

# Sign Language Recognition System Using Deep Learning

Pasumarthi Yaswanth Sai Varun  
Computer Science and Engineering  
Kalasalingam Academy of Research and Education  
krishnankoil, Tamilnadu, India

Sai Kiran Dangeti  
Computer Science and Engineering  
Kalasalingam Academy of Research and Education  
krishnankoil, Tamilnadu, India

Sankatala Sudheer Kumar  
Computer Science and Engineering  
Kalasalingam Academy of Research and Education  
krishnankoil, Tamilnadu, India

A.Robert Singh  
Computer Science and Engineering  
Kalasalingam Academy of Research and Education  
krishnankoil, Tamilnadu, India

**Abstract:-** Communication is demand in human life in every single day. Better communication takes you to better understanding, and it containing everyone in the community, including the deaf as well as dumb. Interfering with other people in sign language. However, most listeners could not acknowledges signing and learning it's not a simple process. Finally, there is still an unacceptable barrier between the deaf and normal. Communicating with those having vocal and hearing disabilities. To overcome this problem we have planned to develop a model to detecting the language of people like dumb usually used. For that we required to work on the large sized data set to acquire this step the only way to resolve this is we have to work this idea with deep learning concepts and with python 3 programming which is very suitable to data set training in very comfortable manner. We are going with CNN algorithm because a large data set is required to complete the task because CNN manages the image processing in a feel good manner after the completion of data set training we have to show our palm and some alphabets in sign language and it captures and showed as a text

**Keywords:-** Sign Language, CNN Algorithm, Kinect, Deep Neural Networks.

## INTRODUCTION

Deep Learning is an integral part of machine learning . For singular definitions. during ongoing years, it is all around in the world it is a technique which is like a neural network system in our brain but here instead of this , we are creating an artificial neural networks (ANN) in reality ANNs may contain thousands of nodes and billions of links. Every node is related to a layer, which is a collection of nodes.. Actually there are input as well as output layers and also layers in between them which are known as hidden layers. By addition of number of nodes, links and layers will increase the accuracy level of the ANN.

The need for this idea's implementation is very less people may understand the sign language in the world for that people's sake especially who are normal it could be more challenging. task to communicate with the people belonging to type of speaking disability people to make this task in a easier way to solve it I and our team planned an idea it is the main agenda of the project. The amount of work in this can be done in training part then it considering we have to work in CNN supervised algorithm in image processing the computer vision is using training is the main part here which is to be critical in all the work.

### 1.1 SIGN LANGUAGE :

Deaf people around the world use a visual language that uses signatures through communication, facial and body movements, perceiving signatures each have differently from their everyday language.as shown in the fig 1. Sign language is not a rampant dialect, and dissimilar sign languages are utilizing in dissimilar nations, just like the many spoken languages used everywhere on the planet.



**Fig 1 Finger Spelling in American Sign Language[18]**

The sign language we are considering is American sign language as I said before each nation follows their own sign language as per their resources considered , but here we can doing work in the American sign language system because it is efficient than others as per our news receiving from outside

## LITERATURE SURVEY

In the recent researches, there has been tremendous research on the hand sign language gesture recognition. The technology for gesture recognition is given below;

The two networks used by Lionel Pego et al, existing of 3 layers, have a maximum turnout after each. 1 of the CNNs were up skilled to take features by hand and the other from the upper body. 2 CNN outputs were concreted and feed to a fully closed layer. He used a data set from CLAP 14 containing 20 Italian gestures through 20 articles. He used both depth and color to develop the system. They obtained 91.7% accuracy from cross-validation from the training set consisting of various users with various backdrops and 95.60% testing accuracy from different back drops but it including the upskill set users and backdrop.[2][3] Al, Alina K. and ET Rahm used a various-veneer arbitrary forest model for datasets collected using the Open NI + NITE system of rules on Microsoft Kinect. He got an perfection of 87 for articles on system insurance for which he was trained and 57 for a new article. Le Lemantick and Ital used glove-based motion tracking sensors to detect signals. Regardless of the drawbacks of containing similar wispy sensors and gauntlets, it puts together it impractical to use. They worn VGG-16 architecture (CNN) to up skilled and categorize fist gestures. However, in this scenario the gesture must set down in front of a everlasting framework and any anomaly in the backdrop will end up in the wrong categorization.

In this discourse, a procedure for identifying the signature of the Argentine Signature (LSA) is suggested. This offer offers two key contributions: Create a hand-shaped database for the (LSA). Following to come up, a technique for image processing, a piece of stored data that indicates how other data is stored taking out and following palm physique differentiation that prescribed self-regulation map observing adjustments. known as probSom. This procedure is contrast to others in the form of the art, similar as Support Vector Machines (SVMs), erratic Forests and Neural connections. Probosome-based neural categorizing, exerting the appropriate detector and also achieved an accuracy rate of 90 to 90 which was excellent work.

In vision-based ways, a computer camera is a data resources appliance for displaying bumf on the palms or fingers pointing. This need a single camera, then be aware of the unprocessed association between men and machines Unless the usage of any extra appliances. These arrangements encomium biotic vision by giving explanation processed vision systems that are takes place in software and / or hardware. This presents a arduous problem because the background of these systems is not variable, the lighting is not sensitive, the individual and the camera are free to get real time performances. Furthermore, such a system needs to be increasing to meet this kind of aspiration, including rightness and sturdy. Al the process can be figured in the fig 2.

The first step is to concoct a three-acreage model of the human palm. This model is suited with hand extractions by more than two cameras, and the orientation of the palm and

the parameters related to the common angles are clearly estimated. These parameters are then used to classify the indicator. The second step is to capture the image using a camera then draw some pointers and the following features will be used as input in the ranking algorithm for classification.



**Fig 2: Block Diagram of vision based recognition system**

Indicator recognition system using OpenCV. It inspects the input signals to match the indicator database, the signal database, the gesture database and contains the information required for pattern matching. It is possible to visualize the vocabulary and the result of a gesture through a completely planned one. Due to the use of Open CV, the accuracy of this system is low. There is an error in the maximum time consumption of the property. The use of hand gloves is mandatory. However, the equipment was expensive and, essentially. They diminished the nature of sign language communication. Some background noise will also be recorded and this can lead to low accuracy yields and this can lead to correction issues. All these problems can be solved through our proposing model.

## I. PROPOSED SYSTEM

Generally, we are planning to work on Huge data sets. Those data sets will be trained by using Deep Learning Concepts. Python 3.6, TensorFlow keeps better outcomes. The hand pattern recognition procedure includes major features. In this We have several advantages like as below. While Using the Huge data sets, there must be effective results. Accurable and fast Visualizing can be possible. Very low expensive Compared to the previous methods either software and hardware usages. These are all done in Run time only that saves the executable time. This project is especially targeted at people who are unable to speak as well as those who have difficulty hearing. With the tremendous increase in the use of smart technologies, generation makes things easier. Users of this project are assured of easy communication with the general public as well as with each other. We are managing some of the best time complexities that are deep learning concepts that can be planned to get the output as the output frequency that comes to the screen causing the audio through which it Easy to understand text can be done with text to speech converter. A type that can hear but is not able to speak and is useful in understanding two ordinary people.

### A. IMAGE ACQUISITION :

It is a integral part of image processing, In the process of retrieving an image from an external source it could be a hardware based system it is the starting step in the training part . If we didn't have an image in the sense the whole process cannot be started as we are expecting the image is a 2 dimensional image to capture it we required single sensor the motion should be in x and y directions.

**B. IMAGE PREPROCESSING :**

Data image preprocessing is generally is two types i.e. analogue and digital image pre processing, Digital image processing is a subset of digital signaling, it is more advantageous than analogue, image processing system it applies a broad range of algorithms to be applied to the input .[19] The agenda of processing is to increase the image features by suppressing background unwanted noise.

**C.FEATURE EXTRACTION :**

It is sub part of reducing the dimension s firstly, we have to divide the raw data and reduced into more groups it is going to be very easy to process when ever we want. The most important factor in large data sets are actually is they have large number of variables these variables are almost consuming more computing resources to process. So feature extracting is very helpful in this problem to avoid some critical obstacles while processing. It is only useful to reduce the redundant data from a given data set.

**D.IMAGE CLASSIFICATION :**

Actually, They are two types of classification i.e. supervised and unsupervised classifications. The main concentration of classification is to dividing all the pixels into groups throughout into a thematic maps. In digital image analysis the classification plays an important role. So the picture may be as a good one which is showing a magnitude of different colors showing different features.

**E.DATA AUGUMENTATION :**

With the Conventional Neural Network (CNN). In fact, it would not be wrong to say that AI re-emerged (after many AI winters) simply because of the availability of giant computing power (GPU) and the vast amount of data on the Internet. More fortunately for me, there is a lot of information out there in pictures and videos like this.

Although the availability of all the data, bringing the right kind of data that matches the exact usage case in our experience is a daunting task. Furthermore, good adversity has been found in the information because the purpose of interest must be present in different sizes, lighting conditions and poses if we want our network to b effectively generalized during the testing (or deployment) phase. In order to overcome this problem of limited quantity and limited diversity of knowledge, we build our data with the usual data that we have. This method of generating our own data is considered data enhancement. with visual cortex.

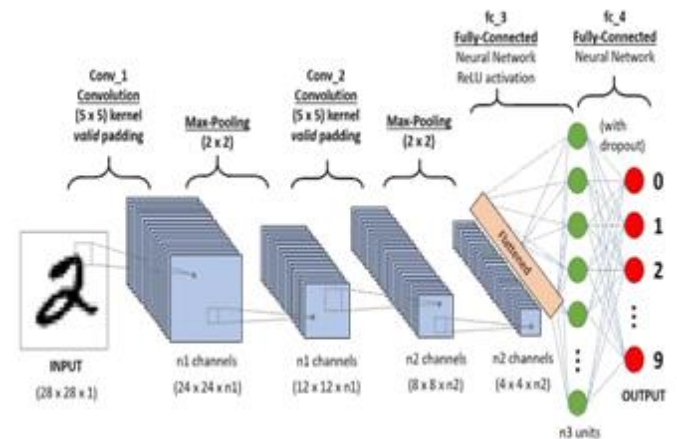
The Convolutional Neural Network is an in-depth learning algorithm that can scan an input image, assign significance (learning weight and bias) to various aspects / objects within the image and prepare for post-dissolution of the opponent. The pre-processing required during a ConvNet is similar to the diagram shown below After the training phase completion, we have to go through the execution process using TensorFlow libraries and also with the help of keras after this phase execute the code in the platform and getting the output file then we have to load into it and in the time of capturing during absorbing it matches the gestures with dataset which is already

trained one. Finally, the one step the text we have got is converting into speech using text to speech Converter.

**F. TYPES OF LAYERS IN CNN :**

Avoid combining SI and CGS units, such as current in amperes and magnetic field in oersteds. This often leads to confusion because equations do not balance dimensional. If you must use mixed units, clearly state the units for each quantity that you use in an equation.. In CNN they are several layers which are clearly shown in fig 3 which are worked on the basis of conditions what the requirements are needed by the programmer.

Let’s take an example by running a covnets on of image of dimension 32 x 32 x 3.



**Fig3: Convolution layer to layer Block diagram[20]**

**Input Layer:** This layer holds the raw input of image with width 32, height 32 and depth 3.

**Convolution Layer:** This layer calculates the output volume by computing dot product among all filters and image patch. Assume we use total 12 filters for this layer we will get output volume of dimension 32 x 32 x 12.

**Activation Function Layer:** This layer will apply element wise initiation function to the output of convolution layer. Some mutual activation functions are RELU: max(0, x), Sigmoid: 1/(1+e^-x), Tanh, Leaky RELU, etc. The volume remains unchanged hence output volume will have dimension 32 x 32 x 12.

**Pool Layer:** This layer is intermittently inserted in the convnets and its main function is to diminish the size of volume which makes the computation fast reduces memory and also prevents from overfitting. Two common types of pooling layers are **max pooling** and **average pooling**. If we use a max pool with 2 x 2 filters and stride 2, the resultant volume will be of dimension 16x16x12.

**Fully-Connected Layer:** This layer is consistent neural network layer which receipts input from the preceding layer and computes the class scores and outputs the 1-D array of size equal to the number of classes.

**G.GPU:**

A graphics processing unit (GPU) can be a specialized, electronic circuit designed to rapidly manipulate and alter memory to accelerate the creation of images during a buffer



intended for output to a display device. GPUs are utilized in embedded systems, mobile phones, personal computers, workstations, and game consoles. Modern GPUs are very efficient at manipulating special effects and image processing. Their highly parallel structure makes them more efficient than general-purpose central processing units (CPUs) for algorithms that process large blocks of knowledge in parallel. In a pc, a GPU are often present on a video card or embedded on the motherboard. In certain CPUs, they're embedded on the CPU. Identification Algorithm In the identification of the palm, Few steps used to get towards the goal, background image and input image are implemented through color transformation. Color change transfers color in RGB projects to YCbCr projects, it is better than skin color detection. Second, the background reduction operation is used to extract the skin color according to the background image and input image.

METHODOLOGY

The project focuses on developing software that converts gestures into text as well as recognizes them as output speech. First of all, we have to make the plan as shown in fig 4 in parts. as there are no other updated versions due to compatibility. And then we have to create a data set for that and we have to take different photographs with different concepts of different alphabets and numbers.. First of all, we have to train the data for which we have to use some algorithms.

Actually, we are planning the project very keenly understanding from point to point for getting the successful output. According to the completion of project we are able to divide the work;

- 1.Test
2.Train

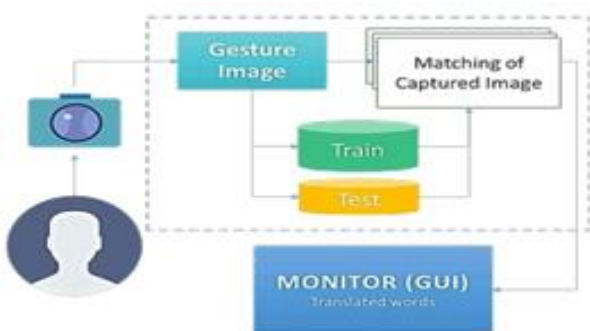


Fig 4. Block Diagram of Project Deliverables

Matplotlib produces publication-quality figures in a variety of hard copy formats and interactive environments across platforms. Matplotlib can be used in Python scripts, the Python and I Python shell, web application servers, and various graphical user interface tool kits.

H. DATA SET :

In our project we are using the sign language which is going to be best for getting maximum accuracy. In sign language normally people are used hand, facial and body movements for communication purpose. Basically they are

more than 135 sign languages in the world. By the way now we are using American Sign Language in this project because it has more benefits when it compared to the others. Although, it is a very famous among all the sign languages what we have till now.

For that we collected various photos of different signs of alphabets in American Sign Language (ASL) The used data set is shown in fig 5 as our Data set from various resources. and we had trained it using CNN required neural networks . This training is set under computer vision because all the datasets which are in heavy sized can be trained in this way only. The training can be done in other ways like the data set would be treated or trained in RNN neural networks but it can be possible only for small sized data sets, For the datasets of large sized like in this project would be followed in CNN ways that too in under computer vision process.

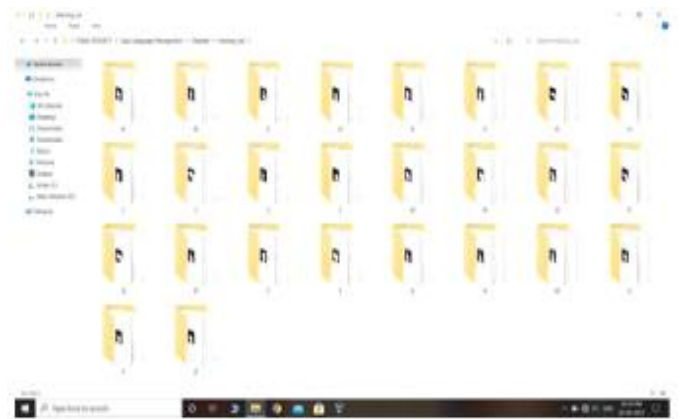


Fig : 5 American Sign Language Data set



Fig 6.Extracted Image of A

Collection of requirements and making Surveys about depth concept of Project Idea mainly concentrate in getting the complete code for execution as shown in the fig 6 Concentrating the Objectives and their scope Project Planning for success full completion of project whether all the training can be passing correctly or not. . At present dataset is collected by own standards and training is going on. with these huge data set training can be done slowly but we have to configure correctly as possible wisely or knowledge in handling data set.

## VIII. CONCLUSION

Hand gestures are a powerful method of human communication, including many potential uses in the field of human-computer interaction. Vision-based hand posture detection methods have several proven pros than traditional appliances. However, hand gesture detection is a arduous problem and the present job is only a tiny part of achieving the desired results in the field of gesture language recognition. We will planning to get an accuracy of 93.27%. This shows that CNN can be used successfully in learning local and temporal features and classifying sign language. From now on, we have to train to achieve the rest of our goals and go through better developed ways to achieve our goals.

## REFERENCES

- [1]. Ronchetti, Franco, Facundo Quiroga, César Armando Estrebow, and Laura Cristina Lanzarini. "Handshape recognition for Argentinian sign language using probsom." *Journal of Computer Science & Technology* 16 (2016).
- [2]. Pigou, Lionel, et al. "Sign language recognition using convolutional neural networks." *Workshop at the European Conference on Computer Vision*. Springer International Publishing, 2014.
- [3]. <https://biblio.ugent.be/publication/5796137/file/5796322.pdf>.
- [4]. D. Metaxas. Sign language and human activity recognition, June 2011. *CVPR Workshop on Gesture Recognition*.
- [5]. M. Ranzato. Efcient learning of sparse representations with an energy-based model, 2006. *Courant Institute of Mathematical Sciences*.
- [6]. Enhanced Local Texture Feature Sets for Face Recognition Under Difficult Lighting Conditions- Xiaoyang Tan and Bill Triggs.
- [7]. Cooper, H., Bowden, R.: Sign language recognition using boosted volumetric features. In: *Procs. of IAPR Conf. on Machine Vision Applications*, pp. 359 – 362. Tokyo, Japan (2007).
- [8]. Ekman, P.: Basic emotions. In: T. Dalgleish, T. Power (eds.) *The Handbook of Cognition and Emotion*, pp. 45–60. John Wiley & Sons, Ltd. (1999).
- [9]. Gao, W., Ma, J., Wu, J., Wang, C.: Sign language recognition based on HMM/ANN/DP. *International journal of pattern recognition and artificial intelligence* 14(5), 587 – 602 (2000)
- [10]. Han, J., Awad, G., Sutherland, A.: Modelling and segmenting subunits for sign language recognition based on hand motion analysis. *PATTERN RECOGN LETTERS* 30(6), 623 – 633 (2009).
- [11]. Holden, E., Lee, G., Owens, R.: Australian sign language recognition. *Machine Vision and Applications* 16(5), 312 – 320 (2005).
- [12]. .Ming, K.W., Ranganath, S.: Representations for facial expressions. In: *Procs. of Int. Conf. on Control, Automation, Robotics and Vision*, vol. 2, pp. 716 – 721 (2002). DOI 10.1109/ICARCV.2002.1238510.
- [13]. Ouhyoung, M., Liang, R.H.: A sign language recognition system using hidden markov model and context sensitive search. In: *Procs. of ACM Virtual Reality Software and Technology Conference*.
- [14]. Zieren, J., Kraiss, K.: Robust person-independent visual sign language recognition. In: *Procs. of IbPRIA*, pp. 520 – 528. Springer, Estoril, Portugal.
- [15]. D.Kelly, J.McDonald and C.Markham, "A person independant system for recognition of hand postures used in sign language," *Pattern Recognition Letters*, Vol.31, pp.1359-1368, 2010.
- [16]. H. K. Nishihara et al., *Hand-Gesture Recognition Method*, US 2009/0103780 A1, date of filing Dec 17, 2008, date of publication Apr 23, 2009.
- [17]. Daniel Martinez Capilla, *Sign Language Translator using Microsoft Kinect XBOX 360*, Master dissertation, EE dept., The University of Tennessee, 2012.
- [18]. <https://www.startasl.com/american-sign-language-alphabet/>
- [19]. <https://towardsdatascience.com/image-pre-processing-c1aec0be3edf>
- [20]. <https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53>