# SPEECH RECOGNITION SYSTEM FOR TAMIL LANGUAGE USING CMUSPHINX

A.R.F.S. Fanoon[134] & G.A.I. Uwanthika[135]

Correspondence: fanoonarfs@gmail.com

## ABSTRACT

Speech recognition is the process of converting an audio signal into sequences of words paving the way for enabling a rich Human-Computer Interaction in many emerging applications. It is a very difficult and a complex task to recognize speech using a computer though several approaches have been made to develop an accurate speech recognition system. Most of the researches have been done for English, Chinese, Mandarin and Arabic languages while a very few have been done for Tamil language. Therefore, this research intends to develop a continuous speech recognition system for Tamil language. Tamil language has 12 vowels and 18 consonants. This paper is focused on building a Tamil Speech Recognition system using CMUSphinx toolkit. The system was developed for several desktop applications with the aim that in the future, this system could be integrated with those applications and executes those applications through Tamil speech input rather than typing and clickin. Audio recordings from several people were recorded in different environment which were later converted to Wav file format. A text corpus, transcription file and fileids file for both training and testing database were prepared accordingly. Three models namely Acoustic model, Language model and Lexicon model have been developed. Hidden Markov Model was employed for building Acoustic Model. The performance of the system over various speaker subsets of different sex, age and dialect was examined. The Word Error Rate was calculated in order to measure the performance of the system while the accuracy of the system was calculated using another formula. The results from both the measurements showed a satisfactory performance.

*Keywords:* human-computer interaction, hidden markov model, acoustic model, word error rate.

## INTRODUCTION

A speech recognition system recognizes the words uttered by a speaker. Speech is one of the best methods to interact with computers rather than typing and clicking. It also acts as a better interface for illiterate as well as physically disabled people. Researchers and computer scientists are working towards making computers understand human speech. As a result of these, attempts are taken to develop a human and machine interface where both can communicate in an unskilful way.

Several researches have been carried out on speech recognition, but most of them have been done only for English, Chinese, Mandarin and Arabic languages. Very few researches have been carried out for Tamil language to recognize isolated speech. Tamil is the second major language spoken in Sri Lanka and therefore, recognition of Tamil speech is beneficiary. Recognition of Tamil language is a challenging task (Chong,et.al, 2008) due to the speech variability in any human's spoken utterance and language nature of Tamil Language (Kalith, et.al, 2016). This research is a first step for implementing an efficient speech recognition

---

[134] Postgraduate Institute of Science, University of Peradeniya, Sri Lanka.
[135] Postgraduate Institute of Science, University of Peradeniya, Sri Lanka.

system for Tamil speaking people in order to increase the human computer interaction in Tamil language.

Tamil is a Dravidian language spoken by 77 million people all over the world. It is a 15th largest language among the world languages. Spoken predominantly in south India Sri Lanka and Malaysia, it is an official language in India and Sri Lanka.

The research is mainly aimed to design and implement a Tamil speech recognition system based on Hidden Markov Model (HMM) using CMUSphinx toolkit. The system is developed for the commands which are used to execute computer applications.

The main activities involved are to feed the spoken utterance into the computer via an audio device as to train the system. Then the system is tested with spoken utterance to evaluate the recognition result. The recognition is a connected word of Tamil speech of about 125 computer commands.

The organization of the paper is given as follows. Some of the work related to speech recognition and methods used previously are discussed in Section 2 under the topic literary study. Section 3 deals with the Methodology of the proposed system. Speech recognition architecture, its design, Technology adopted, and approaches used are described in the 4th section. Section 5 describes about the Evaluation details. Finally, the paper concludes and proposes future work.

## LITERARY STUDY

### Current State of Speech Recognition Technology

Speech recognition technology is keeping on advancing its technology daily, taking many factors into consideration. Speech input may either be isolated speech, connected speech, spontaneous speech or continuous speech (Kalith, et.al, 2016) and the system may handle the speaker utterances accordingly.

Few systems have been developed and they allow only a limited form of natural language input within a very specific domain at any particular point in the interaction. ASR devices have been available from 1970s, but they were very expensive and could recognize only isolated speech. Later in 80s and 90s, ASR technology showed some improvements and reached to an advanced level in the late 90s where the researchers were able to develop software for desktop applications. They also developed two types of ASR: Speaker Dependent and Speaker Independent.

ASR also used various techniques for recognition of speech, and they are listed below.

Template based approaches matching
Words are recorded and saved in a repository. Unknown words uttered by a speaker are compared against these recorded words as to find the best matched word. Dynamic Type warping approach falls into this technique (Tolba & O'Shaughnessy, 2001).
Knowledge based approaches

This approach stores speech an expert knowledge on the variations of speech but this is not much effective.

Statistical based approaches (Hidden Markov Model-HMM)

The speech variations are modelled statistically and based on the probabilistic value; the words are matched.

Learning based approaches

This approach utilizes the neural network and genetic algorithm techniques. The models are studied through emulation or evolutionary process.

Artificial intelligence approach

This is an attempt to mechanize the way the words are recognized according to the user's perspective. This approach widely uses expert systems (Mori et.al., 1987).

Among all the ASR techniques mentioned above, we have adopted statistical based approach which is HMM. HMM is an effective method for developing a high-performance speaker independent speech recognition engine. This technique also is very suitable for large vocabulary (Ehsani & Knodt, 1998).

## Tamil Speech Recognition System

Tamil language is spoken in many parts of the world and the majority is in the South Asia: India and Sri Lanka. A little number of researches has been carried out in the context of Tamil speech recognition. But they have only been done to recognize numbers and some isolated and continuous words or utterances. No systems have been developed for desktop applications for smart devices. Tamil language has 12 vowels and 18 consonants, but the number of phonemes counts more due to different sounds depending on where they occur (Gnanathesigar, 2012).



*Fig. 1 Vowels and Consonants in Tamil Language*

Implementation of ASR for Tamil language is not an easy task as speech recognition design depends on many other fields, namely acoustics, signal processing, pattern recognition, phonetics, linguistics, psychology, neuroscience, and computer science.

**METHODOLOGY**

Figure 1 below shows the process to design and implement the speech recognition system. The main aim of the research is to implement a speech recognition system for Tamil language for windows desktop applications. CMUSphinx toolkit was used for recognition process. CMUSphinx is an open source speech recognition toolkit known that comprises of libraries Sphinxbase, Pocketsphinx, and SphinxTrain.

Pocketsphinx - lightweight recognizer library written in C.
Sphinxbase - support library required by Pocketsphinx
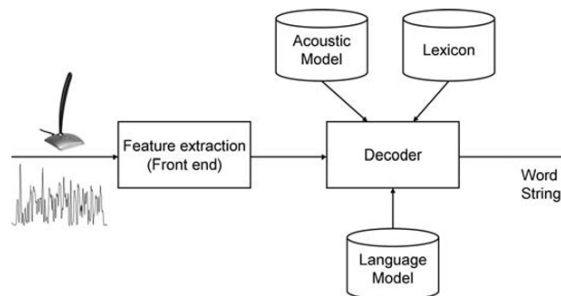SphinxTrain - acoustic model training tool



*Fig. 2 Speech Recognition Process*

**Data Preparation**

Voice recordings of commands used for executing PC applications were first collected from several people speaking Tamil language living in various parts of Sri Lanka were collected. During the collection of voice recordings, age, gender, and language dialect were mainly considered. Additionally, the environmental background and time were also taken into considerations. When all the recordings were taken, as a next step, the audio files were converted into wav files of mono channel at a sampling rate of 16kHz and a bi rate of 16. Later, the voice recordings were segmented into separate voice commands using Audacity. Figure 2 shows a segmented voice command.
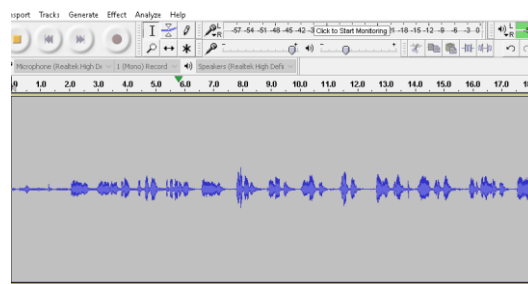


*Fig. 3 Segmentation of Voice Recording using Audacity*

## Language Model

Language model represents the most likelihood occurrence of a word uttered. It works based on a probability and used to search for a word and improve accuracy. It lists 1-,2- and 3-grams along with their likelihood (the first field) and a back-off factor (the third field). SRILM toolkit was used to build the language model for Tamil language.

## Lexicon Model

Lexicon Model also known as Pronunciation dictionary, describes how each word are pronounced. Each and every word the recognizer needs to recognize must be placed in the dictionary file. That is a dictionary file contains the words to be recognized and their respective monophonic sequences. Here, pronunciation dictionary for Tamil language has been constructed manually for each word in the vocabulary. 'Anunaadam' tool, an online IPA converter was used to create the monophonic sequence of the words.
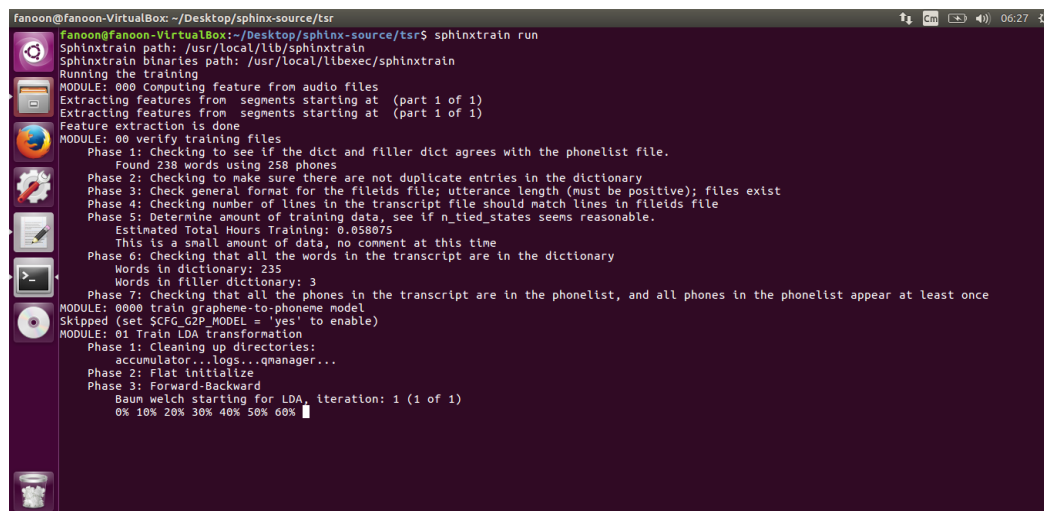
| | |
|---|---|
| மின்னஞ்சல் | mɪ n̪ n̪ʌ ɲd ʒʌ l |
| மின்னஞ்சலை | mɪ n̪ n̪ʌ ɲd ʒʌ lʌː |
| திற | t̪ɪ rə |
| உற்புகு | ʷʊ r pʊ xɨ |
| புகுபதிவு | pʊ xu βʌ ðɪ vʉ |
| வெளியேறு | ʋɛ˯ lˈɪɪ̯ eː rɨ |
| விடுபதி | ʋɪ˯ ɾɨ βʌ ðɪ˙ |
| கீழே | kiˈː ʒeˑ |
| செல் | sɛ˯ l |
| மேலே | meː leˑ |
| நிறுத்து | n̪ɪ rɨ t̪ t̪ɨ |
| செய்தி | sɛ˯ ɪ̯ ðɪ˙ |

## Acoustic Model

An acoustic model was trained for Tamil language for the training as well as testing purpose. In order to train acoustic model, fields file, transcription file, dictionary file, and file containing phone sounds (phone file) were prepared. Sphinxtrain tool was used to obtain the acoustic model.

| | | | | |
|---|---|---|---|---|
| | | | | SIL |
| 1 | speaker1/seg1_f1_c1 | \<s\> கணினி \</s\> (seg1_f1_c1) | | -si: |
| 2 | speaker1/seg1_f1_c2 | \<s\> மின்னஞ்சலை திற \</s\> (seg1_f1_c2) | | -sl |
| 3 | speaker1/seg1_f1_c3 | \<s\> உற்புகு \</s\> (seg1_f1_c3) | | Fo |
| 4 | speaker1/seg1_f1_c4 | \<s\> புகுபதிவு \</s\> (seg1_f1_c4) | | bɑˑː |
| 5 | speaker1/seg1_f1_c5 | | | |
| 6 | speaker1/seg1_f1_c6 | \<s\> வெளியேறு \</s\> (seg1_f1_c5) | | bʉ˯ |
| 7 | speaker1/seg1_f1_c7 | \<s\> விடுபதி \</s\> (seg1_f1_c6) | | c |
| 8 | speaker1/seg1_f1_c8 | \<s\> கீழே செல் \</s\> (seg1_f1_c7) | | cə |
| 9 | speaker1/seg1_f1_c9 | \<s\> மேலே செல் \</s\> (seg1_f1_c8) | | cɨ |
| 10 | speaker1/seg1_f1_c10 | \<s\> நிறுத்து \</s\> (seg1_f1_c9) | | d |
| 11 | speaker1/seg1_f1_c11 | ... செய்தியைத்திற ... | | |

*Fig. 5 fileid file, transcription file and phone file*

*Fig. 6 Acoustic Model Training in ubuntu*

**Speech Recognition on Windows**

All the above processes were carried on in Ubuntu environment and then, the output was implied onto the windows environment. For this purpose, CMUSphinx was installed onto windows with other dependencies such as MS Visual Studio, Swig, Python and many other. After the successful setting up of the environment on windows, two additional files argFile.txt and ctlFile.txt were created inside the Pocketsphinx.

argFile.txt consisted of the path to the acoustic model, language model and the lexicon model.
ctlFile.txt consists of all the utterances for the speech recognition; i.e., it has all the commands recorded.
And finally, the models were integrated with another Speech Recognition System Simon that could control the Windows applications.

**RESULTS AND DISCUSSION**

Hidden Markov Model was used for recognition purpose and both known and unknown speakers were selected. They were asked to utter 15 commands from the recorded list and the number of commands recognized was identified. The result is tabulated below.

Table 1 Results of Speech Recognition

| Speakers | Number of commands | Correct recognition | Incorrect recognition |
|---|---|---|---|
| **Known speaker** | 15 | 12 | 3 |
| **Unknown speaker** | 15 | 8 | 7 |

The Word Error Rate (WER) of the system was calculated with the formulae;

WER = (I+D+S)/N

Where I, D, and S are the number of words inserted, deleted and substituted respectively. This system had nearly three hundred words in total and only fifty of them were used for testing.

Table 2

The number of speakers against the accuracy

| Speaker | Accuracy | Speaker | Accuracy |
|---------|----------|---------|----------|
| **First** | 60% | Sixth | 78% |
| **Second** | 70% | Seventh | 80% |
| **Third** | 65% | Eighth | 85% |
| **Fourth** | 73% | Ninth | 90% |
| **Fifth** | 75% | tenth | 75% |

And the accuracy of the Tamil speech recognition system was calculated with the formula;

Accuracy = (N-D-S)/N

According to table, Table 02, it is very clear that accuracy depends on the training time and also may vary depending on the speaker and the environment where the words are uttered.

**CONCLUSION**

Speech recognition system was built for a very few commands and tested with each 5 male and female speakers. The database size used for this research work is approximately 250 words and produces an average of 85% of accuracy. Our vocabulary size is small and therefore, the system gives lesser word error rate compared to the system with large size vocabulary. Speech recognition is generally easy when vocabularies are small, but word error rate increases as the vocabulary size grows.

In future, this project can be used at a very large scale and be integrated with desktop applications. The speaker amount can be increased up to 200 including both male and female speakers with 5 hours of recording in order to get a speech recognizer with maximum accuracy. Also, the research can be further extended to increase the accuracy by parallelizing the speech recognition tasks, i.e. the search algorithm, using GPGPU processors.

# REFERENCES

Chong, J. et al., 2008. *Data-Parallel Large Vocabulary Continuous Speech,* Berkeley: s.n

CMUSphinx Tutorial For Developers'**,** wiki article, [Online]. Available: http://cmusphinx.sourceforge.net/wiki/tutorial, [Accessed: 24[th] of May 2016]

Ehsani, F, & Knodt, E. 1998, 'Speech Technology In Computer Aided Learning: Strengths And Limitations Of A New Call Paradigm', Langugae Learning and Technology, [Online]. Available: http://hstrik.ruhosting.nl/wordpress/wp-content/uploads/2013/03/Ehsani-Knodt-LLT-1998.pdf

Gnanathesigar.H, 2012, 'Tamil Speech Recognition using Semi-Continuous Models', *International Journal of Scientific and Research Publications,* vol. 2[6]

Installing CMU-Sphinx on Ubuntu, [Online]. Available:

http://jrmeyer.github.io/installation/2016/01/09/Installing-CMU-SPhinx-on-Ubuntu.html  Josh Meyer's Website, [Accessed: 15[th] of August 2016]

Kalith, M., David, A. & Samantha, T. (2016), 'Isolated to Connected Tamil Digit Speech Recognition System Based on Hidden Markov Model', *International Journal of New Technologies in Science and Engineering*, 3(4), 2-6

Lee, K. F., n.d. In: *Automatic Speech Recognition.* Boston/London/Dodrecht: Kluwer Academic Publishers

Mori R.D, Lam L., and Gilloux M. (1987). Learning and plan refinement in a knowledgebased system for automatic speech recognition. IEEE Transaction on Pattern Analysis Machine Intelligence, 9(2):289-305

Picheny, M., (2002). Large vocabulary speech recognition, IEEE Computer, 35(4),42- 50

Reddy D.R., (1976). Speech Recognition by Machine: a Review. Proceeding of IEEE, 64(4),501-531

Shammur, A.C 2010, 'Implementation of Speech Recognition System for Bangla: thesis, BRAC University

SRILM Installation and Running Tutorial, [Online]. Available: http://okapiframework.org/wiki/index.php?title=SRILM_Installation_and_Running_Tutorial [Accessed: 20[th] of August 2016]

Tamil, Omniglot, the online encyclopedia for writing systems and languages, [Online]. Available: http://www.omniglot.com/writing/tamil.htm , [Accessed: 15[th] of November 2016]

Tolba, H., and O'Shaughnessy, D., (2001)., Speech Recognition by Intelligent Machines, IEEE Canadian Review (38).

Vrinda, & Chander,S, T 2003, 'SPEECH RECOGNITION SYSTEM FOR ENGLISH LANGUAGE', *International Journal of Advanced Research in Computer and Communication Engineering,* vol. 2, pp. 919-922.

Waibel, A. & Kai-Fu Lee, n. d. In: M. B. Morgan, ed. *Readings in Speech Recognition*. San Mateo, California: Morgan Kaufmann Publishers Inc...

Y. Anunaadam, [Online]. Available from:

http://anunaadam.appspot.com/, [Accessed: 25th of May 2016].