

## Abstract

Speech processing is still now a very big challenge. Speech processing for Indian language is a broad area of research. In case of Assamese language till now the research work is going very slowly. To collect our Assamese speech corpus we have collected about nine hours of Assamese speech data collected from 27 Assamese native speakers of different parts of the state with different categories. A balanced speech corpus is the basic need for any speech processing task. In this report we also describe on development of Assamese speech corpus. We mainly focused on some issues and challenges faced during development of the corpus. We try to define some syllabification rules of Assamese language. Being a less computationally aware language, this is the first effort to develop speech corpus for Assamese. As corpus development is an ongoing process, in this report we have presented only the initial task. We have split these speech files with less than ten seconds in each file. Then we transcribed these recorded speech files using IPA (**International Phonetic Alphabet**) and ASCII for automatic transcription. During IPA transcription we used thirty four phonemes, where twenty five consonants and nine vowels are used. We reported the frequency of IPA symbols. For automatic transcription we use about three hours of speech data which was transcribed using ASCII (**American Standard Code for Information Interchange**) characters. For automatic transcription we use thirty eight phonemes, where nine vowels and twenty five consonants. During this work we trained about two and half hours speech data and approximately half hour test data for automatic transcription. From our result we found that 55.26 percentage accuracy for automatic transcription of Assamese speech. In our work we also reported about individual accuracy about all phonemes used for transcription. At last we found that if we trained more data then accuracy will more and at a stage accuracy is static .

**Keywords:** *speech, IPA, ASCII, phonetic, corpus, syllabification, automatic, transcription, Assamese, phonemes, frequency*