



# Speech To Sign Translator

<sup>1</sup>Nighila Ashok, <sup>2</sup>Aswin K.S., <sup>3</sup>Nikhil Babu P., <sup>4</sup>Sajith Suresh, <sup>5</sup>Sharun K. Suresh

<sup>1</sup>Assistant Professor, Department of Computer Science and Engineering, Universal Engineering College, Vallivattom, Thrissur, Kerala, India

<sup>2,3,4,5</sup>B. Tech Students, Dept. of CSE, Universal Engineering College, Vallivattom, Thrissur, Kerala, India

**Abstract:** Sign language is a visual language that is mainly used by deaf people and the Hard-of-Hearing community as their primary language. Unlike acoustically expressed sound patterns, sign language uses hand signals, gestures, facial expressions, and body language. Sign language may also help the people with disabilities, namely Autism, Cerebral Palsy, Apraxia of speech and Down Syndrome.

The objective of our project is to reduce the communication gap between normal people and disabled people by providing a virtual animation translator for sign language. Speech to sign translator is a user-friendly web application that converts speech/text to sign language animation and also reduces the amount of effort spent on communication. The system consists of both speech and text conversion using natural language processing (NLP) and machine learning techniques. Input speech is captured through a microphone then it is translated to text through Speech recognition using Mel-frequency cepstral coefficients (MFCC) features. The recognized text preprocessed using NLTK and resultant words looks for a word match in ASL dataset. The database contains a certain number of pre-recorded video animation signs using 3D avatar where mainly there is one video clip per each basic word. If a match occurred for a word, then the equivalent ASL translation will be displayed corresponding to that word. Otherwise, the word will be fingerspelled.

**Index Terms -** Speech recognition, MFCC, Text preprocessing, NLP, ASL.

## I. INTRODUCTION

The most recent studies & research shows that about 5% of the total population are facing hearing problems. Communication via gestures is routinely the essential methodology for correspondence utilized by individuals with hearing-inabilities. Gesture based communication comprises manual communication features including hand shapes, gestures, motion & orientations. Dissimilar to typical language, sign languages don't have obvious design or syntax, hence there is no or exceptionally less worthiness of these signs outside the little universe of these contrastingly abled networks. Research on ASL exhibits that gesture-based communication is an undeniable language with its own syntax, grammar and linguistic attributes.

Nowadays, there is the requirement of a middle person as a go-between for exchanging information. Regardless, the expert interpreter may not commonly be available and may cause a communication limit and deferment in the range of information. There is an enormous heap of exploration and working environments made for Sign language confirmation. But, the Speech to sign translators that assists the hearing disabled people are considered to be very rare. In this unprecedented circumstance, Speech to sign based communication translators goes likely as a help to hearing debilitated humankind by giving a preferable associate over communication with the normal ones.

To assist individuals with such incapacities a great deal of exploration has been led and a few arrangements have been made worldwide up until this point however no significant achievement has been accounted for until now. Sign language communication is crucial as the conference disabled people groups are stressed over their energetic, social, and phonetic turn of events. The principal language for the hearing-impaired people groups is communication of sign language correspondence which proceeds bilingually with the preparation of public correspondence through motions similarly as the public made or conveyed in language. Figuring out a viable method for investigating and making is exceptionally going after for by a wide margin most with hearing disabled. Hearing adversity compels the youth's mentoring, high-level training and impacts future master possibilities. For those that can examine and form, getting a handle on the setting of what is being spoken becomes risky, particularly in conditions where nonverbal sounds or activities are involved. The explanation for such a low proficiency rate can be both of the accompanying Lack of Sign Language translators.

As communication through signing interchanges don't have particular development or sentence structure, subsequently there is no or incredibly less amiability of these signs outside the little universe of these contrastingly abled people. Research on American Sign Language confirms that disclosing through stamping is a certain language with its complement, its sentence structure, and other phonetic characteristics. To approve something practically identical for various changes through movements, there are several undertakings including Indian Sign Language. Correspondence for the gathering run down people likewise puts like rail line stations, transport stands, banks, emergency clinics, and so on, is really difficult in light of the fact that a gathering individual may not get a handle on the signal-based correspondence utilized by the meeting weakened convey. Moreover, a discussion individual can't give

any message to a gathering prevented person as he/she probably won't have even the farthest sign about the correspondence through checking. To make the correlation between hearing debilitated individuals and non-hard hearing territory, language interpretation is a must. To help them with thinking better with the rest of the world, a system is required which will engage the change of text to Sign Language as well as the opposite way around. These frameworks will expand the nature of living in this community.

## II. RELATED WORKS

Here we discuss the various studies based on Sign language translation.

The paper [1] discusses Speech-to-Sign Language (BSL) translation development that expects to help exchanges between a hearing handicapped individual and a specialist in a Post Office by translating the assistant's speech to communication through hand signals. A speech recognizer captures speech from the agent and the system then mixes the fitting progression of signs in British Sign language (BSL) using a remarkably developed image. The system used the Entropic talk recognizer. The recognizer requires a lot of acoustic models for matching the data talk signal and an association that coordinates the chase of the recognizer during affirmation. An anticipated language structure approach is used by the system. The understanding is done using an expression inquiry informational index. Regardless, on the grounds that there are several sentences to use as configurations, the conversation between the individuals is limited.

The author [2] proposed the programmed speech recognition system model the association between acoustic speech-sign and phone classes in two stages, explicitly, extraction of spooky set up features in view of prior data followed by means of arrangement of acoustic models, regularly an Artificial Neural Network (ANN). It was shown that the Convolutional Neural Networks (CNNs) can show phone classes from rough acoustic audio signals, showing up at execution that is comparable to other existing component-based approaches. The paper loosens up the CNN-based method for managing complex language speech recognition tasks. Even more conclusively, the proposed procedure investigates the CNN-based approach against the standard ANN-set up system concerning Wall Street Journal corpus. The assessments show that the CNN-based system achieves favored execution over the normal ANN based approach with various limits. We in like manner show that the features acquired from rough talk by the CNN based technique could summarize across different informational collections.

The paper [3] is about a model sign association application. It translates the input text into corresponding American Sign Language (ASL). Sign mix and talk mix perform essentially a comparative endeavor. The primary differentiation is the outcome. Henceforth the plans of both of these are also for all intents and purposes something almost identical. The system uses Perl scripts through the ordinary entry interface for performing promoting exercises. It has three essential association focuses. The primary association point, MENU CGI offers menus for signs by which the users can decide the phonological limits. Additional menus help the users who know nothing about ASCII-code. They can clearly pick the hand shape, hand position, and hand signs for each hold. Second association point, ASCII-based downsized parser is for additional created users to type with increments for timing and non-manuals. The fingerspelling module helps the user with forming in the Roman letter set. This section yields an ASCII-based tree which transforms into a commitment to the change module. The change module further conveys Web3D turns for joints that are used. After the development of turns, they become the commitment for the Sign-generating module. The Sign Generating module integrates them with Web3D avatar for making all out web archives with action data. Then, with the help of a module, the liveliness is played.

In the paper [4], the framework proposes a unique method for managing modified Sign Language generation including late enhancements in Neural Machine Translation, Generative Adversarial Networks and the development age. The system is prepared for conveying sign accounts imparted in language sentences. Contrary to current techniques that are likely to strongly explain data, this philosophy requires unimportant shimmer and skeletal level remarks for planning. This is achieved by isolating the endeavor into committed sub-processes. In first case, the system translates the imparted language sentences into sign stance progressions by joining a NMT network with a Motion Graph. The resulting stance information is then used to build a generative model that produces photo functional correspondence by means of signals and video groupings. The understanding the deficiencies of the system were evaluated using some Sign Language Translation dataset. The structure further shows the video age capacities with regards to both multi-endorser and top-quality settings abstractly and quantitatively using broadcast quality assessment estimations.

In the work [5], it is wanted to get inputs from an alternate format. The data sources can be of structures: Live audio input, Text input, or Recorded sound file input. The live speech is taken as a commitment from the mouthpiece(mic) of the system. This is done using the PyAudio, which is a Python pack that is used to record audio on a collection of stages. The speech is converted into text using Google Speech Recognizer API. An API helps with converting audio over to text by combining neural network models. In the data association of giving the sound record, the input speech is changed over into text by using this Google Speech Recognizer. For long sound records, the sound is segregated into additional divisions in view of the occasion of pause. The speech segments are then passed into the Google Speech Recognizer to change over into text capably. As a general rule is taken for playing the video progression in the Sign Language mediator. Time taken for the translator to convert from speech to communication by means of hand sign is noted from the request line. Speech recognition takes time dependent upon the length of the sentence that is delivered by the speaker. The message assessment takes an irrelevant proportion of time in changing over that message sentence.

In the paper [6], the cell phone application is proposed and it helps the hearing disabled people to just learn unique dialects abuse their comfortable language severally, client can learn in on the web or in disconnected modes, furthermore makes client to interact and moreover gives more significant level of convenience, improves the understanding capabilities. This can do the constant interpretation from various Indian dialects and English to communication through signing will assist with overcoming this issue generally and can work even without the web. The client signed in to the application, then Sign Translator will translate the word entered by the client, which will change over it into marking double-dealing data. The application comprises a Sign Translator interface, a Mobile Unit, Text Reader, Interpreter, Data Storage, and Maintenance as a database. Sign Translator will perceive the input speech, it coordinates the voice with string, and the suitable picture related with string will be coordinated and subsequently will be meant Sign Language where then, at that point, typical individuals might convey their considerations and denied individuals might envision the result and can undoubtedly comprehend in their agreeable Sign Language. Screen captures of different parts of Sign Translator results are incorporated as proof.

In this paper [7], the design proposed in this paper copies the handling technique for individuals with the hearing handicaps in eating up overpowering press outlets. The paper advances the need for gesture-based communication to help with appreciation and

gaining for students with a social occasion handicap from the beginning forward. From the different basics, they initiate that sign to help engage the youngsters to learn, survey and value the substance better. It is comparably exhaustively supported by the past examination, which outlines the advantage of sign-based correspondence for discernment and learning language and syntax. Considering these experiences proposes a framework that can guarantee solid advancement by offering hint help to the conversation that obstructs students from consuming costly trades, for example, any comparability to YouTube. The construction is versatile and simple to use in a survey hall setting which can make an important augmentation to help the data set of students who in any case feel that it's difficult to appreciate and learn content outside their audit passages. Besides, the evaluation likewise features various modules that consolidate the construction and how each module was made based on student educator affiliations and bits of knowledge to guarantee the most incredible obligation from the students.

In the paper [8], an advantageous application that will assist Ukrainian with shocking individuals to pass on 'in a hurry' for adventurer purposes effectively without help from others is something based on the Voyager sign interpretation framework is portrayed in the paper. As indicated by the assessments of the Ukrainian Society of the hearing impaired, there are in excess of 100 thousand individuals with hearing handicaps. According to the reasonable perspective, the thing has been made in the work, with the assistance of which it is viable to complete not just the interpretation of Ukrainian conferred in and correspondences through stamping, yet in like manner semantic evaluation of these dialects. The course of action and execution of the PC understanding strategy of the Ukrainian sign-based correspondence are one of the major problems of the present, which should be tended to. For the interpretation into Ukrainian Sign Language utilize the technique that contains the going with that is first we will enter the sentence. The information relies on the user of the application (hearing handicapped). Hearing individuals can enter the information sentence by making each word with a screen console or through talk confirmation (Google Speech API). The individual with hearing impedances basically can enter the responsibility by framing words with a screen console. What's truly Separating input information on contemplations and normal words. In this development, the user's analysis sentence is confined into phrases as shown by the standards, Using assessment of understanding. In this development, the standard-based calculation of understanding Ukrainian Sign Language utilizing thought word reference is utilized. Similarly, at this stage, the word interest in the Ukrainian correspondence through stamping sentences is considered, on the grounds that the movement-based correspondence has its highlight and sentence structure. Then, at that point, Using picture portrayal or genuine human records portrayal of signs. At this development, we really want to imitate the signs as per the interpreted sentence in the past development. We can utilize human records portrayal or picture portrayal of signs.

In the paper [9], the work includes the task of making an understanding conveyed in language to signal based correspondence and presents game plans to make a way for two-sided correspondence between the debilitated meeting and the rest of the world. The paper proposes another ISL dataset with speech procedure and avoids the costly sparkle clarifications. The paper proposes performing different assignments utilizing a transformer network joined with a cross-measured discriminator to make presents directly from the speech modules. The underpinning of this proposed model is a changed transformer configuration introduced. In most typical language dealing with tasks, both the information and the goal space comprise of discrete language. In fact, for the endeavor, both data and target have a spot with the steady space. Transformer designing contains three critical parts: (i) a joint talk encoder, (ii) a present decoder and (iii) a message decoder disentangling the complementary modalities. Finally, the system uses a cross-measured coordinating association as a discriminator to help the transformer with acquiring the incredible translations from talk to correspondence through signing.

The paper [10] is the system assembled game plan based concerning Kinect v2 made in this exploration work offers an instrument through which people with hearing and talking insufficiencies can normally interact with the rest of the world. These people currently use various videos and pictures to talk and pass on/acknowledge their messages. However, this doesn't handle the communication issue, as regular language speakers don't grasp the conveyed message through hand gestures. Consequently, there exists a communication gap between these two communities. The developed framework gives the two-fold strategy for communication between the normal people & those with hearing disabilities. It doesn't simply decrease effort and time for a nearly deaf person in correspondence yet would moreover traverse the correspondence opening between incapacitate neighborhood average people. This framework, by and large, is a utility for humanity and particularly for the deaf neighborhood. Since the framework gives the twofold strategy for correspondence it has been requested into two free modules. The first module, signed by a speech translator, records gestures from the deaf individual, grasps these movements, and converts them into speech that can be perceived by normal people easily. Similarly, the subsequent module, speech to sign translator, takes normal speech as input, processes the language, and shows related sign motions on the screen along with captions. These exercises are performed by a 3D humanoid model continuously and a connected subtitle is displayed on the screen to help the client in knowing the structure. This system is strong enough considering the precision figures 84% for motion-based correspondence to talk and 87% for talk to correspondence through marking change. Accommodation is one more notwithstanding the mark of this system as it doesn't require such people to wear any gadget to perform signs and they can perform hand development as actually and deftly as they do in their existence. As the made system relies upon Kinect for Windows v2, so the structure is practical with any contraption which maintains Kinect v2.

In the paper [11], the authors developed an English to British Sign Language translation system. It includes a semantic level of depiction for performing English examination to the BSL age. It consolidates the assessment of motion-based communication conveying using different ways. This system is straightforward. The client enters English text into the system, at this stage, the client can change or change the text as demonstrated by the need. From there on out, at the syntactic stage, the inputted text is parsed with an interface grammar parser. From this parser, a mostly semantic depiction is made as Discourse Representation Structure (DRS). From this depiction, morphology and accentuation of sign age are portrayed in Head-Driven Phrase Structure Grammar (HPSG).

In the paper [12] the Proposed structure gives a viable strategy to assist correspondence between an individual with hearing and talk impedence. It was a field that had little advancements all through the long time, particularly ineffectual execution in the Python programming language. The system can similarly continue like an enlightening gadget to learn ISL. To ensure consistency like the video, breathed new life into characters, finding the correspondence by means of motions can be combined. Since there is a slight variety in the syntax of ISL depending upon the district of India, the system can be adjusted to suit the necessities of the meetings and talk debilitated around there. This ought to be conceivable by making custom corpora by taking the help of schools and various workplaces for the in need of hearing disabled people. In like manner, Google Translate APIs can be utilized to help translate from any Indian regional language to ISL. This would chip away at the correspondence between an individual not taught in English and

an individual who uses ISL. The proposed system can in like manner be used in applications, like video conference, to translate the substance into ISL logically. This would assist the hearing handicapped community acquiring a predominant perception of the one-of-a-kind situation and that implies behind the substance being shown.

The author [13] developed a system subject to a machine learning methodology. It acknowledges English text as data and produces signs contrasted with the inputted text. The system involves four major modules which are: input message preprocessor and parser, LFG f-structure depiction, Transfer Grammar Rules, ISL Sentence Generation and ISL mix. Normal English sentences are inputted to the parser. Clear sentence suggests which sentence has recently a solitary guideline activity word in it. The parser module parses the sentence and makes a dependence tree. An articulation inquiry table of around 350 articulations and common enunciations is made prior to parsing. English morphological analyzer perceives most of things. LFG common sense development (f-structure) encodes syntactic association of the information sentence. It furthermore integrates the higher syntactic and utilitarian information depiction of a sentence. Where quality is for the name of a syntactic picture and worth is for incorporating moves by the constituent. It transforms into a commitment to the age module which applies move language, so it could move source sentences into target structure. Lexical assurance and word demand correspondence are two crucial exercises that are performed during the age stage. Lexical assurance is done with the help of English to ISL bilingual word reference. ISL uses Subject-Object-Verb (SOV) word demands.

This paper [14] proposes a framework that maintains ISL translator talk. Since the hearing handicapped feels that it is challenging to communicate with people and their world in any way possible, this program ought to maintain them. This article settles an issue with correspondence the next day and therefore proposes a program that helps the neighborhood translating the communications through motions including a PC recipient or a media plan in PDAs. This paper contains information about such communications. Speech to Indian Sign Language Translator is a proposed programming structure executed using programming language, AI, ML and NLP. Sound data using python PyAudio module. The sound which has been recorded with the help of Pyaudio it is then converted using speech library or this structure can use the Google talk API (application program interface) Dependency parser helps in seeing the grammar or separating the semantic plan of a particular sentence, it helps in setting up association among words and the words which change those words. The message divider then separates the message using NLP techniques which kills weakness or disambiguates input sentences to convey machine depiction language python.

In the paper [15], The application acknowledges speech as input, transforms it into text and subsequently shows the Indian Sign Language pictures. The principal objective is to help with individuals encountering the issue of hearing. There have been numerous exercises done on motion-based interchanges that convert correspondence through marking as commitment to text or sound as result. However, sound to motion based correspondence change systems have only from time to time advanced. It is significant to both standard and in need of hearing aid people. This paper presents another development that is sound to motion-based correspondence translators. In this it acknowledges sound as data, glance through that recording using google programming connection point, shows the text on screen finally it offers the clue code of the given data using ISL (Indian Sign Language) generator. This endeavor didn't revolve around looks anyway prominent looks pass on a huge piece of correspondence through marking. This system can be done in various regions including Accessing Government Websites wherein no video cut for almost deaf and calm is open or wrapping up structures online where no interpreter is accessible to help. The consequence of this structure will be a catch of ISL words. The predefined database will have video for each and every different word and the outcome video will be a mixed video of such words. Besides Google Speech-to-Text features, it changes sound over to message by applying brain association models in an easy-to-use API. The text is then pre-taken care of using NLP (Natural Language Processing). Natural Language Processing is the limit of the machine where it processes the text and plans it. It fathoms the meaning of the words said and moreover makes the outcome. NLP can add additional abilities to our language. We will get our information ensuing to give sound data reliant upon the NLP devices to grasp human language. Word reference-based machine translation is finally done.

In this paper [16] The Proposed framework gives an effective strategy to help correspondence between a person with hearing and discourse hindrance. The structure would additionally foster permission to information for the gathering that blocked the general population of countries like India. Also, the structure can in like manner act like an informational contraption to learn ISL. This paper presents an open web stage for creating, assembling, and running guideline-based discourse to communication through signing interpretation applications. Talk affirmation is performed using the Nuance Recognizer 10.2 device compartment, and checked outcome, including both manual and non-manual parts, is conveyed using the JA Signing image system. The stage is intended to make the part innovations promptly available to communication through sign language specialists who are not PC researchers. Understanding accentuations are written in a variation of Synchronous Context-Free Grammar changed in accordance with the attributes of gesture-based communication. All handling is completed on a distant server, with content transferred and gotten to through a web interface. Starting encounters show that basic interpretation syntaxes can be executed in a period of a couple of hours to a couple of days and produce marked results promptly intelligible to Deaf witnesses. Generally speaking, the stage definitely brings the boundary down to passage for scientist keen on building applications that create top-notch marked language.

This paper [17] has proposed an assistive contraption that executes two critical engineering highlights to meet steady taking care of requirements: pipelined dealing with and point recognizing evidence got together with equivalent taking care of. The improvement utilized eye-following-based client fulfillment assessment to stimulate figuring out accuracy moreover. The possibility of the device was assessed utilizing a significant assessment of predictable speech to-motion based correspondence understanding that uses subject allocating and continuous extraction dependent upon conveyed words or expressions. The evaluation results showed that our eye-following-based client fulfillment analysis had the choice to decrease the blunder paces of the translation by 16% (per the SER metric) and chipped-way at the precision by 5.4% (per the BLEU metric) while moreover changing in accordance with unending time limitations. These outcomes recommend that this specific continuous application could profit from the use of pipelined handling and subject extraction to meet the fundamental consistent cutoff times. The solace of the made gadget was additionally researched by hard of hearing clients. The possible outcomes of the comfort study showed that our assistive gadget was satisfying and fulfilling to them, and may add to the more fundamental responsibility of need of a hearing aided peoples in commonplace activities.

In this framework [18], The Automatic Speech Recognizer (ASR) used is a condition-of-the-craftsmanship talk affirmation structure made at GTH-UPM. It is a without speaker incessant talk affirmation system reliant upon HMMs (Hidden Markov Models). The component extraction consolidates CMN and CVN (Cepstrum Medium and Variance Normalization) techniques. The speech recognizer system offers one sureness benefit for each word apparent in the word gathering. The Phrase-set up understanding structure

is based on the item conveyed to help the normal task at the 2010 NAACL Workshop on Statistical Machine Translation. The Moses decoder is used for the translation cycle. This program is a bar examine decoder for express-based authentic machine understanding models. To get a 3-gram language model, the SRI language exhibiting apparatus compartment has been used.

In this paper [19], The plan, improvement, and evaluation of a preliminary translation framework that expects to help exchanges between a hard of hearing individual and an agent in a Post Office is depicted. The framework utilizes a discourse recognizer to perceive discourse from the Post Office assistant and afterward blends the perceived expression in British Sign language (BSL) utilizing an uncommonly evolved symbol. The essential objective in cultivating this model framework was to conclude the manner by which it would be to a client whose first language was BSL and to find which district of the system expected more imaginative work to make it more feasible. The framework was tested by six pre-lingually altogether deaf people and three Post Office specialists. Deaf clients and Post Office specialists were solid of the framework, at this point the past social occasion required a more prominent check from the image and the last a framework that was less constrained in the articulations it could see: both these areas are being tended to in the accompanying time of the development.

In the paper [20], the new system contains a few pieces for getting, playing, and changing the Arabic text over to hail-based correspondence as well as an opposite procedure for getting around. The hidden segment is utilized to show correspondence through movements for needing a hearing aided and normal people. Learning decision is interfacing with and working with into a few social gatherings. Every social event contains things that are related to one another. Following picking the social affair, an outline view will show up for a long time and sign portrayal. The huge informational arrangement is a certain dataset that contains all of the photographs of the signs. The resulting dataset contains a near significance for every single sign. The new application gives an instinctual Android game, which depends right after seeing the equivalent characteristics between four pictures looking out for four verifiable things in Arabic correspondence through signals. The proposed nearly deaf game gives one more methodology to fabricate the data on the correspondence through marking in a tomfoolery and charming way. In encouraging this application, the Android SDK and nearby progression pack (NDK) were used.

### III. EXISTING SYSTEM

Most of the existing system deals with accepting input in the form of text from the user and displaying the still images of the words in the given sentence as still gesture images for each alphabet. Some other works deals with motion pictures of hand-signs for each word in a text. Most of them are accepting the input in the form of text but receiving input in the form of speech is found to be very rare.

### IV. PROPOSED SYSTEM

The proposed system receives input as speech/text and equivalent sign languages animations are displayed using AI and ML techniques.

#### 4.1) Methodology

Our system is subdivided into three modules. The first module consists of the conversion of speech to English sentences using speech recognition techniques. Then, the text is undergoing preprocessing in the second module, and some natural language processing techniques are applied to the text to prepare them ready for the next steps. Then, a word match is checked for the Animation dataset with the tokenized sentence after the text preprocessing steps. Finally, if a word match is found, the corresponding animations are displayed, otherwise the words are finger spelled using Animation character/avatar.

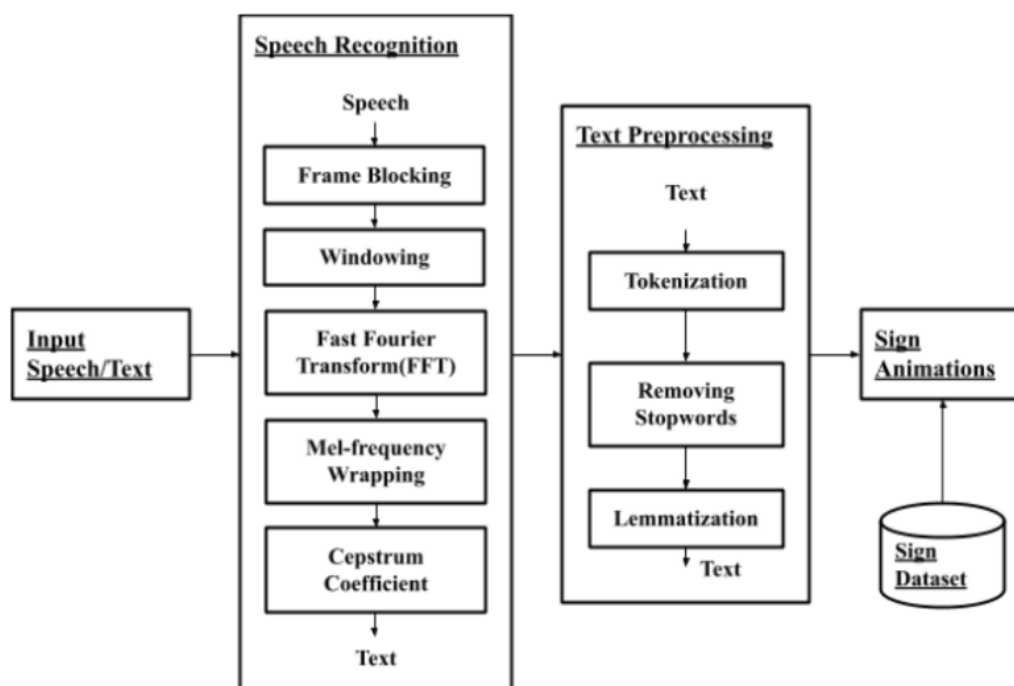


Fig.1: Block Diagram of Speech to Sign Translator

## 4.2) Speech Recognition

Speech Recognition, which is our most memorable module, is a supervised learning process. Here, info will be the sound signal and we need to anticipate the text from the sound signals. We can't take the crude sound signal as a contribution to our model since there will be a great deal of noise & ambiguity in the sound signal. It is seen that extricating features from the sound signal and using it as input to the base model will produce much better performance than directly considering raw audio signal as input. MFCC is the broadly involved procedure for extracting the features from the sound sign.

### 4.2.1) Building Speech to Text Model

We use Tensor flow and Keras for Building the model of speech to text engine and for the training. Tensor Flow provides a collection of workflows that we require to develop and train models using Python.

We will fabricate the speech to text model utilizing conv1d in Kaggle. Conv1d is a convolutional neural network which performs the convolution along only one dimension, which is often used in pattern recognition models and extract the feature from the vectors. We execute the model utilizing Keras utilitarian API.

#### 4.2.1.1) Dataset collection

Tensor Flow recently released the Speech Commands Datasets. It includes 65,000 one-second-long utterances of daily used short words, by a huge number of various individuals. The datasets are downloaded from kaggle to assemble a discussion insistence framework that sorts out spoken orders. These included instances of individuals talking single-word orders, instead of conversational sentences, so they were provoked for individual words throughout the span of a brief meeting. The documents are facilitated into coordinators, with each list name denoting the word that is communicated in every one of the contained sound records. No details were kept of any of the members, gender, age or place, and irregular ids were allocated to every person. The sound clips haven't been defined as test, validation and training sets separately, yet by show a hashing limit is used to give out each record to a set consistently. To assist the training network to adjust to noisy circumstances, it might be valuable to mix in sensible realistic establishment sound. The `_background_noise_` folder in the dataset contained a bunch of longer sound bites that are either accounts or numerical reenactments of noise.

Speech Commands Data Set v0.01 is a lot of one-second. wav sound documents, each containing a solitary message in English. These words are from a little arrangement of orders, and are spoken by a wide range of speakers. The sound records are coordinated into organizers in light of the word they contain, and this informational index is intended to assist with preparing basic AI models. We also contributed to the data collection with some daily-used words spoken in our own voice.

#### 4.2.1.2) Preprocessing

LibROSA and SciPy are the Python libraries that we used for processing audio signals.

**Librosa:** We use this python package for audio analysis to work with the audio data for the Recognition of speech. Librosa is used to visualize the audio signals and also do the feature extractions & feature matchings in it using different signal processing techniques.

**SciPy:** We use SciPy for plotting audio files as time series since the most basic step in signal processing of audio files is to visualize an audio sample file as time-series data. The most basic step in signal processing of audio files is to visualize an audio sample file as time-series data. Sound sounds can be considered a one-layered vector that stores mathematical qualities relating to each example. The time-series plot is a two layered plot of those example values as an element of time.

Testing the signal is a course of switching a simple signal over completely to an advanced signal by choosing a specific number of tests each second from the simple signal. We are changing over a sound signal to a discrete signal through testing with the objective that it will in general be taken care of and dealt with successfully in memory. We test the sample rate of audio signals and resample it by defining the sampling rate or sampling frequency as 16,000 Hz as the number of samples selected per second.

The initial phase in discourse acknowledgment is to remove the highlights from a sound signal which we will contribute to our model later.

'Mel-Frequency Cepstral Coefficients'(MFCC) is that one feature being used in any machine learning experiment involving audio files. librosa is an API for including extraction and handling information in Python. librosa feature MFCC is a technique that works on the most common way of acquiring MFCCs by giving contentions to set the quantity of edges, jump length, number of MFCCs, etc. In view of the contentions that are set, a 2D exhibit is returned.

#### Mel Frequency Cepstral Coefficients (MFCC):

Mel Frequency Cepstral Coefficients (MFCC) is coefficients that address sound, in view of perception. It is deduced from the Fourier Transform (FFT) or the Discrete Cosine Transform (DCT) of the sound clip. The basic differentiation between the FFT/DCT and the MFCC is that in the MFCC, the recurrence groups are arranged logarithmically (on the mel scale) which approximates the human hear-able structure's response more eagerly than the directly dispersed recurrence groups of FFT or DCT. This considers better handling of data. The principal motivation behind the MFCC processor is to copy the way of behaving of the human ears.

The most-clear method includes deciding the normal energy of the sound signal. This measurement, alongside all out energy in the signal, demonstrates the volume of the speaker. The signal is processed in the frequency domain through the (Fast) Fourier Transform. The windowed test is utilized to get exact portrayals of the frequency content of the signal at various moments. By taking the square of the signal value at every window test, power spectrum can be inferred. Then, at that point, the upsides of the power spectrum as the features. The three biggest recurrence tops for every window are acquired and add those to the feature vector.

Here, the sound sign is addressed by the adequacy as a component of time. In straightforward words, it is a plot among plentiful and time. Time-space investigation totally overlooks the recurrence part though recurrence area examination gives no consideration to the time part. We can get the time-subordinate frequencies with the assistance of a spectrogram.

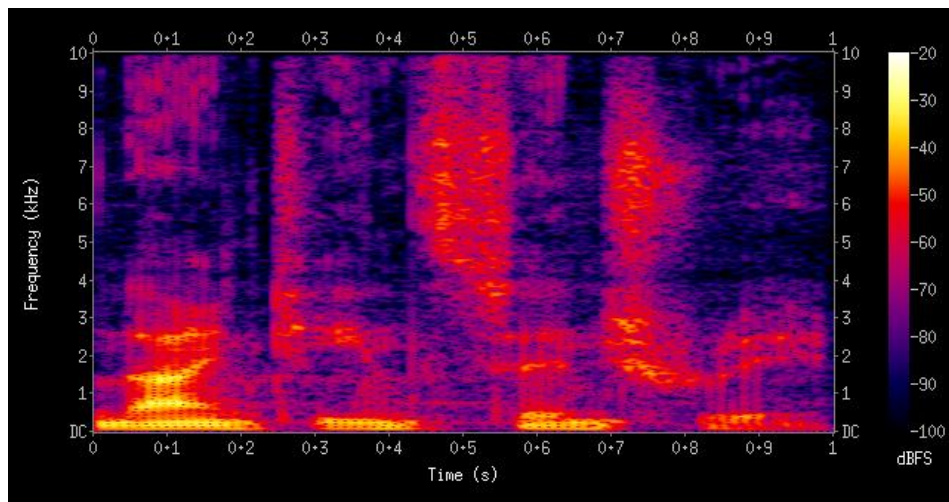


Fig.2: Spectrogram

It's a 2D plot among time and recurrence where each point in the plot addresses the sufficiency of a specific recurrence at a specific time with regards to force of variety. In basic terms, the spectrogram is a range (expansive scope of shades) of frequencies as it fluctuates with time.

**4.2.1.3) Training**

To begin with, we Split information into a train and approval set. Then, we will prepare the model on 80% of the information and approve on the leftover 20%:

We build the speech-to-text model architecture using conv1d for training. Conv1d is a convolutional neural network which performs the convolution along only one dimension, which is often used in pattern recognition models and extract the feature from the vectors. We implement the model using Keras functional API. In the CNN, we define the activation functions & the loss function to be categorical cross-entropy since it is a multi-classification problem.

The audio signals that are given as input are passed to the primary convolution layer; moreover, the outcome is obtained as an initiation work. The channels that are applied in the convolution layer extract the required features from the info picture. Each channel gives an alternate feature component to help the right prediction of the class.

As we go more profound in the network, more unambiguous elements are separated when contrasted with a shallow network where the features extricated are more nonexclusive. The result layer in a CNN as referenced already is a completely associated layer, where the contribution from different layers is joined & sent so as to change the result into the quantity of classes as wanted by the network. Early pausing and model designated spots are the callbacks that empower it to quit preparing the neural network with impeccable timing and to save the best model after each epoch.

**4.2.1.4) Validation**

The visualization to understand the performance of the model is given below:

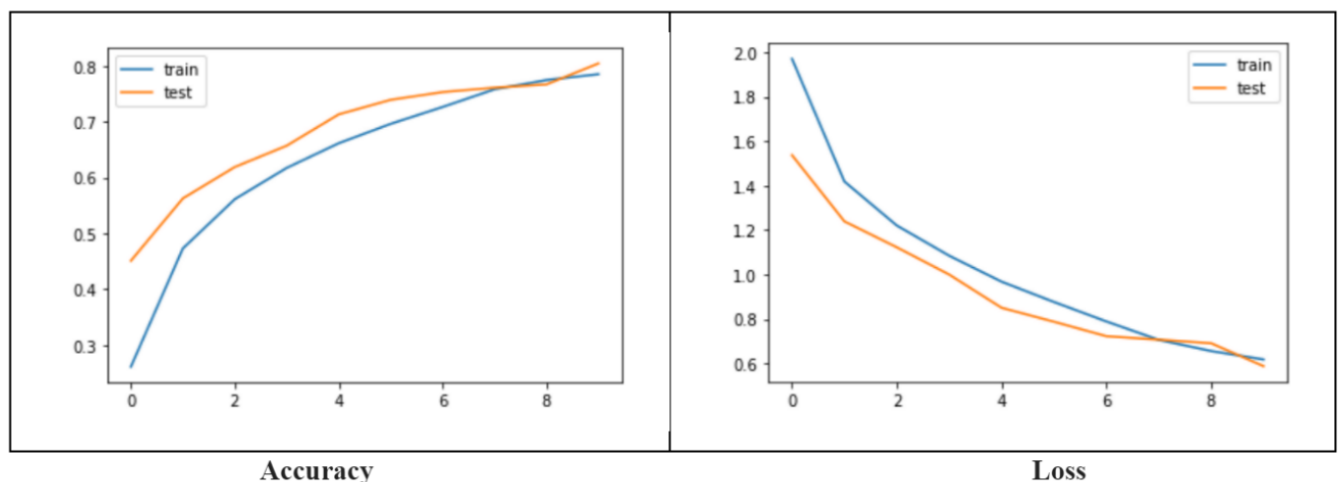


Fig.3: Model Performance

**4.2.1.5) Deployment**

From CNN, we load the best model that causes the best expectations on the approval data. Then we characterize the capacity that predicts text for the given audio.

**4.3) Text Preprocessing**

Text preprocessing is a significant and one the most fundamental steps prior to building any model utilizing Natural Language Processing (NLP). Before we analyze that data programmatically, we first need to preprocess it. Natural Language Processing (NLP) is a field that spotlights on making normal human language usable by computer programs.

The steps involved in the change of interaction from English sentence to its related ASL message are remembered for this segment. The words in the ASL sentence have been identified to generate corresponding sign movements. For the conversion from English to ASL, we use the Natural Language Toolkit (NLTK), which is a set of open-source Python modules used to work with human language data for applying statistical natural language processing. A rough text corpus, accumulated from one or many sources, may be stacked with anomalies and ambiguity that requires preprocessing for cleaning it up. For sure, even directly following cleaning, more text preprocessing is supposed to reshape the data such that it might be dealt with clearly to the model.

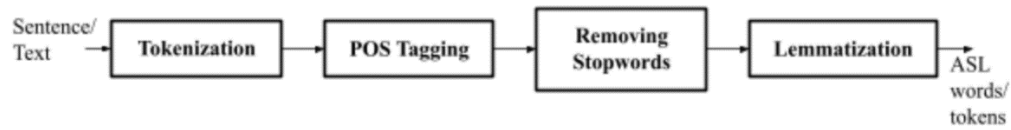


Fig.4: Text Preprocessing

#### 4.3.1) Tokenization

Tokenization is the most common way of parting text into pieces called tokens. A corpus of text can be changed into tokens of sentences, words, or even characters. This will allow us to work with smaller pieces of text that are still relatively coherent and meaningful even outside of the context of the rest of the text. We pass our English sentences from the information message box and messages from the discourse acknowledgment module through the tokenization interaction since it's our initial phase in transforming unstructured information into organized information, which is simpler to break down. For our ASL linguistic structure, we really want to isolate each word and fabricate a jargon with the end goal that we can address all words particularly in a rundown from the sentence. Numbers, words, etc. all fall under tokens.

Here, we convert the text into word tokens during preprocessing as they are prerequisites for some NLP operations. The tokenization of sentences also helps us to look for a word-match in our sign animation data set.

Example of word tokenization:

Word\_tokenize ("Welcome to the world of sign")

Output = ['Welcome', 'to', 'the', 'world', 'of', 'sign']

#### 4.3.2) POS Tagging

A POS tag is a unique name appointed to every token or word in the text corpus to show the grammatical form and frequently additionally other linguistic classifications, for example, tense, number—either plural or single case, etc. After the tokenization, we will initially utilize the POS label depicted as the cycle where each word is worked with its phonetic component. Additionally, generally known as POS tagging. Grammatical features incorporate adjectives, conjunctions, verbs, adverbs, nouns, pronouns, and sub-categories. POS tags are useful for lemmatization and extracting relationships between words.

Some of POS tags used are given:

<u>POS TAGS USED</u>	<u>DESCRIPTION</u>
MD	modal (could, will)
VBP	verb, present tense not 3rd person singular(wrap)
VBZ	verb, present tense with 3rd person singular (bases)
VBG	verb, gerund(judging)/present participle taking
VBD	verb, past tense took(pleaded)
VBN	verb, past participle taken

Table: POS Tags

Example of POS tagging:

Sentence = "I am going to school"

Output after POS tagging:

[('I', 'PRP'), ('am', 'VBP'), ('going', 'VBG'), ('to', 'TO'), ('school', 'NN')]

#### 4.3.3) Removing Stop words

Whenever we utilize the features from a text to model, we will experience a ton of commotion or noise. Stop words are typical language words which have not any particular meaning, for instance, "a", "an", "the", "or, etc. that appear so routinely in the text that they could reshape different NLP practices without adding a huge load of critical data. So almost always we have to remove stop words from the corpus as part of our preprocessing. For this, we import stop words from nltk.corpus.

NLTK library keeps a rundown of around 179 stop words that can be utilized to channel stop words from the text. We can similarly add or dispense with stop words from the default list.

Example of Stop words removal:

Example: "This", "is", "an", "apple"

After stop word removal: "This", "is", "apple"



#### 4.3.4) Lemmatization

Lemmatization is changing the word completely to its base structure or lemma by eliminating affixes from the inflected words. It assists with making better elements for AI and NLP models since it is an important step in preprocessing.

Here, we need to extract the base form of the word for ASL syntax. The word extracted here is known as Lemma and it is available in the dictionary. We have the WordNet corpus and the lemma created will be accessible in this corpus. There are many Lemmatizers available in NLTK that employ different algorithms. In our project, we have used the WordNetLemmatizer that makes use of the WordNet Database to lookup lemmas of words.

Example of Lemmatization:

[ change/ changing/ changes/ changed] = change

#### 4.4) ASL Animations

After processing the English sentence, we proceeded to generate the sign movement for the ASL tokens. The avatar generates animated sign movements for each word matches, otherwise, those tokens are fingerspelled. Here, we used the blender software for bringing motion to the signing avatar & created more than 150 animation clips for the dataset (also adding new clips) with some words used in one's daily-life.

The entire process for making the 3d avatar is isolated into three stages. To start with, the skeleton and face of the animation character are made. In the next step, we characterize the perspective or direction of the model. In the third step, we characterize the development joints and facial looks of the 3d-character.

Then, we give the frame sequences that decide the motion of the avatar over the given sequence of words according to the ASL. At last, movement (showing gestures, moving hands, arms and so on) is characterized by giving solid animation rendering.

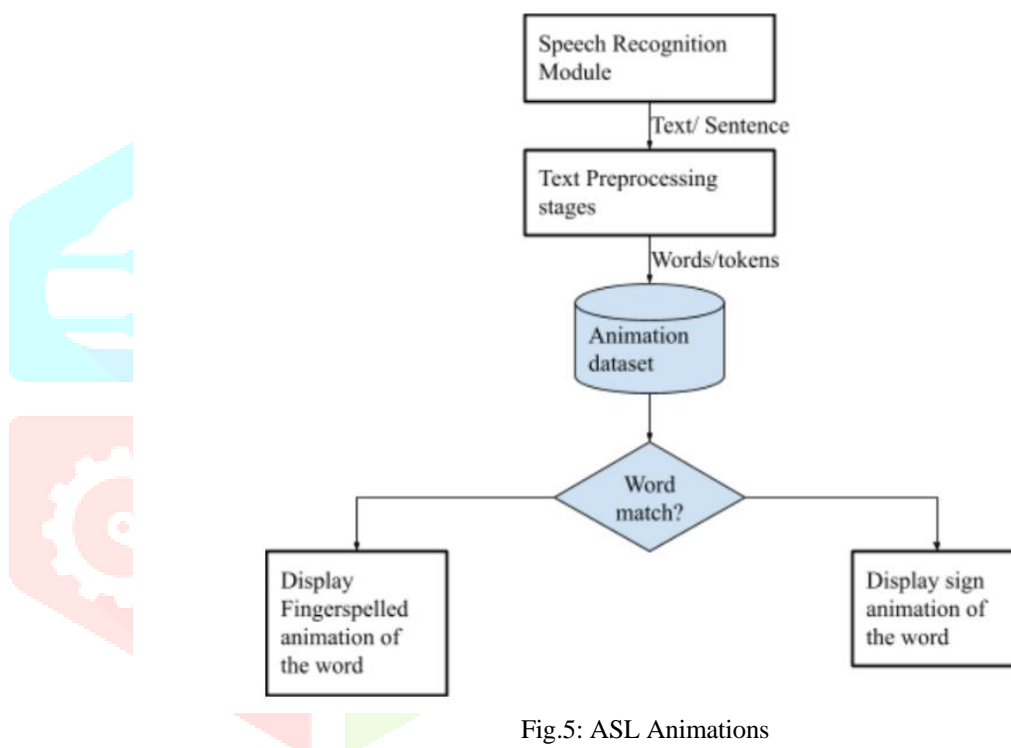


Fig.5: ASL Animations

## V. RESULTS AND DISCUSSION

The Speech to Sign Translator has three modules. That is speech recognition, text preprocessing and 3D animation. Speech recognition offers a good accuracy currently, but wrong recognition of words, and misspellings are common in speech recognition systems in proportion to the clarity of speech delivered by the speaker. That is, Cepstral coefficients are said to be accurate in certain pattern recognition problems relating to human voice. They are used mostly in speech recognition & speaker identification. MFCC highlights are not exactly accurate in that frame of mind of a lot of foundation commotion and probably won't be appropriate for speculation. In order to reduce these problems, we use training dataset under different environment conditions including noise.

The ASL animation dataset contains a set of sign animations for the ASL words and they are animated manually. The minute movements of fingers and sign require a lot of time for bringing animations to the 3d character. The contribution and availability of suitable ASL animations can enhance the usability and efficiency of this Speech to sign translation system. It generates realistic 3D animations.

In the proposed system, The Speech to Sign Translator system intakes input speech and delivers the respective American Sign Language (ASL) sign language animations in live mode. It generates realistic 3D animations. It assists the active participation of deaf people in their day to day professional and personal lives and doesn't limit their participation in any field of work. Moreover, it takes out the need of going to costly Sign Language classes, as it tends to be advanced at home effectively utilizing the created framework. The model furnishes us with great exactness and can likewise act as a way to give fundamental instruction to kids. In addition, it can additionally be carried out in different fields like gaming, simulations, and various other sectors.

## VI. CONCLUSION AND FUTURE SCOPE

The hard of hearing and hearing-handicapped form a community with explicit requirements that administrators and innovation have as of late designated. There is no such transparently available application that translates speech to sign animations, with a sensible cost. Speech to Sign interpreter systems can work in the presence of almost deaf people as motion-based correspondence is the main method for data sharing for them.

In the proposed system, The Speech to Sign Translator system intakes input speech and delivers the respective American Sign Language (ASL) sign language animations in live mode. It generates realistic 3D animations. It assists the active participation of deaf people in their day to day professional and personal lives and doesn't limit their participation in any field of work. In addition, it discards the need of going to expensive Sign Language classes, as it might be progressed at home really using the made structure. The model gives us good accuracy and can likewise act as a way to provide basic sign language education to children and the hearing handicapped community. Plus, it can moreover be executed in various fields like gaming, simulations, and different areas. In schools, universities, medical clinics, colleges, air terminals, courts anyplace anybody can involve this framework for understanding the gesture-based communication to impart. It makes correspondence between a run of the mill hearing individual and a hard to hearing individual easier.

The future work is to implement this system in different media platforms, where the corresponding sign animations are displayed in one corner of the screen while the media or platform is playing. At this moment some news channels are utilizing this sort of show yet they are utilizing individual appearance signs as indicated by the speech of the individual giving news live. So, this will be a superior thought which we can provide for news channels by replacing humans as they need more experience and expertise in sign language. Moreover, the availability of such human translators is very difficult today. We expect to broaden the errand by furthermore including looks into the framework.

## VII. ACKNOWLEDGMENT

We have great pleasure in publishing paper on Speech to Sign Translator. We take opportunity to express our sincere thanks towards our guide Mrs. Nighila Ashok; Department of Computer Science & Engineering, for providing the technical guidelines and suggestions regarding the line of work. We would like to express our gratitude towards her constant encouragement, support and guidelines through the development of the project. We thank Dr. Sreeraj. R; Head of Department of Computer Science & Engineering, for his encouragement during progress meetings and for providing guidelines to write this project. We thank Mrs. Chinju Poulose; project coordinator, Department of Computer Science & Engineering, for being encouraging throughout the course and for guidance. Lastly, we would like to thank our college principal, Dr. Jose. K. Jacob; for providing lab facilities and permitting us to go on with our project. We would also like to thank our colleagues who helped us directly or indirectly during this project.

## VIII. REFERENCES

- [1] Stephen Cox, Michael Lincoln, Judy Tryggvason, Melanie Nakisa, Mark Wells, Marcus Tutt, Sanja Abbott: "TESSA, a system to aid communication with deaf people" Conference Paper, January 2002. [2]. Palaz, Dimitri, Mathew Magimai Doss, and Ronan Collobert: "Convolutional neural networks-based continuous speech recognition using raw speech signal." In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4295-4299. IEEE, 2015.
- [2] Palaz, Dimitri, Mathew Magimai Doss, and Ronan Collobert: "Convolutional neural networks-based continuous speech recognition using raw speech signal." In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 4295-4299. IEEE, 2015.
- [3] Angus B. Grieve-Smith: SignSynth: A Sign Language Synthesis Application Using Web3D and Perl, <https://www.researchgate.net/publication/2488471>, Article · November 2001.
- [4] Stephanie Stoll, Necati Cihan Camgoz, Simon Hadfield, Richard Bowden: "Text2Sign: Towards Sign Language Production Using Neural Machine Translation and Generative Adversarial Networks", International Journal of Computer Vision, 2020.
- [5] Ezhumalai P, Raj Kumar M, Rahul A S, Vimalanathan V, Yuvaraj: "A Speech to Sign Language Translator For Hearing Impaired".
- [6] Indian and English Language to Sign Language Translator- an Automated Portable Two Way Communicator for Bridging Normal and Deprived Ones, 2nd international conference on power, energy, control and transmission systems ,2020.
- [7] Prof. Abhishek Mehta, Dr. Kamini Solanki, Prof. Trupti Rathod: Automatic Translate Real-Time Voice to Sign Language Conversion for Deaf and Dumb People, ICRADL - 2021 Conference Proceedings.
- [8] Lozynska Olga, Savchuk Valeriia, Pasichnyk Volodymyr- "The Sign Translator Information System for Tourist", CSIT 2019, 17-20 September, 2019, Lviv, Ukraine.
- [9] Parul Kapoor 1 , Rudrabha Mukhopadhyay2 , Sindhu B Hegde2 , Vinay Namboodiri1,3 , C V Jawahar- "Towards Automatic Speech to Sign Language Generation", on 11 July 2021.
- [10] Mateen Ahmed, Mujtaba Idrees, Zain ul Abideen, Rafia Mumtaz, Sana Khaliq: "Deaf Talk Using 3D Animated Sign Language", SAI Computing Conference, 2016.
- [11] JA Bangham, SJ Cox, R Elliott, JRW Glauert, I Marshall: "Virtual Signing: Capture, Animation, Storage and Transmission" – an Overview of the ViSiCAST Project, Conference: Speech and Language Processing for Disabled and Elderly People (Ref. No. 2000/025), February 2000.
- [12] Alisha Kulkarni, Archith Vinod Kariyal, Dhanush V, Paras Nath Singh: "Speech to Indian Sign Language Translator", Proceedings of the 3rd International Conference on Integrated Intelligent Computing Communication & Security, 2021.
- [13] Tirthankar Dasgupta, Sandipan Dandpat, Anupam Basu: "Prototype Machine Translation System from Text-To-Indian Sign Language", 2008.
- [14] Kajal Jadhav, Shubham Gangdhar, Viraj Ghanekar: "SPEECH TO ISL (INDIAN SIGN LANGUAGE) TRANSLATOR", International Research Journal of Engineering and Technology, 2021.

- [15] Ankita Harkude#1, Sarika Namade#2, Shefali Patil#3, Anita Morey #4- “Audio to Sign Language Translation for Deaf People”, Volume 9, Issue 10, April 2020.
- [16] An Open Web Platform for Rule-Based Speech-to-Sign Translation
- [17] Mwaffaq Ootom PhD & Mohammad A. Alzubaidi PhD: Ambient intelligence framework for real-time speech-to-sign translation”, Assistive Technology, DOI: 10.1080/10400435.2016.1268218, (2017).
- [18] V. López-Ludeña, R. San-Segundo, R. Córdoba, J. Ferreiros, J.M. Montero, J.M. Pardo: “Factored Translation Models for improving a Speech into Sign Language Translation System”, DOI: 10.21437/Interspeech.2011-481 · Source: DBLP, Conference Paper · August 2011.
- [19] Stephen James Cox, Mike Lincoln, Judy Tryggvason, Mel Nakisa: The Development and Evaluation of a Speech-to-Sign Translation System to Assist Transactions, International Journal of Human-Computer Interaction · October 2003.
- [20] Emad E. Abdallah, Ebaa Fayyoumi: “Assistive Technology for Deaf People Based on Android Platform”, The 11th International Conference on Future Networks and Communications, 2016

