is actively maintained and accessible.

Additionally, not only is Kaldi able to support conventional models such as gaussian mixture models (GMMs) but is also able to implement deep neural network based structures. For example, Kipyatkova and Karpov (2016) developed a “DNN-Based Acoustic Modeling for Russian Speech Recognition Using Kaldi.” The paper mentioned using DNN implementations in Kaldi, ultimately choosing Dan’s implementation because of its support for parallel training on multiple CPUs.

2.5 The Hiligaynon Language

Hiligaynon, also known as Ilonggo, is an Austronesian language spoken in the Western Visayas region of the Philippines, particularly in the provinces of Iloilo, Guimaras, Negros Occidental, and Capiz. It is one of the major languages of the Philippines, spoken by millions of people as a first or second language.

Hiligaynon has a rich and varied vocabulary, with many loanwords from Spanish, English, and other languages. According to Hiligaynon Reference Grammar by Wolfenden (2019) has a complex verb conjugation and tense system, with a range of tense markers including markers for past, present, and future tense, as well as markers for perfective and imperfective aspect. The book also notes that Hiligaynon has a number of mood markers, including markers for indicative, imperative, and subjunctive mood.

Additionally, Hiligaynon Reference Grammar by Wolfenden (2019) describes the phonemic alphabet of Hiligaynon as consisting of 28 letters: A, B, C, D, E, F, G, H, I, J, K, L, M, N, Ñ, O, P, Q, R, S, T, U, V, W, X, Y, and Z. The book notes that the letters C, F, J, Q, V, X, and Z are not used as frequently in Hiligaynon as in other Philippine languages, and that the letter Ñ is used to represent the Spanish sound "ny."

Table 2.1: Table of Hiligaynon-specific phonemes used in training the system’s acoustic model (Gavieta, et al., 2022, p. 20)

Phone Class	Phones/Diphone
Bilabial stops	/p/, /b/
Dental stops	/t/, /d/
Velar stops	/k/, /g/
Africate	/j/
Fricatives	/s/, /sh/, /v/, /z/, /f/
Nasals	/m/, /n/m /ng/
Liquids	/l/, /r/
Semivowels/Glides	/w/, /y/
Vowels	/i/, /e/, /a/, /o/, /u/
Diphones	/ha/, /he/, /hi/, /ho/, /hu/, /at/, /aw/, /ay/, /oy/

Chapter 3

Research Methodology

This chapter lists and discusses the specific steps and activities that will be performed to accomplish the project.

3.1 Research Activities

3.1.1 Acoustic Modelling

For this paper, we will be utilizing the Kaldi ASR toolkit for modelling the acoustic model. Kaldi is popular when it comes to automatic speech recognition because it is flexible and accessible compared to the other ASR toolkits.

For the audio file, three speakers will utter two-to-three syllables of Hiligaynon words. These audio files will be grouped and classified as testing and training data. Each audio file has a corresponding transcription text document that will help in achieving an accurate model result. Speakers are also required to record the given words clearly and audibly. Different speakers for the testing and training audio data will be observed to yield an unbiased result.

One important key component for the Kaldi to function is its various phonemes. These phonemes were paired with words from the system's local dictionary. Table 2.1 from page 8 shows the phonemes that will be used for this study's acoustic modelling.

3.1.2 Forced Alignment

Pascual and Guevara (2017) implemented an HMM-based Viterbi-forced alignment method to produce a likelihood score that tells whether there are reading miscues to the input speech in reference to the target speech. A threshold-based classification using a threshold likelihood score was used to determine this decision. Similarly, Rasmussen et al. (2009) also implemented a forced-alignment method to detect reading miscues. This shows how time-aligned transcriptions are useful for application related to speech recognition (Dimzon & Pascual, 2020).

In this study, reading miscue detection is aimed to be achieved also by forced alignment method. The Kaldi ASR toolkit includes alignment scripts, which the study aims to use for the mentioned objective. Kaldi's alignment process outputs a sequence of alignment ids which tell what was spoken in a given frame.

This information will then be used as inputs for a logistic regression model to calculate the probability of an utterance being acceptable, in reference to a target speech or audio data. This is similar to the approach of Pascual and Guevara (2017), where their study used threshold-based classification to determine the likelihood of detecting reading miscues.

3.1.3 Evaluation

Evaluation of the system will be done by measuring the word error rate across different iterations of a 5-fold cross validation technique to maximize the items gathered for the dataset.

3.2 Calendar of Activities

Table 3.1 shows a Gantt chart of the activities. Each bullet represents approximately one week worth of activity.

Table 3.1: Timetable of Activities

Activities (2009)	Sept	Oct	Nov	Dec	Jan	Feb	Mar	Apr	May	Jun	Jul
Prerequisite knowledge re- search	••	••									
Identification of potential proposal and features		••	•••								
Writing of proposal paper			•	••••							
Recording of audio files						••					
Modeling						•	•••	•			
Development of the system							•	•••	••	•	
Analyzing and interpreta- tion of the results									•••	•	
Documentation	••	••••	••••	••••		•••	••••	•••	••••	•	•

Chapter 4

Preliminary Results/System Prototype

This chapter presents the preliminary results or the system prototype of your SP. Include screenshots, tables, or graphs and provide the discussion of results.

Appendix A

Appendix

Appendix B

Resource Persons

Dr. Firstname1 Lastname1

Adviser

Affiliation1

emailaddr@domain.com

Mr. Firstname2 Lastname2

Role2

Affiliation2

emailaddr2@domain.com

Ms. Firstname3 Lastname3

Role3

Affiliation3

emailaddr3@domain.net

References

- Aquino, A., Tsang, J. L., Lucas, C. R., & de Leon, F. (2019, 8). G2P and ASR techniques for low-resource phonetic transcription of Tagalog, Cebuano, and Hiligaynon. *2019 International Symposium on Multimedia and Communication Technology (ISMAC)*. Retrieved from <http://dx.doi.org/10.1109/ismac.2019.8836168> doi: 10.1109/ismac.2019.8836168
- Billones, R. K. C., & Dadios, E. P. (2014, 11). Hiligaynon language 5-word vocabulary speech recognition using Mel frequency cepstrum coefficients and genetic algorithm. *2014 International Conference on Humanoid, Nanotechnology, Information Technology, Communication and Control, Environment and Management (HNICEM)*. Retrieved from <http://dx.doi.org/10.1109/hnicem.2014.7016247> doi: 10.1109/hnicem.2014.7016247
- Chan, W., Jaitly, N., Le, Q. V., & Vinyals, O. (2015, 8). Listen, Attend and Spell. *arXiv: Computation and Language*.
- Dimzon, F. D., & Pascual, R. M. (2020, 12). An Automatic Phoneme Recognizer for Children's Filipino Read Speech. *2020 IEEE International Conference on Teaching, Assessment, and Learning for Engineering (TALE)*. Retrieved from <http://dx.doi.org/10.1109/tale48869.2020.9368399> doi: 10.1109/tale48869.2020.9368399
- Hernandez, J. (2020, 10). *Literacy rate estimated at 93.8PSA*. Retrieved from <https://www.bworldonline.com/economy/2020/10/29/325932/literacy-rate-estimated-at-93-8-among-5-year-olds-or-older-psa/>
- Kipyatkova, I., & Karpov, A. (2016). DNN-Based Acoustic Modeling for Russian Speech Recognition Using Kaldi. *Speech and Computer*, 246–253. Retrieved from http://dx.doi.org/10.1007/978-3-319-43958-7_29 doi: 10.1007/978-3-319-43958-7_{_}29
- Pascual, R., & Guevara, R. (2017). Experiments and Pilot Study Evaluating the Performance of Reading Miscue Detector and Automated Reading Tutor for Filipino: A Children's Speech Technology for Improving Literacy. *Science Diliman*, 29(1), 5–36. Retrieved from <https://journals.upd.edu.ph/index.php/sciencediliman/article/view/5622>

- Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N. K., ... Vesely, K. (2011, 1). The Kaldi Speech Recognition Toolkit. *IEEE Automatic Speech Recognition and Understanding Workshop*. Retrieved from https://publications.idiap.ch/downloads/papers/2012/Povey_ASRU2011.2011.pdf
- Rasmussen, M. H., Tan, Z.-H., Lindberg, B., & Jensen, S. H. (2009, 9). A system for detecting miscues in dyslexic read speech. *Interspeech 2009*. Retrieved from <http://dx.doi.org/10.21437/interspeech.2009-448> doi: 10.21437/interspeech.2009-448
- UNICEF, UNESCO, & Bank, W. (2022, 3). *Where are we on Education Recovery?* (Tech. Rep.). Retrieved from <https://www.unicef.org/reports/where-are-we-education-recovery>
- Upadhyaya, P., Farooq, O., Abidi, M. R., & Varshney, Y. V. (2017, 3). Continuous hindi speech recognition model based on Kaldi ASR toolkit. *2017 International Conference on Wireless Communications, Signal Processing and Networking (WiSPNET)*. Retrieved from <http://dx.doi.org/10.1109/wispnet.2017.8299868> doi: 10.1109/wispnet.2017.8299868
- Wolfenden, E. (2019). *Hiligaynon Reference Grammar* (Open Access ed.). University of Hawaii Press. Retrieved from <https://core.ac.uk/display/211329359>