

## Development of TTS Engine for Indian Accent using Modified HMM Algorithm

S. S. Gantayat<sup>#</sup>

<sup>#</sup> Department of CSE, GMR Institute of Technology, Rajam, Andhra Pradesh, India  
E-mail: sasankosekhar.g@gmr.it.org

**Abstract**— A text-to-speech (TTS) system converts normal language text into speech. An intelligent text-to-speech program allows people with visual impairments or reading disabilities, to listen to written works on a home computer. Many computer operating systems and day to day software applications like Adobe Reader have included text-to-speech systems. This paper is presented to show that how HMM can be used as a tool to convert text to speech.

**Keywords**— K-means, Text-to-speech, Speech synthesis, HMM Algorithm.

### I. INTRODUCTION

A text to speech system is composed of two parts, a front-end and a back-end. The front end has two major tasks. First, it converts the raw text containing symbols like numbers and abbreviations into the equivalent of written-out words. This process is often called text normalization pre-processing or tokenization. The front-end then assigns phonetic transcriptions to each word and divides and marks the text into prosodic units, like phrases, clauses and sentences. The process of assigning phonetic transcriptions to words is called text-to-phoneme or grapheme-to-phoneme conversion. Phonetic transcriptions and prosody information together make up the symbolic linguistic representation that is output by the front-end. The back-end often referred to as synthesizer, it converts the symbolic linguistic representation into sound.

HMM-based synthesis is a synthesis method based on Hidden Markov Models, also called Statistical Parametric Synthesis, which is a widely used model for any Speech Recognition and Synthesis applications. In this system, the frequency spectrum (vocal tract), fundamental frequency (vocal source), and duration (prosody) of speech are modeled simultaneously by HMMs. In the proposed project work it is considered to improve the HMM technique to the Indian accents.

Statistical parametric speech synthesis based on hidden Markov models (HMMs) is now well-established and can generate natural-sounding synthetic speech. In this framework, we have pioneered the development of the HMM Speech Synthesis System, HTS (H Triple S). This research started by developing algorithms for generating a smooth parameter trajectory from HMMs. Next, to simultaneously

model the excitation parameters of speech as well as the spectral parameters, the multispace probability distribution (MSD) HMM was developed. Using the logarithm of the fundamental frequency and its dynamic and acceleration features as the excitation parameters, the MSD-HMM enabled us to treat the sequence, which is a mixture of one-dimensional real numbers for voiced regions and symbol strings for unvoiced regions, in a probabilistic framework. To simultaneously model the duration parameters for the spectral and excitation components of the model, the MSD hidden semi-Markov model (MSD-HSMM) was developed. The HSMM is an HMM having explicit state duration distributions instead of transition probabilities, to directly model duration; it can generate more appropriate temporal structures for speech. These basic systems employed a mel-cepstral vocoder with simple pulse or noise excitation, resulting in synthetic speech with a “buzzy” quality. To reduce buzziness, mixed or multi-band excitation techniques have been integrated into the basic systems to replace the simple pulse or noise excitation and have been evaluated. These basic systems also had another significant problem: the trajectories generated from the HMMs were excessively smooth due to statistical processing, resulting in synthetic speech with a “muffled” quality. To alleviate this problem, a parameter generation algorithm that considers the global variance (GV) of a trajectory to be generated was developed.

### II. MODULES

Module 1: Creation of Speech Corpus (dictionary).

Module 2: Labeling.

Module 3: Prototype HMM (Training of HMMs).

Module 4: Speech Waveform Synthesizer (Model Set).

## A. Creation of Speech Corpus

CMU Arctic databases designed for the purpose of speech synthesis research. These single speaker speech databases have been carefully recorded under studio conditions and consist of nearly 1150 phonetically balanced English utterances. They are distributed as free software, without restriction on commercial or non-commercial use. The Arctic corpus consists of four primary sets of recordings (3 male, 1 female), plus several ancillary databases. Each database is distributed with automatically segmented phonetic labels. These extra files were derived using the standard voice building scripts of the Festvox system. In addition to phonetic labels, the databases provide complete support for the Festival Speech Synthesis System, including pre-built voices that may be used as is. Festival and Festvox are available at <http://www.festvox.org>. The Arctic speech corpus is available at [http://www.festvox.org/cmu\\_arctic](http://www.festvox.org/cmu_arctic).

CMU Arctic is a set of single speaker databases that have been carefully recorded under studio conditions, packaged with associated information such as phonetic labels and pitchmark files. An Arctic “database” is a reading of the Arctic prompt set (plus associated files) by a single speaker in a specified style of delivery. This release of Arctic contains recordings by four separate speakers. When referring to the Arctic “corpus” we mean the entire collection of databases, including test sets. The databases have version numbers. As with computer code, version numbers indicate the level of maturity and stability. Numbers with a zero after the decimal point (e.g. version 1.0) are major releases intended to serve as a reference point for system development and evaluation. Minor releases are subject to change, allowing for more frequent additions, deletions, and improvements.

arctic_a0001.lab	arctic_a0001.sl	arctic_a0002.lab	arctic_a0002.sl
arctic_a0003.lab	arctic_a0003.sl	arctic_a0004.lab	arctic_a0004.sl
arctic_a0005.lab	arctic_a0005.sl	arctic_a0006.lab	arctic_a0006.sl
arctic_a0007.lab	arctic_a0007.sl	arctic_a0008.lab	arctic_a0008.sl
arctic_a0009.lab	arctic_a0009.sl	arctic_a0010.lab	arctic_a0010.sl
arctic_a0011.lab	arctic_a0011.sl	arctic_a0012.lab	arctic_a0012.sl
arctic_a0013.lab	arctic_a0013.sl	arctic_a0014.lab	arctic_a0014.sl
arctic_a0015.lab	arctic_a0015.sl	arctic_a0016.lab	arctic_a0016.sl
arctic_a0017.lab	arctic_a0017.sl	arctic_a0018.lab	arctic_a0018.sl
arctic_a0019.lab	arctic_a0019.sl	arctic_a0020.lab	arctic_a0020.sl
arctic_a0021.lab	arctic_a0021.sl	arctic_a0022.lab	arctic_a0022.sl
arctic_a0023.lab	arctic_a0023.sl	arctic_a0024.lab	arctic_a0024.sl
arctic_a0025.lab	arctic_a0025.sl	arctic_a0026.lab	arctic_a0026.sl
arctic_a0027.lab	arctic_a0027.sl	arctic_a0028.lab	arctic_a0028.sl
arctic_a0029.lab	arctic_a0029.sl	arctic_a0030.lab	arctic_a0030.sl
arctic_a0031.lab	arctic_a0031.sl	arctic_a0032.lab	arctic_a0032.sl
arctic_a0033.lab	arctic_a0033.sl	arctic_a0034.lab	arctic_a0034.sl
arctic_a0035.lab	arctic_a0035.sl	arctic_a0036.lab	arctic_a0036.sl
arctic_a0037.lab	arctic_a0037.sl	arctic_a0038.lab	arctic_a0038.sl
arctic_a0039.lab	arctic_a0039.sl	arctic_a0040.lab	arctic_a0040.sl
arctic_a0041.lab	arctic_a0041.sl	arctic_a0042.lab	arctic_a0042.sl
arctic_a0043.lab	arctic_a0043.sl	arctic_a0044.lab	arctic_a0044.sl
arctic_a0045.lab	arctic_a0045.sl	arctic_a0046.lab	arctic_a0046.sl
arctic_a0047.lab	arctic_a0047.sl	arctic_a0048.lab	arctic_a0048.sl
arctic_a0049.lab	arctic_a0049.sl	arctic_a0050.lab	arctic_a0050.sl
arctic_a0051.lab	arctic_a0051.sl	arctic_a0052.lab	arctic_a0052.sl
arctic_a0053.lab	arctic_a0053.sl	arctic_a0054.lab	arctic_a0054.sl
arctic_a0055.lab	arctic_a0055.sl	arctic_a0056.lab	arctic_a0056.sl
arctic_a0057.lab	arctic_a0057.sl	arctic_a0058.lab	arctic_a0058.sl
arctic_a0059.lab	arctic_a0059.sl	arctic_a0060.lab	arctic_a0060.sl
arctic_a0061.lab	arctic_a0061.sl	arctic_a0062.lab	arctic_a0062.sl

Fig. 1 Label of Dictionary

1	7																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																				
---	---	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--

Fig. 2 Contents of Label of Speech File

## Markov Model Procces

Classify weather into three states

- State 1: rain or snow
- State 2: cloudy
- State 3: sunny

By carefully examining the weather of some city for a long time, we found following weather change pattern.

TABLE 1  
WEATHER PROBABILITY

		Tomorrow		
		snow	cloudy	sunny
today	snow	0.4	0.3	0.3
	cloudy	0.2	0.6	0.2
	sunny	0.1	0.1	0.8

At every state tomorrow weather is depending upon today state Visual representation.

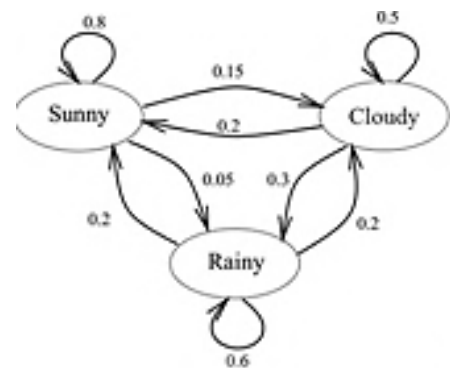


Fig. 3 State Diagram of Weather

Each state corresponds to one observation,

Sum of outgoing edge weights is one

Observable states,  $\{1, 2, \dots, N\}$

Observed sequence,  $\{q_1, q_2, \dots, q_n\}$

1st order Markov assumption

$P(q_t = j \mid q_{t-1} = i, q_{t-2} = k, \dots)$

$= P(q_t = j \mid q_{t-1} = i)$

We denote the time instants associated with the state changes as  $t=1,2,3,\dots,n$  and actual state at time  $t$  as  $q_t$ .

### B. Markov Model: Sequence Probability

• Question: What is the probability that the weather for the next 7 days will be “sun-sun-rain-rain-sun-cloudy-sun” when today is sunny?

$S$  : rain,  $S$  : cloudy,  $S$  : sunny

$$P(O | \text{model}) = P(S_3, S_3, S_3, S_1, S_1, S_3, S_2, S_3 | \text{model}) =$$

$$P(S_3).P(S_3|S_3).P(S_3|S_3).P(S_1|S_3).P(S_1|S_1).P(S_3|S_1).$$

$$P(S_2|S_3).P(S_3|S_2) = \pi_3.a_{33}.a_{33}.a_{31}.a_{11}.a_{13}.a_{32}.a_{23}$$

$$=1.(0.8)(0.8)(0.1)(0.4)(0.3)(0.1)(0.2) =1.536 \times 10^{-4}$$

### C. HMM Algorithms

There are three algorithms used in HMM

1. Viterbi Algorithm,
2. Baum-Welch Algorithm,
3. K-Means.

#### 1. Viterbi Algorithm

- Purpose – An analysis for internal processing result, The best, the most likely state sequence for Internal segmentation
- Alignment of observation and state transition
- Dynamic programming technique

Key idea – Span a lattice of  $N$  states and  $T$  times, to keep the probability and the previous node of the most probable path coming to each state  $i$  at time  $t$

Introduction Phase:

$$\delta_1(i) = \pi_i b_i(x_1), \psi_1(i) = 0$$

Recursion Phase:

$$\delta_{t+1}(j) = \max_{1 \leq i \leq N} \delta_t(i) a_{ij} b_j(x_{t+1})$$

$$\psi_{t+1}(j) = \arg \max_{1 \leq i \leq N} \delta_t(i) a_{ij}$$

Termination Phase:

$$P^* = \max_{1 \leq i \leq N} \delta_T(i)$$

$$q_T^* = \arg \max_{1 \leq i \leq N} \delta_T(i)$$

Path backtracking Phase:

$$q_t^* = \psi_{t+1}(q_{t+1}^*), t = T-1, \dots, 1,$$

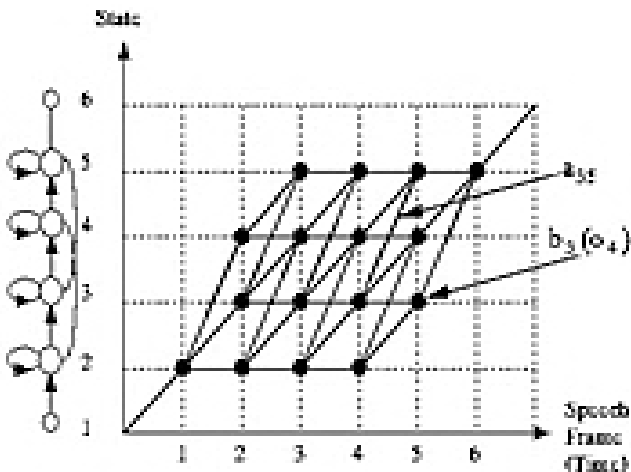


Fig 4. The Viterbi Algorithm for Isolated Word Recognition

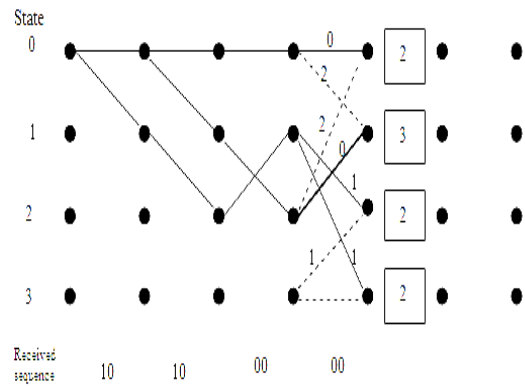


Fig. 5 Mapping of Model

#### 2. Baum-Welch Algorithm

1. The Baum-Welch algorithm can be used to train an HMM to model a set of sequence data.
2. The algorithm starts with an initial model and iteratively updates it until convergence.
3. The algorithm is guaranteed to converge to an HMM that locally maximizes the probability of the training data given the model.

Baum-Welch Reestimation (EM)

- Iterative procedures that locally maximizes
- Convergence proven
- MLE statistic estimation

The following table shows a one run of the Baum-Welch algorithm:

‘ $a$ ’ is considered as start and ending symbols of the sequence string and through multiple iterations the probability values will be converges.

$a \rightarrow a$	$a \rightarrow b$	$b \rightarrow a$	$a \rightarrow b$	$b \rightarrow b$	$b \rightarrow a$	$P(\text{path})$	$q \xrightarrow{a} r$	$r \xrightarrow{b} q$	$q \xrightarrow{a} q$	$q \xrightarrow{b} q$
q	r	q	r	q	q	0.00077	0.00154	0.00154	0	0.00077
q	r	q	q	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	r	q	q	0.00442	0.00442	0.00442	0.00442	0.00884
q	q	q	q	q	q	0.02548	0.0	0.000	0.05096	0.07644
Rounded Total $\rightarrow$							0.035	0.01	0.01	0.06
New Probabilities (P) $\rightarrow$							0.06	1.0	0.36	0.581
State sequences										

Fig 6 Example of Baum-Welch Algorithm

#### 3. K-Means Clustering

HMM kind of algorithm is well known to be sensitive to its initialization point. The problem of this initialization point choice is addressed in this paper: a model with a very large number of states which describe training sequences with accuracy is first constructed. The number of states is then reduced using a k-mean algorithm on the state. This algorithm is compared to other methods based on a k-mean algorithm on the data with numerical simulations.

The clustering algorithm partitions a dataset into a fixed number of clusters supplied by the user. Hidden Markov Model (HMM) based clustering method, which identifies a

suitable number of clusters in a given dataset without using prior knowledge about the number of clusters. Initially, the dataset is partitioned into windows of fixed size based on the HMM log likelihood values. This provides a framework for identifying the most appropriate number of clusters (windows of varying sizes). After determining the number of clusters, the data values are then labelled and allocated to clusters.

#### *D. Features of Audio*

The following features will be extracted from the audio file for training the HMM to get the speech.

1. Energy Entropy Standard Deviation (std)
2. Signal Energy Std by Mean(average) Ratio
3. Zero Crossing Rate Std
4. Spectral Rolloff Std
5. Spectral Centroid Std
6. Spectral Flux Std by Mean Ratio

### III. CONCLUSION

In this paper it is shown how HMM can be used to generate speech from the text file. HMM uses dictionaries, labelling and the three different algorithms like K-means, Viterbi Algorithm and Baum-Welch Algorithm to train and after that the speech will be generated from the text files. The HMM can be used for image processing and forecasting also.

### REFERENCES

- [1] Junichi Yamagishi, Takashi Nose, Heiga Zen, Zhen-Hua Ling, Tomoki Toda, Keiichi Tokuda, Simon King and Steve Renals, "Robust Speaker-Adaptive HMM-Based Text-to-Speech Synthesis", IEEE Transactions on Audio, Speech, and Language Processing, Vol.17, No. 6, August 2009
- [2] Agni Dika1, Adnan Maxhunil, Avni Rexhepi, "The principles of designing of algorithm for speech synthesis from texts written in Albanian language", IJCSI International Journal of Computer Science Issues, Vol. 9, Issue 3, No 3, May 2012
- [3] L. R. Rabiner, "A tutorial on hidden markov models and selected applications in speech recognition," Proceedings of the IEEE, vol. 77, no. 2, pp. 257–286, February 1989.
- [4] A. W. Black, K. Lenzo, Building voices in the Festival speech synthesis system, 2000, <http://festvox.org/bsv>.
- [5] Kevin Murphy, "HMM toolbox for Matlab", freely downloadable SW written in Matlab, <http://www.cs.ubc.ca/~murphyk/Software/HMM/hmm.html>
- [6] Juang BH, Rabiner LR, "Mixture Autoregressive Hidden Markov Models for Speech Signals", IEEE Trans Acoustics, Speech and Signal Processing 33: 1404-1413,1985
- [7] Qystein Birkenes, Tomoko Matsui, Kunio Tanabe, Sabato Marco Siniscalchi, Tor Andre Myrvoll, and Magne Hallstein Johnsen, "Penalized Logistic Regression with HMM LogLikelihood Regressors for Speech Recognition", IEEE Transactions on Audio, Speech, and Language Processing Vol. 18, No. 6, pp. 1440-1454, August 2010.
- [8] Tokuda, K., Nankaku, Y., Toda, T., Zen, H., Yamagishi, J., Oura, K., "Speech synthesis based on hidden markov models", In: Proceedings of the IEEE, Vol. 101(5), pp. 1234–1252 (2013)
- [9] HMM-based Speech Synthesis System (HTS). <http://hts.sp.nitech.ac.jp/>
- [10] L.E. Baum, T. Petrie, "Statistical inference for probabilistic functions of finite state Markov chains", Ann. Math. Stat., 37 (1966), pp. 1554-1563
- [11] Olivier Cappé Eric Moulines Tobias Rydén, Inference in Hidden Markov Mode,l Springer Series in Statistics, 2005.
- [12] Y. Ariki, M.A. Jack, "Enhanced time duration constraints in hidden Markov modelling for phoneme recognition", Electronics Letters, 25 (13) (22 June 1989), pp. 824-825.