# PROJECT REPORT

# EXECUTION OF A VOICE - BASED ATTENDENCE SYSTEM

BY

## SIDDHARTH SAGAR BARPANDA

## 111EI0477

UNDER THE GUIDANCE OF

## PROF. U. K. SAHOO



DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING

NIT ROURKELA

# ACKNOWLEDGEMENTS

# SYNOPSIS

Speech Recognition is the methodology of consequently perceiving a certain word talked by a specific speaker taking into account singular data included in speech waves. This system makes it conceivable to utilize the speaker's voice to confirm his/her personality and give controlled access to administrations like voice based biometrics, database access administrations, voice based dialling, phone message and remote access to PCs.

Speech processing front end for extricating the feature set is a critical stage in any voice recognition system. The ideal list of capabilities is still not yet chosen however the limitless endeavours of scientists. There are numerous sorts of highlights, which are determined distinctively and have great effect on the acknowledgment rate. This project shows one of the strategies to extract the feature from a voice signal, which can be utilized as a part of speech acknowledgment system.

The key is to change the speech wave to some kind of parametric representation (at an impressively lower data rate) for further examination and processing. This is frequently known as the voice processing front end. An extensive variety of potential outcomes exist for parametrically speaking to the discourse signal for the speaker acknowledgment undertaking, for example, Mel-Frequency Cepstrum Coefficients (MFCC), Linear Prediction Coding (LPC), and others. MFCC is maybe the best known and generally prominent, furthermore, these will be utilized as a part of this undertaking. MFCCs are in view of the known variety of the human ear's discriminating transmission capacities with recurrence channels dispersed sprightly at low frequencies and logarithmically at high frequencies have been utilized to catch the phonetically essential qualities of discourse. Nonetheless, another key normal for discourse is semi stationary, i.e. it is brief time stationary which is contemplated and investigated utilizing brief time, recurrence space examination.

In this project work, I have built a straightforward yet completed and agent automatic speaker recognition (ASR) framework, as connected to a voice based attention framework, i.e., a speech based access control system. To attain to this, I had to first

made a relative investigation of the MFCC approach with the Time space approach for acknowledgment by simulating both these strategies utilizing MATLAB 7.0 and investigating the consistency of acknowledgment utilizing both the procedures.
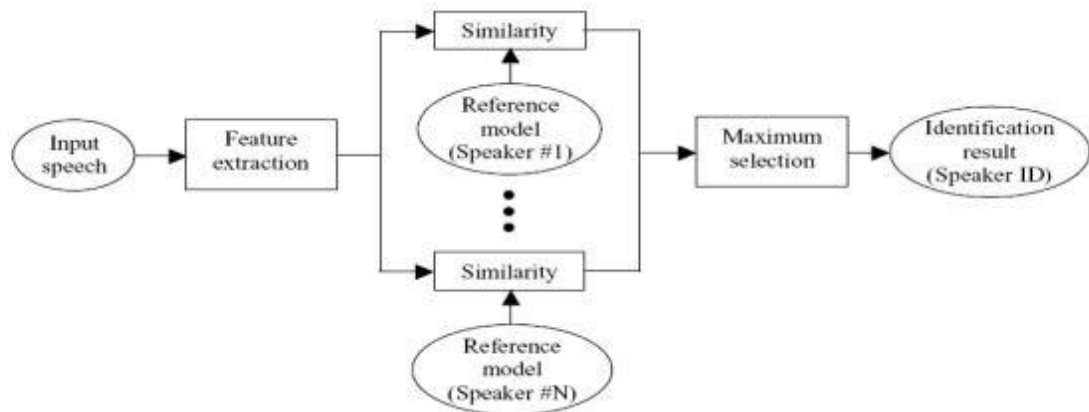
The voice based attendance system is based with respect to confined or one word recognition. A specific speaker articulates the secret word once in the instructional course so as to prepare and store the highlights of the entrance word. While in the testing session the speaker articulates the secret key again to accomplish acknowledgment if there is a match. The highlight vectors interesting to that speaker are acquired in the preparation stage and this is made utilization of later on to allow validation to the same speaker who at the end of the day expresses the same word in the testing stage. At this stage a gate crasher can likewise test the framework to test the inalienable security include by expressing the same word.
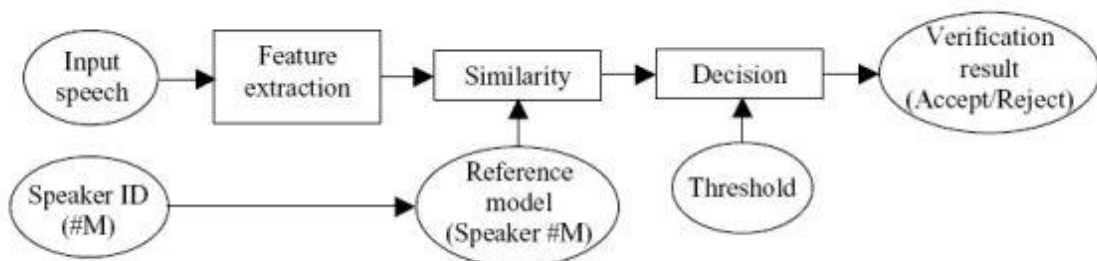
# INTRODUCTION

From an innovative viewpoint it is conceivable to recognize two wide sorts of ASR, i.e., Direct Voice Input (DVI) and Large Vocabulary Continuous Speech Recognition (LVCSR). DVI gadgets are essentially gone for voice order and control, though LVCSR frameworks are utilized for structure filling or voice-based archive creation. In both cases the basic innovation is pretty much the same. DVI frameworks are commonly designed for little to medium estimated vocabularies (up to a few hundred words) and may utilize word or expression spotting strategies. Likewise, DVI frameworks are typically needed to react promptly to a voice order. LVCSR frameworks include vocabularies of maybe countless words, and are commonly designed to decipher persistent discourse. Likewise, LVCSR require not be performed continuously - for instance, no less than one seller has offered a phone based transcription benefit in which the deciphered archive is messaged back to the user.

From an application perspective, the advantages of utilizing ASR get from giving an additional correspondence direct in hands-occupied eyes-occupied Human Machine Interaction (HMI), or basically from the way that talking can be speedier than writing. Likewise, whilst addressing a machine can't be portrayed as common, it can by and by be considered instinctive; as one ASR notice pronounced "you have been learning since conception the just aptitude expected to utilize our framework".

At the most normal, all speaker acknowledgment frameworks contain two principle modules: feature extraction and feature matching. Extraction is the methodology that concentrates a little measure of information from the voice flag that can later be utilized to speak to every speaker. Feature matching includes the real methodology to recognize the obscure speaker by contrasting removed highlights from his/her voice info with the ones from an arrangement of known speakers. We will examine this two modules in subtle in later parts.

Speaker Identification (Training)



Speaker verification (Testing)

Voice recognition is a hard assignment and it is still a dynamic research range. Automatic Speech Recognition works in light of the reason that a person's voice shows attributes that are special to the speaker itself. However this undertaking has been tested by the exceptionally variation of information speech waves. The rule wellspring of change is the speaker himself. Voice motions in training and testing phase can be extraordinarily diverse because of numerous truths, for example, individuals voice change with time, wellbeing conditions (e.g. the speaker has a cool), talking rates, and so forth. There are additionally different components, past speaker variability, that present a test to voice recognition innovation. Examples of these are acoustical clamor and varieties in recording situations (e.g. speaker utilizes diverse phone handsets). The test would be make the framework "Robust".

# 1.    ASR SYSTEMS

Speech processing is the investigation of speech signs and the processing strategies for these signals. These signals are typically transformed in a computerized representation whereby speech processing can be seen as the cooperation of advanced sign processing and regular dialect processing.

Speech coding is the pressurization of speech (into a code) for communication with speech codecs that utilize sound voice processing and speech processing strategies. The strategies utilized are like that as a part of sound information pressure and sound coding where learning in psycho-acoustics is utilized to transmit just information that is pertinent to the human sound-related framework. Case in point, in tight band speech coding, just data in the recurrence band of 400 Hz to 3500 Hz is transmitted however the reproduced sign is still satisfactory for comprehensibility.

Nonetheless, speech coding contrasts from voice coding in that there is a considerable measure more measurable data accessible about the properties of speech. Furthermore, some sound-related data which is applicable in sound coding can be pointless in the speech coding connection. In speech coding, the most imperative measure is conservation of clarity and charm of speech, with an obliged measure of transmitted information.

Speech synthesis is the simulated creation of human voice. A text-to-speech (TTS) framework changes over typical dialect content into discourse; different frameworks render typical phonetic representations like phonetic translations into audio signal. Synthesized speech can likewise be made by connecting bits of recorded voice that are put away in a database. Frameworks contrast in the measure of the put away discourse units; a framework that stores telephones or diphones gives the biggest yield range, yet may need clarity. For particular utilization areas, the capacity of whole words or sentences considers amazing yield. Then again, a synthesizer can fuse a model of the vocal tract and other human voice attributes to make a totally "synthetic" voice yield.

Voice issues that oblige voice analysis most normally start from the vocal cords since it is the audio source and is in this manner most effectively subject to tiring. In any case, analysis of the vocal ropes is physically hard. The area of the vocal lines viably precludes direct estimation of development. Imaging strategies, for example, x-beams or ultrasounds don't work in light of the fact that the vocal ropes are encompassed via ligament which mutilates picture quality. Developments in the vocal lines are fast, essential frequencies are ordinarily somewhere around 80 and 300 Hz, subsequently counteracting utilization of standard feature. Rapid features give a choice however so as to see the vocal lines the cam must be situated in the throat which makes talking rather troublesome.

Most vital aberrant techniques are converse sifting of sound recordings and electroglottographs (EGG). In converse sifting techniques, the discourse sound is recorded outside the mouth and after that sifted by a scientific system to evacuate the impacts of the vocal tract. This strategy delivers an assessment of the waveform of the weight beat which again contrarily shows the developments of the vocal strings. The other sort of opposite evidence is the electroglottographs, which works with terminals connected to the subject's throat near to the vocal strings. Changes in conductivity of the throat demonstrate conversely how extensive a segment of the vocal strings are touching one another. It accordingly yields one-dimensional data of the contact region. Neither reverse sifting nor is EGG consequently adequate to totally portray the glottal development and give just backhanded confirmation of that development.

Speech recognition is the methodology by which a PC (or other sort of machine) distinguishes talked words. Essentially, it means conversing with your PC, and having it accurately perceive what you are stating. This is the way to any speech application.

# 2.    SPEECH RECOGNITION BASICS

## 2.1 Utterance

An utterance is the vocalization (talking) of a word or words that speak to a solitary intending to the PC. Utterances can be a solitary word, a couple of words, a sentence, or even different sentences.

## 2.2 Speaker Dependence

Speaker dependent frameworks are planned around a particular speaker. They for the most part are more precise for the right speaker, however substantially less exact for different speakers. They expect the speaker will talk in a reliable voice and rhythm. Speaker independent frameworks are intended for a mixture of speakers. Versatile frameworks generally begin as speaker independent frameworks and use preparing procedures to adjust to the speaker to build their acknowledgment exactness.

## 2.3 Vocabularies

Vocabularies (or word references) are arrangements of words or articulations that can be perceived by the SR framework. By and large, littler vocabularies are less demanding for a PC to perceive, while bigger vocabularies are more hard. Dissimilar to typical lexicons, every entrance doesn't need to be a solitary word. They can be the length of a sentence or two. Littler vocabularies can have as few as 1 or 2 perceived expressions (e.g." Wake Up"), while vast vocabularies can have a hundred thousand words or more!

## 2.4 Accuracy

The capacity of a recognizer could be inspected by measuring its accuracy - or how well it perceives expressions. This incorporates effectively distinguishing an articulation as well as recognizing if the talked expression is not in its vocabulary. Great ASR frameworks have an accuracy of 98% or more! The worthy accuracy of a framework truly relies upon the application.

**2.5 Training**

Some voice recognizers can be adjusted to a speaker. At the point when the system has this capacity, it may permit training to occur. An ASR framework is prepared by having the speaker rehash standard or normal expressions and changing its correlation calculations to match that specific speaker. Training a recognizer for the most part enhances its precision.

Training can likewise be utilized by speakers that experience issues talking, or affirming certain words. The length of the speaker can reliably rehash an articulation, ASR frameworks with training ought to have the capacity to adjust.

# 3.   SPEECH ANALYZER

Speech examination, likewise alluded to as front-end investigation or feature extraction, is the initial phase in a programmed voice recognition framework. This methodology plans to concentrate acoustic highlights from the discourse waveform. The yield of front-end investigation is a conservative, effective arrangement of parameters that speak to the acoustic properties saw from information discourse signals, for ensuing usage by acoustic demonstrating.

There are three noteworthy sorts of front-end handling systems, specifically Mel-Frequency Cepstral Coefficients (MFCC), Linear Predictive Coding (LPC), and Perceptual Linear Prediction (PLP), where MFCC and PLP are most ordinarily utilized as a part of cutting edge ASR frameworks.

## 3.1 Linear predictive coding

LPC begins with the presumption that a voice signal is delivered by a signal toward the end of a tube (voiced sounds), with periodic included murmuring and popping sounds. Albeit obviously rough, this model is really a nearby close estimation to the truth of speech creation. The glottis (the space between the vocal strings) delivers the buzz, which is described by its power (clamor) and recurrence (pitch). The vocal tract (the throat and mouth) shapes the tube, which is described by its resonances, which are called formants. Murmurs and pops are created by the activity of the tongue, lips and throat amid sibilants and plosives.

LPC investigates the speech motion by assessing the formants, expelling their belongings from the speech flag, and evaluating the force and recurrence of the remaining buzz. The procedure of uprooting the formants is called opposite separating, and the staying flag after the subtraction of the separated demonstrated sign is known as the deposit.

The numbers which portray the power and recurrence of the buzz, the formants, and the deposit sign, can be put away or transmitted someplace else. LPC blends the speech motion by turning

around the methodology: utilize the buzz parameters and the deposit to make a source sign, utilize the formants to make a channel (which speaks to the tube), and run the source through the channel, bringing about speech.

Since speech signs change with time, this procedure is done on short pieces of the speech signal, which are called edges; by and large 30 to 50 casings every second give clear speech with great pressure.

## 3.2 Mel Frequency Cepstrum Coefficients

These are gotten from a sort of cepstral representation of the sound clasp (a "spectrum of a spectrum"). The distinction between the cepstrum and the Mel-frequency cepstrum is that in the MFC, the recurrence groups are situated logarithmically (on the mel scale) which approximates the human sound-related framework's reaction more nearly than the straightly dispersed recurrence groups got straightforwardly from the FFT or DCT. This can take into consideration better handling of information, for instance, in sound pressure. Then again, dissimilar to the sonogram, MFCCs do not have an external ear model and, consequently, can't speak to saw din precisely.

MFCCs are usually determined as takes after:

First take the Fourier transform of a windowed passage of a voice signal. Then mapping the log amplitudes of the spectrum acquired above onto the Mel scale, utilizing triangular covering windows. After taking the Discrete Cosine Transform of the rundown of Mel log-amplitudes, as though it were a signal, the MFCCs are found, i.e., it is the amplitudes of the subsequent spectrum.

## 3.3 Perceptual Linear Prediction

Perceptual linear prediction, just like LPC analysis, is in light of the fleeting range of speech. As opposed to unadulterated linear predictive investigation of speech, perceptual linear prediction (PLP) alters the transient range of the speech by a few psychophysically based

changes. This procedure utilizes three ideas from the psychophysics of hearing to determine an appraisal of the sound-related range: Firstly, the discriminating band spectral determination is taken. At that point we can locate the equivalent loudness bend, and the power uproar power law.

The sound-related range is then approximated by an autoregressive all-post model. In correlation with ordinary linear predictive (LP) examination, PLP investigation is steadier with human's hearing.

# 4.    PREPROCESSING

Firstly the info speech signal is pre-emphasized to misleadingly intensify the high frequencies. Voice signals are non-stationary in nature, intending to say the exchange capacity of the vocal tract; which produces it, changes with time, however it changes slowly. It is safe to expect that it is piecewise stationary. Subsequently the speech signal is confined with 30ms to 32ms casings with a cover of 20ms. The requirement for the cover is that data may be lost at the casing limits, so outline limits need to be inside another edge. Numerically, confining is identical to reproducing the signal with a progression of sliding rectangular windows. The issue with rectangular windows is that the force contained in the side flaps is essentially high and thusly may offer ascent to otherworldly spillage. So as to dodge this we utilize a Hamming window given by:

$$w(n) = 0.54 - 0.46\cos\left(\frac{2\pi n}{N-1}\right), \quad 0 \le n \le N-1$$

# 5. MEL FREQUENCY CEPSTRAL COEFFICIENTS

MFCCs are ordinarily processed by utilizing a bank of triangular-formed channels, with the middle recurrence of the channel separated directly for frequencies under 1000 Hz and logarithmically over 1000 Hz. The transmission capacity of every channel is controlled by the inside frequencies of the two nearby channels and is reliant on the recurrence scope of the channel bank and number of channels picked for configuration. In any case, for the human sound-related framework it is assessed that the channels have a transmission capacity that is identified with the inside recurrence of the channel. Further it has been demonstrated that there is no confirmation of two areas (straight and logarithmic) in the tentatively decided Mel recurrence scale.

## 5.1 Silence detection

This is an essential venture in any front end speaker acknowledgment framework. At the point when the client says out any word the framework will never have control over the moment the word is articulated. It is consequently that we require a quiet recognition step which essentially decides when the client has really begun articulating the word any along these lines deciphers the casing of reference to that moment.

## 5.2 Windowing

The following venture in the handling is to window every individual edge to minimize the signal discontinuities toward the starting and end of the casing. The idea here is to minimize the ghostly twisting by utilizing the window to decrease the signal to zero toward the starting and end of the casing. In the event that we characterize the window as,

$$w(n), 0 \leq n \leq N-1.$$

Where N is the quantity of tests in every casing, then the aftereffect of windowing is the signal

We are using Hamming window in this case, whose equation is:

$$w(n) = 0.54 - 0.46 \cos\left(\frac{2\pi n}{N-1}\right) \quad 0 \leq n \leq N-1$$

## 5.3 Fast Fourier Transform

The following transforming step is the Fast Fourier Transform, which changes over every casing of N tests from the time area into the recurrence space. The FFT is a quick

calculation to execute the Discrete Fourier Transform (DFT) which is characterized on the situated of N tests {Xn}, as takes after:

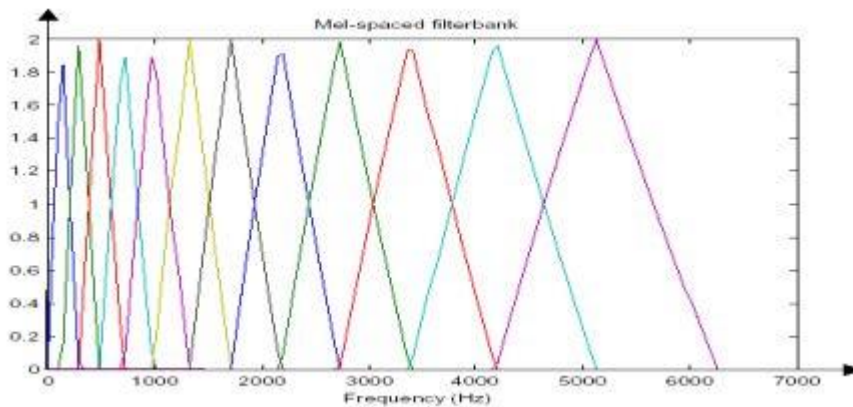$$X_n = \sum_{k=0}^{N-1} x_k e^{-2\pi j k n / N}, \qquad n = 0,1,2,\dots, N-1$$

Note that we utilize j here to mean the fanciful unit. When all is said in done Xn's are complex numbers. The subsequent grouping {Xn} is translated as takes after: the zero recurrence compares to n = 0, positive frequencies: $0 < f < Fs/2$ compare to values $1 <= n <= N/2-1$ while negative frequencies $-Fs/2 < f < 0$ relate to $N/2+1 <= n <= N-1$. Here, Fs indicates the testing recurrence. The outcome got after this step is regularly alluded to as signal's periodogram or spectrum.
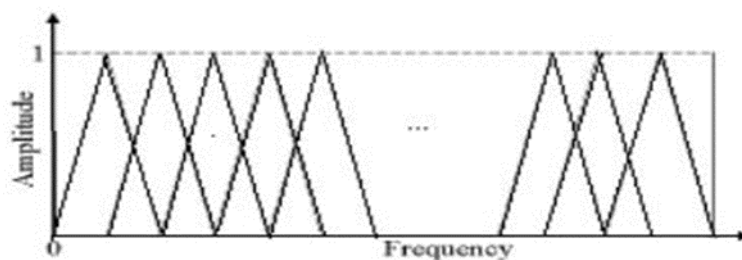
## 5.4 Mel-frequency Warping

As specified above, psychophysical studies have demonstrated that human view of the recurrence substance of sounds for speech signals does not take after a straight scale. Accordingly for every tone with a genuine recurrence, f, measured in Hz, a subjective pitch is measured on a scale called the "mel" scale. The mel-recurrence scale is a straight recurrence dispersing beneath 1000 Hz and a logarithmic dividing over 1000 Hz. As a kind of perspective point, the pitch of a 1 kHz tone, 40 dB over the perceptual listening to edge, is characterized as 1000 mels. Along these lines we can utilize the accompanying estimated equation to process the mels for a given recurrence f in Hz:

$$mel(f) = 2595 * \log_{10}(1 + f/700)$$

One way to deal with reenacting the subjective range is to utilize a channel bank, one channel for every sought mel-recurrence part. That channel bank has a triangular bandpass recurrence reaction, and the separating and additionally the data transfer capacity is controlled by a consistent mel-recurrence interim. The changed range of S (w) in this manner comprises of the yield force of these channels when S (w) is the info. The quantity of mel cepstral coefficients, K, is regularly picked as 20. Note that this channel bank is connected in the recurrence space; consequently it basically adds up to taking those triangle-shape windows in the Figure on the range. A helpful state of mind about this mel-twisted channel bank is to view every channel as a histogram container (where canisters have cover) in the recurrence space. A valuable and effective method for executing this is to consider these triangular channels in the Mel scale where they would as a result be similarly dispersed channels. In linear frequency scale, figure:

Mel-spaced filterbank

In Mel frequency scale, the filter bank would be:



## 5.5 Cepstrum

In the last step, we change over the log mel range back to time. The outcome is known as the mel recurrence cepstral coefficients (MFCC). The cepstral representation of the speech range gives a decent representation of the neighborhood unearthly properties of the signal for the given casing investigation. Since the mel range coefficients (thus their logarithm) are genuine numbers, we can change over them to the time space utilizing the Discrete Cosine Transform (DCT). Accordingly on the off chance that we indicate those mel power range coefficients that are the consequence of the last step are Sk , where k =1, 2,… ..,K. Figure the MFCC's, Cn as,
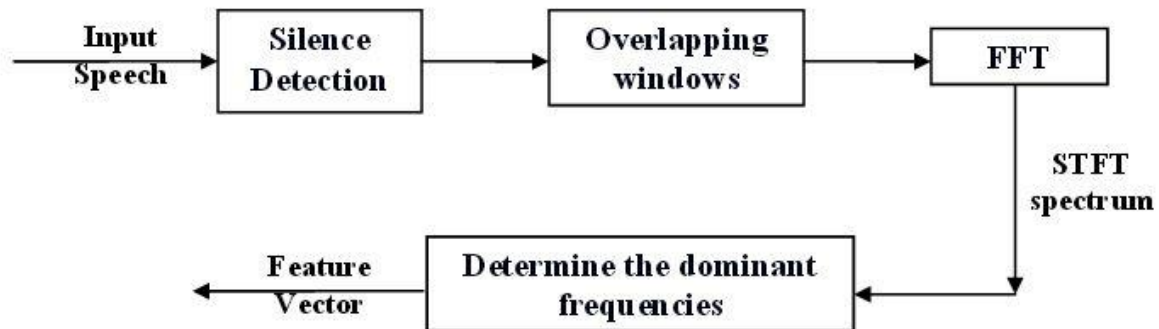
$$\tilde{c}_n = \sum_{k=1}^{K} (\log \tilde{S}_k) \cos\left[ n\left(k - \frac{1}{2}\right)\frac{\pi}{K}\right], \qquad n = 1, 2, ..., K$$

Note that we bar the first segment c0 from the DCT since it speaks to the mean estimation of the information signal which conveyed little speaker particular data. When we have the MFCCs, these portray the specific word so amid the preparation stage the coefficients for the word are dead set and put away. Amid the acknowledgment stage, the coefficients are again decided for the articulated word and acknowledgment is completed by breaking down the MSE

concerning the put away coefficients and characterizing a proper edge relying upon the level of security needed for the application.

**5.6 Approach**

A square outline of the structure of a MFCC processor is as demonstrated in Fig.



The speech information is ordinarily recorded at an examining rate over 10000 Hz. This testing recurrence was decided to minimize the impacts of associating in the simple to-advanced change. These tested signals can catch all frequencies up to 5 kHz, which cover most vitality of sounds that are produced by people. The fundamental motivation behind the MFCC processor is to copy the conduct of the human ears. Also, instead of the speech waveforms themselves, MFCC's are demonstrated to be less defenceless to said varieties.
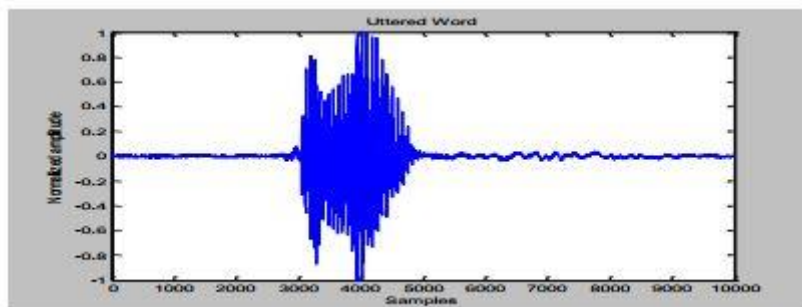
We initially put away the speech signal as a 10000 example vector. It was seen from our trial that the genuine expressed speech dispensing with the static bits came up to around 2500 examples, thus, by utilizing a straightforward edge system we completed the hush discovery to concentrate the real articulated speech.

It is clear that what we needed to accomplish was a voice based biometric framework fit for perceiving disengaged words. As our investigations uncovered all the segregated words were expressed inside 2500 examples. At the same time, when we passed this speech signal through a MFCC processor, it spilt this up in the time space by utilizing covering windows each with around 250 examples. Hence when we change over this into the recurrence space we simply have around 250 range values under every window. This inferred that changing over it to the Mel scale would be excess as the Mel scale is straight till 1000 Hz. Along these lines, we dispensed with the square which did the Mel distorting. We specifically utilized the covering triangular windows as a part of the recurrence space. We acquired the vitality inside every
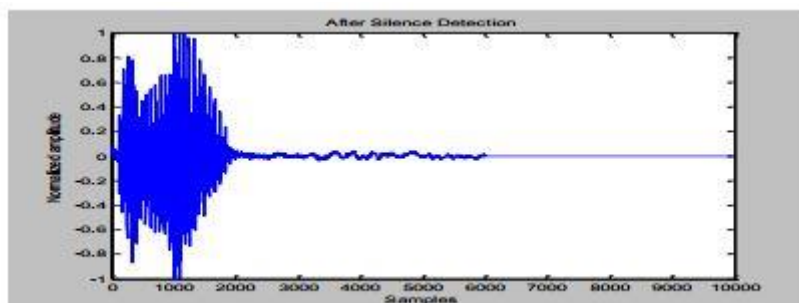
triangular window, trailed by the DCT of their logarithms to attain to great compaction inside a little number of coefficients as depicted by the MFCC approach.

This calculation be that as it may, has a downside. As disclosed before the way to this methodology is utilizing the energies inside every triangular window, notwithstanding, this may not be the best approach as was found. It was seen from the trials that on account of the conspicuousness given to vitality, this methodology neglected to perceive the same word expressed with distinctive vitality. Additionally, as this takes the summation of the vitality inside every triangular window it would basically give the same estimation of vitality independent of whether the range crests at one specific recurrence and tumbles to lower values around it or whether it has an equivalent spread inside the window. This is the reason we chose not to continue with the execution of the MFCC approach.
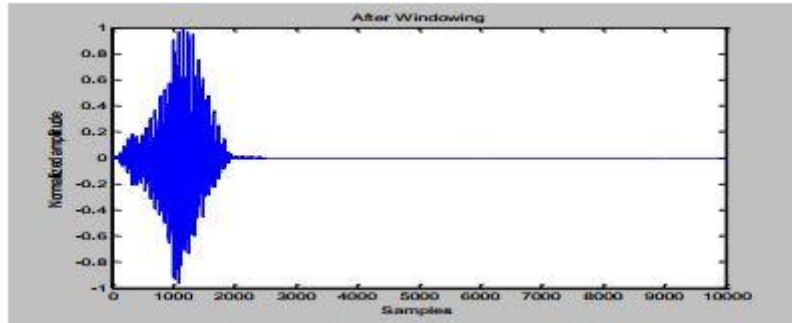
The reproduction was completed in MATLAB. The different phases of the reproduction have been spoken to as the plots indicated. The information ceaseless speech signal considered as an illustration for this task is "Hi".
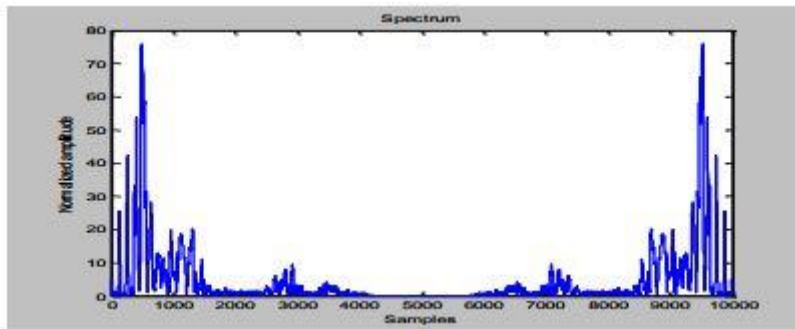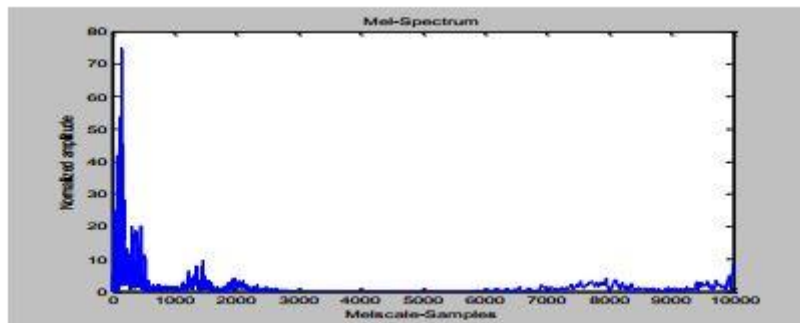


Uttered word



After silence detection

After windowing



FFT Spectrum



Mel-frequency Spectrum

# TIME DOMAIN ANALYSIS

Despite the fact that discourse is non-stationary in a genuine sense, study has demonstrated that it can be thought to be semi stationary i.e. gradually time fluctuating. At the point when analysed over an adequately brief time of time (somewhere around 20 and 40 msec), its attributes are genuinely stationary. The ramifications of this is that when we consider such little portions of discourse they have a predominant recurrence inside the windowed fragment. Hence, this can be handled through a brief while recurrence examination. The figure underneath shows how we can utilize covering time area windows to cater for the brief while Fourier investigation.

The short-time Fourier transform (STFT) of a speech signal s(t) is given by

$$S(f,t) = \int_{-\infty}^{\infty} s(\tau)w(t-\tau)e^{-j2\pi f\tau}d\tau$$

Where, w(t) is a window capacity of span Tw. In discourse transforming, the Hamming window capacity is ordinarily utilized and its width Tw is regularly 20 - 40 msec. We can deteriorate S(f, t) as takes after:

$$S(f,t) = |S(f,t)|e^{j\psi(f,t)}$$

where |S(f, t)| is the brief while size range and ψ(f, t) = S(f, t) is the brief while stage range. Our undertaking spotlights on the brief while size range as connected to the prevailing recurrence inside every covering portion.

In the STFT-based framework, there are various configuration issues that must be tended to. To start with, what kind of window capacity w(t) ought to be utilized for figuring the STFT. A decreased window capacity (Hamming is normally favored in discourse applications) has been utilized as a part of all studies on the grounds that with the slow move off of such a window, the impact of cover is kept under check to verify that every window offers unmistakable quality to a free fragment of discourse while keeping up a smooth move starting with one section of discourse then onto the next. Second, what ought to be the term tw. The modifier „short-time" suggests a limited time window over which the properties of discourse may be accepted stationary; it doesn't allude to the real term of the window. Nonetheless, it is realized that

discourse fundamentally includes phonemes, and tests directed coordinated with our deductions that Tw could be between 20 msec to 50 msec.

**Time Domain Approach**

The embodiment of STFT as disclosed is to concentrate the recurrence range inside every window in the time area to investigate the commanding frequencies inside every window comparing to different phonemes. In our venture, rather than specifically applying the STFT comparison, we have utilized a hamming window to concentrate out portions of the discourse and inside each of these fragments we utilized FFT to get the wanted recurrence range.

The plots in Fig demonstrate a case of how covering time area hamming windows could be spoken to. The sample considers hamming windows of width 500 examples which is around 50 ms and a cover of 100 specimens which would associate with 10 msec. This has been changed in our

# MATLAB CODES

**MFCC APPROACH**

**TRAINING PHASE**

```
fsn = 10000;

tn= hamming(4000);

wn = [tn ; zeros(6000,1)];

fn = (1:10000);

mel(fn) = 2595 * log(1 + f n/ 700);

trin = triang(100);

win1n = [trin ; zeros(9900,1)];

win2n = [zeros(50,1) ; trin ; zeros(9850,1)];

win3n = [zeros(100,1) ; trin ; zeros(9800,1)];

win4n = [zeros(150,1) ; trin ; zeros(9750,1)];

win5n = [zeros(200,1) ; trin ; zeros(9700,1)];

win6n = [zeros(250,1) ; trin ; zeros(9650,1)];

win7n = [zeros(300,1) ; trin ; zeros(9600,1)];

win8n = [zeros(350,1) ; trin ; zeros(9550,1)];

win9n = [zeros(400,1) ; trin ; zeros(9500,1)];

win10n = [zeros(450,1) ; trin ; zeros(9450,1)];

win11n = [zeros(500,1) ; trin ; zeros(9400,1)];

win12n = [zeros(550,1) ; trin ; zeros(9350,1)];

win13n = [zeros(600,1) ; trin ; zeros(9300,1)];

win14n = [zeros(650,1) ; trin ; zeros(9250,1)];

win15n = [zeros(700,1) ; trin ; zeros(9200,1)];
```

```
win16n = [zeros(750,1) ; trin ; zeros(9150,1)];

win17n = [zeros(800,1) ; trin ; zeros(9100,1)];

win18n = [zeros(850,1) ; trin ; zeros(9050,1)];

win19n = [zeros(900,1) ; trin ; zeros(9000,1)];

win20n = [zeros(950,1) ; trin ; zeros(8950,1)];

xn = wavrecord(1 * fsn, fsn, 'double');

plot(xn);

wavplay(xn);

in = 1;

while abs(x(in)) < 0.05

in = in + 1;

end

xn(1 : in) = [];

xn(6000 : 10000) = 0;

x1n = xn. * w;

mxn = fft(x1n);

nxn = abs(mxn(floor(mel(fn))));

nxn = nxn. / max(nxn);

nx1n = nxn. * win1n;

nx2n = nxn. * win2n;

nx3n = nxn. * win3n;

nx4n = nxn. * win4n;

nx5n = nxn. * win5n;

nx6n = nxn. * win6n;
```

```
nx7n = nxn. * win7n;

nx8n = nxn. * win8n;

nx9n = nxn. * win9n;

nx10n = nxn. * win10n;

nx11n = nxn. * win11n;

nx12n = nxn. *win12n;

nx13n = nxn. * win13n;

nx14n = nxn. * win14n;

nx15n = nxn. * win15n;

nx16n = nxn. * win16n;

nx17n = nxn. * win17n;

nx18n = nxn. * win18n;

nx19n = nxn. * win19n;

nx20n = nxn. * win20n;

sxn1 = sum(nx1n. ^ 2);

sxn2 = sum(nx2n. ^ 2);

sxn3 = sum(nx3n. ^ 2);

sxn4 = sum(nx4n. ^ 2);

sxn5 = sum(nx5n. ^ 2);

sxn6 = sum(nx6n. ^ 2);

sxn7 = sum(nx7n. ^ 2);

sxn8 = sum(nx8n. ^ 2);

sxn9 = sum(nx9n. ^ 2);

sxn10 = sum(nx10n. ^ 2);
```

```
sxn11 = sum(nx11n. ^ 2);

sxn12 = sum(nx12n. ^ 2);

sxn13 = sum(nx13n. ^ 2);

sxn14 = sum(nx14n. ^ 2);

sxn15 = sum(nx15n. ^ 2);

sxn16 = sum(nx16n. ^ 2);

sxn17 = sum(nx17n. ^ 2);

sxn18 = sum(nx18n. ^ 2);

sxn19 = sum(nx19n. ^ 2);

sxn20 = sum(nx20n. ^ 2);

sxn = [sxn1, sxn2, sxn3, sxn4, sxn5, sxn6, sxn7, sxn8, sxn9, sxn10, sxn11, sxn12, sxn13, sxn14,

sxn15, sxn16, sxn17, sxn18, sxn19, sxn20];

sxn = log(sxn);

dxn = dct(sxn);

fidn = fopen('hello.dat', 'w');

fwrite(fidn, dxn, 'real*8');

fclose(fidn);
```

**TESTING PHASE**

```
fsn = 10000;

tn = hamming(4000);

wn = [tn ; zeros(6000,1)];

fn = (1:10000);

mel(fn) = 2595 * log(1 + fn / 700);

trin = triang(100);

win1n = [trin ; zeros(9900,1)];

win2n = [zeros(50,1) ; trin ; zeros(9850,1)];

win3n = [zeros(100,1) ; trin ; zeros(9800,1)];

win4n = [zeros(150,1) ; trin ; zeros(9750,1)];

win5n = [zeros(200,1) ; trin ; zeros(9700,1)];

win6n = [zeros(250,1) ; trin ; zeros(9650,1)];

win7n = [zeros(300,1) ; trin ; zeros(9600,1)];

win8n = [zeros(350,1) ; trin ; zeros(9550,1)];

win9n = [zeros(400,1) ; trin ; zeros(9500,1)];

win10n = [zeros(450,1) ; trin ; zeros(9450,1)];

win11n = [zeros(500,1) ; trin ; zeros(9400,1)];

win12n = [zeros(550,1) ; trin ; zeros(9350,1)];

win13n = [zeros(600,1) ; trin ; zeros(9300,1)];

win14n = [zeros(650,1) ; trin ; zeros(9250,1)];

win15n = [zeros(700,1) ; trin ; zeros(9200,1)];

win16n = [zeros(750,1) ; trin ; zeros(9150,1)];

win17n = [zeros(800,1) ; trin ; zeros(9100,1)];

win18n = [zeros(850,1) ; trin ; zeros(9050,1)];

win19n = [zeros(900,1) ; trin ; zeros(9000,1)];

win20n = [zeros(950,1) ; trin ; zeros(8950,1)];

yn = wavrecord(1 * fsn, fsn, 'double');

in = 1;

while abs(y(in)) < 0.05 \
```

```matlab
in = in + 1;
end
yn(1 : in) = [];
yn(6000 : 10000) = 0;
yn1 = yn. * wn;
myn = fft(yn1);
nyn = abs(my(floor(mel(fn))));
nyn = nyn. / max(nyn);
nyn1 = nyn. * win1n;
nyn2 = nyn. * win2n;
nyn3 = nyn. * win3n;
nyn4 = nyn. * win4n;
nyn5 = nyn. * win5n;
nyn6 = nyn. * win6n;
nyn7 = nyn. * win7n;
nyn8 = nyn. * win8n;
nyn9 = nyn. * win9n;
nyn10 = nyn. * win10n;
nyn11 = nyn. * win11n;
nyn12 = nyn. * win12n;
nyn13 = nyn. * win13n;
nyn14 = nyn. * win14n;
nyn15 = nyn. * win15n;
nyn16 = nyn. * win16n;
nyn17 = nyn. * win17n;
nyn18 = nyn. * win18n;
nyn19 = nyn. * win19n;
nyn20 = nyn. * win20n;
syn1 = sum(nyn1. ^ 2);
syn2 = sum(nyn2. ^ 2);
```

```
syn3 = sum(nyn3. ^ 2);

syn4 = sum(nyn4. ^ 2);

syn5 = sum(nyn5. ^ 2);

syn6 = sum(nyn6. ^ 2);

syn7 = sum(nyn7. ^ 2);

syn8 = sum(nyn8. ^ 2);

syn9 = sum(nyn9. ^ 2);

syn10 = sum(nyn10. ^ 2);

syn11 = sum(nyn11. ^ 2);

syn12 = sum(nyn12. ^ 2);

syn13 = sum(nyn13. ^ 2);

syn14 = sum(nyn14. ^ 2);

syn15 = sum(nyn15. ^ 2);

syn16 = sum(nyn16. ^ 2);

syn17 = sum(nyn17. ^ 2);

syn18 = sum(nyn18. ^ 2);

syn19 = sum(nyn19. ^ 2);

syn20 = sum(nyn20. ^ 2);

syn = [syn1, syn2, syn3, syn4, syn5, syn6, syn7, syn8, syn9, syn10, syn11, syn12, syn13, syn14, syn15, syn16, syn17, syn18, syn19, syn20];

syn = log(syn);

dyn = dct(syn);

fidn = fopen('hello.dat','r');

dxn = fread(fidn, 20, 'real*8');

fclose(fidn);

dxn = dxn.';

MSE=(sum((dxn - dyn). ^ 2)) / 20;

if MSE<1

fprintf('\n\n It's Ok\n\n');

Ok=wavread('Ok.wav');
```

```
wavplay(Ok);
else
fprintf('\n\n Not Ok\n\n');
Not_Ok=wavread('Not_Ok.wav');
wavplay(Not_Ok);
end
```

# CONCLUSION

We effectively recreated the MFCC methodology and the Time domain methodology utilizing MATLAB. We finished up from our tests that for a disengaged word acknowledgment framework like the one we meant to execute in our venture, the time domain methodology ended up being more powerful. The explanation behind this is that, the MFCC methodology has a downside. As disclosed before the way to this methodology is utilizing the energies, notwithstanding, this may not be the best approach as was found.

It was seen from the investigations that due to the noticeable quality given to vitality, this methodology neglected to perceive the same word articulated with distinctive vitality. Likewise, as this takes the summation of the vitality inside every triangular window it would basically give the same estimation of vitality independent of whether the range crests at one specific recurrence and tumbles to lower values around it or whether it has an equivalent spread inside the window. Then again, as the time domain methodology is not taking into account vitality but rather absolutely in light of the overwhelming frequencies inside little portions of discourse, it makes great utilization of the semi – stationary property of discourse. We hence closed from our recreation comes about that the time domain methodology is more suitable for the disconnected work recognizer needed for a voice based biometric framework. An examination between the calculations we recreated is given in the tables underneath.

# APPLICATIONS

After almost sixty years of exploration, discourse acknowledgment innovation has come to a generally abnormal state. Then again, most best in class ASR frameworks run on desktop with intense chip, abundant memory and an ever-present power supply. In these years, with the quick evolution of equipment and programming advancements, ASR has gotten to be more convenient as an option human-to-machine interface that is required for the accompanying application regions:

1. Stand-alone customer gadgets, for example, wrist watch, toys and sans hands cellular telephone in auto where individuals are not able to utilize different interfaces or huge information stages like consoles are not accessible.

2. Single reason order and control framework, for example, voice dialing for cell, home, and office telephones where multi-capacity (PCs) are repetitive.

A portion of the utilizations of speaker check frameworks are:

Time and Attendance Systems

Access Control Systems

Phone Banking/Broking

Biometric Login to phone helped shopping frameworks

Data and Reservation Services

Security control for secret data

Criminological purposes

Voice based Telephone dialing is one of the applications we mimicked. The key center of this application is to help the physically tested in executing an unremarkable undertaking like phone dialing. Here the client at first prepares the framework by expressing the digits from 0 to 9. Once the framework has been prepared, the framework can perceive the digits expressed by the client who prepared the framework. This framework can likewise include some natural security as the framework taking into account cepstral methodology is speaker subordinate. The calculation is run on a specific speaker and the MFCC coefficients decided. Presently the calculation is connected to an alternate speaker and the bungle was unmistakably watched. In this manner the natural security gave by the framework was affirmed.

Instantly frameworks have likewise been planned which fuse Speech and Speaker Recognition. Ordinarily a client has two levels of check. He/She needs to at first talk the right watchword to obtain entrance to a framework. The framework not just checks if the right secret key has been said additionally centered around the genuineness of the speaker. A definitive objective is do have a framework which does a Speech, Iris, Fingerprint Recognition to actualize access control.

# REFERENCES

[1] Lawrence Rabiner, Biing-Hwang Juang – 'Fundamentals of Speech Recognition'

[2] Wei Han, Cheong-Fat Chan, Chiu-Sing Choy and Kong-Pang Pun – „An Efficient MFCC Extraction Method in Speech Recognition', Department of Electronic Engineering, The Chinese University of Hong Kong, Hong, IEEE – ISCAS, 2006

[3] Leigh D. Alsteris and Kuldip K. Paliwal – „ASR on Speech Reconstructed from Short- time Fourier Phase Spectra', School of Microelectronic Engineering Griffth University, Brisbane, Australia, ICLSP - 2004

[4] Waleed H. Abdulla – „Auditory Based Feature Vectors for Speech Recognition Systems', Electrical & Electronic Engineering Department, The University of Auckland

[5] Pradeep Kumar P and Preeti Rao – „A Study of Frequency-Scale Warping for Speaker Recognition', Dept of Electrical Engineering, IIT- Bombay, National Conference on Communications, NCC 2004, IISc Bangalore, Jan 30 -Feb 1, 2004

[6] Beth Logan – „Mel Frequency Cepstral Coefficients for Music Modeling', Cambridge Research Laboratory, Compaq Computer Corporation

[7] MIT Open Courseware