# Different Methods Review for Speech to Text and Text to Speech Conversion

Deep Kothadiya
Post Graduate Student
MIT, Poud Road

Nitin Pise, PhD
Professor
MIT, Poud Road

Mangesh Bedekar
Professor
MIT, Poud Road

## ABSTRACT

In the instant corporation, transmission is the primary fundamental to momentum. Transitory information, to the correct person with the correct aspect is very essential, not just on an industry level, but also on an individual position. Nature is inspiring in the direction of digitization and the mechanisms of intercommunication. Telephone calling, e-mails, text memorandums belong to a fundamental element of signal communication in this tech-intellect nature. In procedures to distribute the intention of adequate transmission intervening two endpoints without obstacles, numerous utilizations have shown up the impression, which operates as an intermediary and helps in efficiently transmitting signals in the scheme of text or speech messages accomplished huge structure of webs. Most of these implementations discover the Usage of tasks essentially articulatory and acoustic-positioned speech recognition, reorganization from audio messages to text, and then text to artificial speech signals, vocabulary interpretation amidst individual leftovers. Researchers will be penetrating distinct algorithms and techniques that are enforced to obtain the specified utilitarian.

## General Terms
Machine learning, Algorithm

## Keywords
Speech to Text, Text to Speech, Conversion

## 1. INTRODUCTION

Cellular Phones have become an essential origin of transmission for modernized civilization. Authors can make text messages and calls from the origin to a goal efficiently. It is acknowledged that spoken transmission is the ultimate suitable phase of briefing on and comprehending the legitimate knowledge, averting inappropriate citations. To fulfill the inconsistency over a expanded interval, spoken transmission can occur conveniently on cellular phone calling. A track-cracking deviation has currently appear to show in the SMS automation utilizing the voice perception technology, where speech messages are being transformed to text messages. Absolutely a minor application serviced to help the wounded make usage of STT, TTS, and interpretation. They can also be recycled for additional implementations, taking an illustration: alexa an inventive computerized associate achieved on a photoelectric appliance, to expedite customer interaction with an appliance, and to assist the customer more efficiently enlist with regional and/or distant utility [10] builds usage of distinction Communications speech perception and text-to-speech (TTS) technology. Authors will pay attention to the distinctive forms of audio speech, voice recognition, speech to text reconstruction, text to speech reorganization and voice interpretation. Under speech the recognition Authors will pursue the mechanism that is pre-prominence of semaphores, recognition of the semaphores and feature extraction which boost us in training and testing mechanisms. There are different prototypes used for this aspiration but Dynamic time cape, which is utilized for distance measurement and feature extraction between appearances of semaphores and Hidden Markov Model (HMM) which is a hypothetical miniature and is used to associate distinct articulates of evolution with each other is mainly serviced. Correspondingly for transformation of speech to text authors use HMM and DTW archetypes, onward with different Neural Network miniatures since they endeavor properly with speaker adaptation phoneme categorization and segregated word perception. Point to point ASR is also essentially approved as of late 2014 to obtain identical outcomes. Voice fusion works correctly in comforting modified tokenized words to artificial human voice. Various gadget adaptation schemes, as good as appliances will also be compared and reviewed. Ensuring are the pieces of voice manufacture, which are considered up to while implementation use distinct speech associated range of capabilities [5].

 - Enunciation
 - Angle deviation
 - Voice (containing aeromechanical elements of respiration)
 - Phonation (Producing sound)
 - Fluency

## 2. LITERATURE REVIEW

In this literature review authors have determined the actual procedure for voice recognition, speech to text transformation and techniques of machine-learning.

## 2.1 SPEECH RECOGNITION

Speech Recurrence is the capability of a computer to establish phonemes and remarks in verbal vocabulary and change them into engine-coherent composition. Speech Recognition scheme can be confidential on fundamental of the succeeding frameworks [10]:

a. Speakers: All kind of speakers have various speech types. The designs are drafted for an independent speaker or a particular speaker.

b. Articulate Sound: The technique the speaker utters also shows an aspect in voice acknowledgment. Some ideals can identify either single assertion or split assertion with a halt in between.

c. Terminology: The amount of the terminology portrays an essential position in establishing the intricacy, efficiency, and accuracy of the structure.

### 2.1.1 Primitive Voice Identification Form
Each voice recognition scheme follows some ideal measures
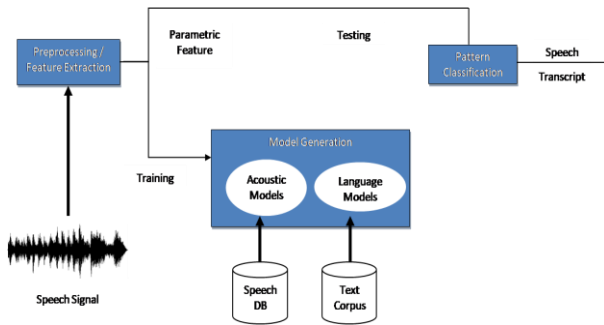
as exhibited in Diagram 1 [10].



**Figure 1: Architecture for Speech Recognition System**

### 2.1.1.1 Pre-processing
The voice analog beckon is revolutionized into digital beckons for afterward transforming. This digital beckon is transported to the initial form and refined to spectrally prostrate the beckons. This stimulates in expanding the gesture's efficiency at an excessive density.

### 2.1.1.2 Feature Extraction
This process discovers collection of criterions of assertion that accept an interaction with voice gestures. These criterions, acknowledged as features, are estimated over transforming of the acoustic waveform. The major focal point is to figure out a series of component vectors contributing a brief portrayal of the delivered input gesture. Frequently applied feature extraction methods are considered beneath::

### A. Linear Predictive Coding (LPC)
The essential opinion is that the audio vocal sound fragment can be roughed as a precise consolidation of previous speech cases. Diagram 2 displays the LPC method [9]. The computerized beckon is held up into frameworks of N cases. Then every sample structure is framed to decrease beckon discontinuation. Each fabricated window is then auto-corresponding. The final phase is the LPC examination and determination, which brings each skeleton of auto interactions into LPC criterion collection.
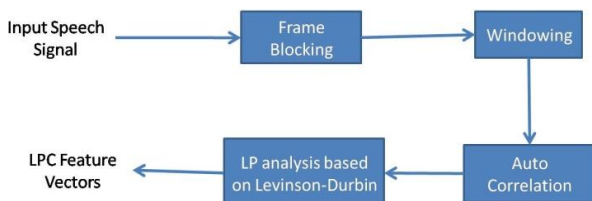


**Figure 2: LPC Feature Extraction Process**

### B. Mel-Frequency Cepstral Coefficient (MFCC)
It is an exact dynamic craft and usage human hearing opinion scheme. MFCC appeals positive phases to the input signal: fabricate: voice flutter- constitute is pruned to discard conflict if current; Windowing: reduces the discontinuation in the signal; various Fourier Transform: reforms each framework from time territory to frequency territory; Mel refine Bank Algorithm: the gesture is operated against the Mel range to mime human listening [9].

### C. Dynamic Time Warping
This method is serviced for quoting the correlation between two-time sequences which may change in momentum, positioned on productive prioritizing. It is objective at coordinating two series of characteristic vectors iteratively until an optimum match between them is established.

### 2.1.1.3 Acoustic Models
The basic components of Automated Speech Recognition (ASR) scheme where a relation between the utterance and acoustic info is fixed. Preparation equivalence between the fundamentals voice entities and the sound inspections establishes.

### 2.1.1.4 Jargon Models
This standard persuades the possibility of a discussion existence subsequently a word order. It consists of the anatomical restraints feasible in the vocabulary to achieve the contingency of existence. The vocabulary miniature determines word and remark with identical speech.

### 2.1.1.5 Design Distribution
It is the method of correlating the anonymous design with actual voice quotation design and estimated correlation between them. Subsequently determine the discipline of the scheme at the time of verification. Designs are confidential to identify the voice. For Design identical various approaches are [3]:

### A. Pattern Based Proposition
This way has a set of speech designs which are gathered as a citation characterizing language words. Using the reference pattern speech is perceived by corresponding the spoken word [14].

### B. Proficiency Established Proposition
This technique takes a collection of features from the voice and then prepares the scheme to create a collection of manufacturing guidelines undoubtedly from the cases.

### C. Neural Network Based Proposition
This method is accomplished to determine further complex acknowledgment exercise. The essential plan is to integrate knowledge and compile from a variation of ability origin with the complication at hand [2].

### D. Analytical Based Proposition
In this method, different audio is formed statistically utilizing practice mechanisms.

### 2.1.2 Speech to Text Transformation Methods
The technique of modifying uttered words into drafted texts. It is compatible with speech concession but the recent is adopted to interpret the expanded operation of voice comprehension. STT pursues the equivalent fundamentals and actions of speech data perception, with various associations of various methods individually step. Some extensively used transformation schemes are examined below.

### 2.1.2.1 Hidden Markov Model (HMM)
This model is an analytical ideal usage in voice recognition because a voice gestures can be considered as a short-time stagnant gesture signal or compose reasonable stagnant gesture. HMM patterns are beneficial for actual time voice to text transformation for cell phone customers [3]. It relied on the ensuing specifications:

### A. Recognizance Efficiency
Recognizance is the procedure of correlating the anonymous analysis design with every voice track collection citation standard and estimating a determination of correlation between the evaluation design and each mentioned system. It is the exceedingly substantial aspect of any perception scheme, exquisitely independent of the speaker and it should be 100%.

*B. Recognition Momentum*

Consumers feel anxious and the scheme drops its implication if the scheme acquires a more than enough amount of time to identify the voice signal. The signals endure the ensuing acts: [6]

*C. Pre-refining*

The intake voice audio signals indicator are transformed into speech structures and contribute a singular sampling, trimming speech clatter.

*D. HMM Preliminaries*

Preparation associates establishing a relevant vessel delegate of the characteristics of a collection utilizing one or more experiment designs that resemble to voice sounds of the identical course.

*E. HMM Acknowledgment*

It is the method of correlating the anonymous examination design with each voice collection citation design and estimating determinability (distance). Maximum possibility is used for recognition.

### 2.1.2.2 *Artificial Neural Network elegant (ANN) based Cuckoo Search Optimization*

This method is used for improved conversation, improved recognition and to discard rejected noise. Automatic Recognition Speech is built for an improved association of machine and individual cooperation. For the identical, a 3-stage case is pursued: [7]

Pre-transforming of the voice signals is the remarkably valuable chunk of voice appreciation which is performed to eliminate the preventable waveform of the indicators. The signals are filled to the tremendous occur penetrates to eliminate the environment clatters.

2 sets of acoustic aspects are obtained from the voice gesture. They are Mel Frequency Cep-strum Coefficients (MFCC) and Linear Predictive Coding coefficients (LPCC).

*A. Categorization*

ANN is usage as the classifier. The auditory structure is a 3-tiered classifier with n intake nodes, l invisible nodes and k output nodes. In CSO (Cuckoo Search Optimization), ANN is achieved by two-tiered Feed Forward Back propagation Neural Network (FFBNN) with 3 entities; 2 input entities, 3 covered entities and 1 output entity. Here, the input tier consists of 2 inputs accepting 2 characteristics obtained which are MFCC and LPCC characteristics. These emphases are given as input in which structures grab skilled and it generates a reciprocal output.
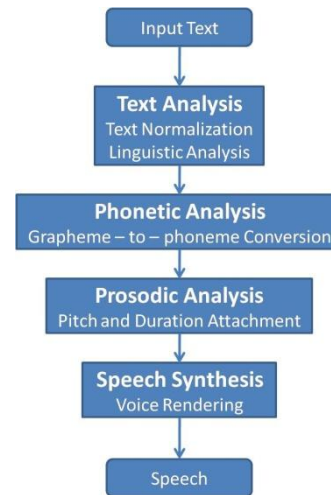


**Figure 3: Text to speech system flow**

## 2.2 TEXT TO SPEECH CONVERSION

In this action first intake content is evaluated, then prepared and accepted, and then the text is transformed to digital audile and then uttered. Diagram 3 displays the Figure of TTS. The diagram displays all the acts concerned in the text to speech changeover but the essential aspects of TTS schemes are [2]:

### 2.2.1 *Content Transforming*

The intake text is evaluated, distributed (manages abbreviation, acronyms and equal antagonist the content) and interpreted into linguistic or phonetic delegation.

### 2.2.2 *Speech Integration*

Some of the speech synthesis approaches are [2]:

### 2.2.2.1 *Coherent fusion*

Use acoustic model and mechanical for speech creation. It cultivates comprehensible fabricated voice but it is distant from distant sound and hence not broadly used.

### 2.2.2.2 *Perpetual fusion*

In this structure, portrayals of particular speech sections are reserved on a parametric fundamental. There are two elemental frameworks in perpetual synthesis, correlation and avalanche, but for improved efficiency, some sort of association of these 2 formations is used. An avalanche perpetual synthesizer resides of symphony-passage vibrators associated in sequences. The output of each perpetual resonator is enforced to the input of the subsequent one. The avalanche architecture obligation only perpetual densities as restraint instruction. A correlation perpetual synthesizer exists of resonators affiliated in complementary. The excitement signal is enforced to whole perpetual together and their outputs are compiled. [2]

### 2.2.2.3 *Integrative fusion*

This procedure incorporates sound by integrating precise cases of voice hailed entities. It is worn in speech fusion to develop a customer's definitive order of sound from a collection of data assembled from the documenting of additional strings. Entities for Integrative synthesis are [2]: Telephone- a particular entity of sound; Ditelephone is designated as the indication from either centriole of a telephone or mark of least modification within the telephone to the identical mark in the next telephone; Tritelephone- is a segment of the beckon contagious in a continuity expatriation from intermediate of a telephone absolutely by way of the afterward one to the intermediate of a third.

## 3. OBSERVATION

**Table 1. Different Structure with Obstacles**

| MODELS | METHODS | OBSTACLES |
|---|---|---|
| Feature Extraction | Linear Predictive Coding (LPC) | Equal weighted signal on linear extent while human ear receptive is numerical |
| | Mel-Frequency Cestrum Co-efficient (MFCC) | Values are not vigorous in the occupancy of obsessive noise |
| | Dynamic Time Warping (DTW) | Complication in selection of the pattern |
| Pattern Recognition | Pattern established | Pre-documented templates are steady. Continuously speech matching is not attainable |
| | Knowledge situated | categorical form alteration in speech is crucial to acquire so this technique is absurd |
| | Neural Based | |
| | Statistical positioned | deficient efficiency of preceding form |
| | Hidden Markov Model (HMM) | inadequacy in bias possessions for allocation |
| Speech to Text Conversion | Artificial optimization | Time fluctuation of voice |
| Text to Speech Conversion | Coherent fusion | Result is distant from original speech |
| | Perpetual fusion | Does not use human voice fragments at executed runtime |
| | Integrative fusion | Intricate procedure |
| Machine Interpretation | Hybrid , statistical | Require initial source of speech data |

## 4. CONCLUSION

Here Authors learned about different methodology of speech to text and text to speech conversion and their application usability. From this various aspects Authors found that HMM works better in conversion just with computational feasibility drawback. Apparently in text to speech cascade fusion is best solution authors got. For future perspective authors have focus of ability to learn fast, smoothness in word correction and data accretion.

## 5. REFERENCES

[1] K. Dutta and K. K. Sarma, "Multiple Feature Extraction for RNN-based Assamese Speech Recognition for Speech to Text Conversion Application", International Conference on Communications, Devices and Intelligent Systems (CODIS), IEEE, 2012.

[2] F. Seide, G. Li, D. Yu,Conversational Speech Transcription Using Context-Dependent Deep Neural Networks, In Interspeech, pp. 437440, 2011.

[3] y Keiichi Tokuda, Yoshihiko Nankaku, Tomoki Toda, Heiga Zen, Speech Synthesis Based on Hidden Markov Models, Proceedings of the IEEE — Vol. 101, No. 5, May 2013. Junichi Yamagishi, Member IEEE, and Keiichiro Oura

[4] Mathias De Wachter, Mike Matton, Kris Demuynck, Patrick Wambacq, Template Based Continuous Speech Recognition,IEEE Transs. On Audio, Speech Language Processing, vol.15, issue 4,pp 1377-1390, May 2007.

[5] Lawrence Rabiner, Biing-Hwang Juang, B.Yegnanarayana, Fundamentals of Speech Recognition.

[6] Ms. Anuja Jadhav, Prof. Arvind Patil, Real Time Speech to Text Con- verter for Mobile Users, National Conference on Innovative Paradigms in Engineering Technology (NCIPET-2012) Proceedings published by International Journal of Computer Applications (IJCA)

[7] Sunanda Mendiratta, Dr. Neelam Turk, Dr. Dipali Bansal, Speech Recognition by Cuckoo Search Optimization based Artificial Neural Network Classifier, 2015 International Conference on Soft Computing Techniques and Implementations- (ICSCTI) Department of ECE, FET, MRIU, Faridabad, India, Oct 8-10, 2015.

[8] Suhas R. Mache, Manasi R. Baheti, C. Namrata Mahender, Review on Text-To-Speech Synthesizer, International Journal of Advanced Research in Computer and Communication Engineering Vol. 4, Issue 8, August 2015.

[9] Pratik K. Kurzekar, Ratnadeep R. Deshmukh, Vishal B. Waghmare, Pukhraj P. Shrishrimal, A Comparative Study of Feature Extraction Techniques for Speech Recognition System, International Journal of Innovative Research in Science, Engineering and Technology (An ISO 3297: 2007 Certified Organization) Vol. 3, Issue 12, December 2014.

[10] Suman K. Saksamudre, P.P. Shrishrimal, R.R. Deshmukh, A Review on Different Approaches for Speech Recognition System, International Journal of Computer Applications (0975 8887) Volume 115 No. 22, April 2015.

[11] M. Vyas, "A Gaussian Mixture Model Based Speech Recognition System Using Matlab", SIPIJ, Vol.4, No.4, August 2013.

[12] N. Srivastava, "Speech Recognition using Artificial Neural Network", IJESIT, Volume 3, Issue 3, May 2014.