

# Enhance Text Creature through Video

Mr. Prashant Mhetre, Mr. Kedar Devkar, Mr. Amol Bhagat, Mr. Pramod Patil,  
Mr. K.B. Vayadande

**Abstract**— The large quantity of education video out there on web, a rise the usability of video data are growing rapidly. Video transcription within which should be conversion video lecture into text information. This can be manner of produce document or notes through the video. This paper present an ASR technique supported Hidden Markov Model. First of all, extract audio from video and transforms speech wave form into multiple frame used by recognition, applying Automatic Speech Recognition on audio track and extract raw data from audio. Then analysis of data in order to get the phonetic dictionary, the pronunciation of every word must be represent phonetically. And represent text document as output of video file

**Index Terms**— Lecture videos, HMM, Large Vocabulary.

## I. INTRODUCTION

Now a day's video documentary is that the most significant issue in education sector to be told new things through the video. During this globe several peoples are grasping there information through videos from you tube, nptel etc. many education institutes record knowledgeable video lecture and upload on internet site. Some institute provide on-line teaching to student in fundamental quantity of institute and student will get information through the video documentary this is often straightforward way to understanding conception of related subjects. The business sector also are using on-line coaching to the employee.

Large number of video obtainable on the World Wide Web and student are learn the conception of related subject however hard copy of document isn't obtainable. It's tedious job to convert video documentary to text document. Main goal of our project is give softcopies of video documentary. It's one quite personal assistant that is we are able to simply carry anyplace with US and done this sort of tedious jobs. This personal assistant is additionally useful for deaf student  
Essentially we are applying approach for extract text using Automatic Speech Recognition (ASR) for the long audio.

**Manuscript received March 17, 2015**

**Prashant Tryambak Mhetre**, Department of Information Technology, Bharati Vidyapeeth's college of Engineering, Kolhapur, India

**Kedar Jyotiba Devkar**, Department of Information Technology, Bharati Vidyapeeth's college of Engineering, Kolhapur, India.

**Amol Dattatraya Bhagat**, Department of Information Technology, Bharati Vidyapeeth's college of Engineering, Kolhapur, India.

**Pramod Shivaji Patil**, Department of Information Technology, Bharati Vidyapeeth's college of Engineering, Kolhapur, India..

**Kuldeep B. Vayadande**, Department of Information Technology, Bharati Vidyapeeth's college of Engineering, Kolhapur, India.

This approach has work on acquire high quality input file means that speech is clean and well-structured [3] for ASR engine. Sphinx-4 could be a simply manageable, standard and interchangeable framework to assist develop new idea within the core analysis of Hidden Markov Model (HMM) recognition systems [2]. Sphinx-4 framework and therefore the implementations are open source. Therefore, Sphinx-4 libraries are using in our application to develop own system. In video to text analysis we have a tendency to apply the open supply Sphinx is ASR tool. To make the acoustic and language model, we have to collected speech coaching data from open-source corpora and lecture videos [1].

## II. EXTRACTION OF TEXT FROM MULTIMEDIA DOCUMENTARY

A large quantity of text information are going to be generated by using ASR technique, that extract content of lecture videos. We are extract information from audio resources of lecture videos automatically by applying applicable analysis technique must be a HMM-based speech recognizer. HMM stands for Hidden Markov Models, That is a form of statistical model. In HMM-based speech recognizers, every unit of sound is represented by a statistical model that represents the distribution of all the data for that phoneme [2].

In general manner of speech recognize is that the following, the frontend parameters depending on the frequency usually collect the audio data within the sphinx. Frequency wave form is split on utterance and that begin when the every silences. It should take all potential mixtures of word and take a look at to match them with the utterance. The best score is chosen in matching combination. The options that Sphinx produce are known as cepstrum. Cepstrum are sometimes 13-dimensional. The features that Sphinx-4 produce are over cepstrum. It's 39-dimensional, and consists of the cepstrum, the delta of the cepstrum, and therefore the double delta of the cepstrum. Since, these feature would take an extended time, compare with acoustic model.

There are connected process step that are follows for big quantity text extraction using long audio that is predicated on the Hidden Markov Model.

Table 1: The following table compares the performance of Sphinx 3.3 with Sphinx-4.

Test	S3.3 WER	S4 WER	S3.3 RT	S4 RT(1)	S4 RT (2)	Vocabulary Size	Language Model
TI46	1.217	0.168	0.14	.03	.02	11	isolated digits recognition
TIDIGITS	0.661	0.549	0.16	0.07	0.05	11	continuous digits
AN4	1.300	1.192	0.38	0.25	0.20	79	trigram
RM1	2.746	2.88	0.50	0.50	0.41	1,000	trigram
WSJ5K	7.323	6.97	1.36	1.22	0.96	5,000	trigram
HUB4	18.845	18.756	3.06	~4.4	3.95	60,000	trigram

**Key:**

- **WER** - Word error rate (%) (lower is better)
- **RT** - Real Time - Ratio of processing time to audio time - (lower is better)
- **S3.3 RT** - Results for a single or dual CPU configuration
- **S4 RT(1)** - Results on a single-CPU configuration
- **S4 RT(2)** - Results for a dual-CPU configuration

**A. Audio Classification and Segmentation**

The ASR are essentially counting on the Hidden Markov Model. Initial of all take quality input audio content for classification and segmentation, within which Associate in audio file is metameric in keeping with audio kind. Audio classification is essentially counting on utterances i.e. utterances are composed of quiet and non-quiet sounds. We tend to differentiate quiet and non-quiet sound within the audio and it's to pick non-quiet sound save into a files is termed as segmentation. In segmentation part we are using MFCC that is signify Mel-Frequency Cepstral Coefficients. Basic approach of MFCC is assessed into silence and speech frames therefore so as to get original utterances from audio track and Classification and segmentation necessary for extraction of text from the audio.

**B. Analysis phase**

Segmentation part is important for speech recognition. Analysis model have 3 part acoustic model, dictionary and language model. These part use in sphinx-4, the search manager is main part decoder and also the task of searching through the graph is completed by the search manager. Acoustic model used for the speech signals are first remodeled into a sequence of vectors that represent bound characteristics of the signal, and also the parameters of the acoustic model are mapping between unit of speech. This method is termed coaching the acoustic models. Incoming speech are divided into piece frames and these are scored against the acoustic model. The score got the actual set of frames belongs to the speech sound of the corresponding acoustic model.

In HMM based mostly recognizers are the graph represents all potential sequences of phonemes within the entire language of the task into account. Construction the

on top of graph needs for data from dictionary and for acoustic model that maps the word to the phonemes. Great amount of word vocabulary needed as a result of we tend to use the Sphinx-4 could be a terribly versatile system capable of activity many alternative varieties of recognition tasks. In table one show comparison of the vocabulary in sphinx with the word error rate is low in sphinx four as compare to sphinx three further as increase size of vocabulary of sphinx four.

The search graph conjointly give data concerning however can occur bound words. This sort of data provided by language model. This model that provides chance between the entry node and first node of HMM. The trail can consequently have the next score.

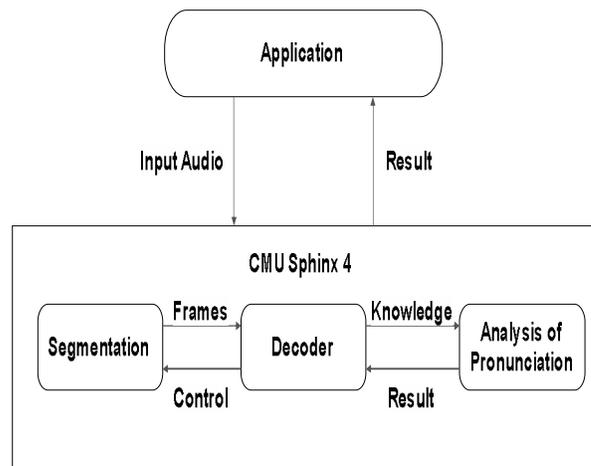


Fig 1: Enhance the text extraction using CMU Sphinx 4

**C. Text Recognition**

The CMU Sphinx design utilized in our application is

show in figure one and improve the projected work. Within the decoder is employed for word recognition in CMU Sphinx tool within which basic input of decoder is segmentation of frame and together with the search graph from the linguist to get result. Decoder maintain the every state for generating consequence because the output given input audio files. There for decoder main produce result object for displaying consequence of ASR engine pattern search manager performs a straightforward breadth-first search through the search graph throughout the coding method to search out the simplest path.

### III. CONCLUSION

This paper represent AN extract text from long audio file that why we tend to needed classification and segmentation of audio file pattern Mel-Frequency Cepstral Coefficients to differentiate frames are classified into silence and speech frames. The text recognition is essentially looking on the ASR tool that base of the HMM based mostly decoder. Using this approach we tend to produce text document for long video lecture.

In future we are able to offer facility for recognizing text in varied languages.

### REFERENCES

- [1] Haojin Yang, Christoph Meinel "Content Based Lecture Video Retrieval Using Speech and Video Text Information," IEEE Transactions on Learning Technologies, Vol. 7, No. 2, April-June 2014.
- [2] Willie Walker, Paul Lamere, Philip Kwok, Bhiksha Raj, Rita Singh, Evandro Gouvea, Peter Wolf, Joe Woelfel, "Sphinx-4: A Flexible Open Source Framework for Speech Recognition", [www.cmusphinx.sourceforge.net/sphinx4](http://www.cmusphinx.sourceforge.net/sphinx4).
- [3] Jigar Patel, Kailash Singh Maurya, Sameer Kulkarni "Multimedia Keyword Spotting (MKWS) Using Training and Template Based Techniques" IJETAE, Volume 4, Issue 2, February 2014.
- [4] R. Ordelman, F. de Jong, and M. Larson, "Enhanced Multimedia Content Access and Exploitation Using Semantic Speech Retrieval," in Proc. ICSC '09, 2009, pp. 521-528.
- [5] A. Park and J. R. Glass, "Unsupervised Word Acquisition from Speech using Pattern Discovery," in Proc. ICASSP, 2006, pp. 1409 - 1412.
- [6] J. G. Wilpon, L. R. Rabiner, C. H. Lee, and E. R. Goldman, "Automatic recognition of keywords in unconstrained speech using hidden Markov models," IEEE Trans. Acoustics, Speech and Signal Processing, vol. 38, pp. 1870-1878, 1990.
- [7] Li S. Z. "Content-based classification and retrieval of audio using the nearest feature line method [J]," IEEE Transactions on Speech and Audio Processing. 2000, vol (8). pp: 619-625.
- [8] A. Jansen and P. Niyogi, "Point Process Models for Spotting Keywords in Continuous Speech," IEEE Trans. Audio, Speech, and Language Processing, vol. 17, pp. 1457-1470, 2009.
- [9] S. J. Young, N. H. Russell, and J. H. S. Russell, "Token passing: A simple conceptual model for connected speech recognition systems," Cambridge University Engineering Dept, UK, Tech. Rep. CUED/F-INFENG/TR38, 1989.