# Hand Gesture Recognition Using Deep Learning

I Hemanth and P Siva Prasad

May 15, 2024

# Hand Gesture Recognition Using Deep Learning

I.Hemanth
Department of Computer Science &
Engineering
Vignan's Foundation for Science,
Technology and Research (Deemed to
be University)
Guntur, India
201fa04374@gmail.com

Dr. P. Siva Prasad
Department of Information Technology
& Computer Applications
Vignan's Foundation for Science,
Technology and Research (Deemed to
be University)
Guntur, India
pusapatisivaprasad@gmail.com

*Abstract*— **A vital component of human-computer interaction, hand gesture recognition has applications in a wide range of industries, including virtual reality, robotics, healthcare, and sign language translation. However, extensively annotated datasets and substantial computational resources are frequently needed for training deep learning models for hand gesture detection from scratch. By using the information that pre-trained models have gained from working with huge datasets and tailoring it to specific tasks using smaller datasets, Deep learning proves to be a potent way to address these issues. This work presents a thorough analysis of Deep learning methods used in hand gesture detection tasks, emphasizing their effectiveness, drawbacks, and potential applications. The key components of this approach include data acquisition, pre-processing, model selection, and evaluation. The proposed approach was evaluated using various metrics, including accuracy, precision, recall, and F1 score. The results demonstrated that Deep learning significantly enhanced the models' ability to differentiate between healthy and diseased leaves, with high accuracy and reduced false positives. Moreover, the model's ability to generalize across different plant species and disease types was assessed, highlighting its versatility.**

*Keywords—Hand Gesture recognition, Hand and finger movement analysis, Machine learning for human interaction, Computer-Science, Indian Sign Language(ISL), Deep Learning,VGG16.*

## I. INTRODUCTION

Hand gesture recognition involves understanding and interpreting hand movements to infer human intent or communication. It plays a vital role in enhancing human-computer interaction by enabling intuitive and natural communication channels. Traditional methods for hand gesture recognition relied on handcrafted features and machine learning algorithms, which often struggled with the complexity and variability of hand movements. With the advent of deep learning, particularly convolutional neural networks (CNNs), the field has witnessed significant advancements. This initiative integrates state-of-the-art technology with agricultural knowledge to provide academics, agronomists, and farmers with a helpful tool. Here are the project's objectives: The first step is to gather and pre-process plant leaf photos so that the dataset is varied and of good quality. This will include photographs of all hand gestures showcasing different letters and numbers and the way they are shown as a gesture. The Second objective is to find the best model that suits the complete dataset model and then pick and optimize pre-trained deep learning models to classify the gestures using the selected model. As a third step , we must train and tune the model and its hyperparameters over the selected dataset images to achieve sound recognition using metrics such as accuracy, recall and F1 Score

The Fourth step is to test the system/model thoroughly to see how well it is going to work under different circumstances and make possible recognition of the gestures.

As a last step, we need to design an intuitive interface for the system's launch that will facilitate simple video or image uploads and deliver rapid, accurate hand gesture recognition findings.

The human hand is a marvel of dexterity, capable of performing exact movements and gestures to convey complex meaning. Harnessing this complexity in computational systems is a formidable task, but it opens up a plethora of possibilities for improving user experience and accessibility. By precisely understanding hand movements, computers can understand and respond to human intentions in real time, resulting in seamless and immersive interactions.

Over time, improvements in computer vision, deep learning, and sensor technologies have increased the bar for hand gesture identification. Deep learning models have demonstrated amazing results in reliably identifying and comprehending complicated hand movements, thanks to the training of large quantities of annotated data. Additionally, the performance of gesture recognition systems has improved due to the increased accuracy and sensitivity of hand movement capture made possible by the widespread usage of depth sensors such as Microsoft Kinect and Leap Motion.

This research seeks to address this urgent issue through the use of a sophisticated technique known as the Deep Learning Method. Deep learning, a powerful method that successfully adapts knowledge from one domain to another and often outperforms starting from zero, has become increasingly popular in machine learning. With this approach, image identification, natural language processing, and an increasing number of sensitive fields have experienced great success.

Applications for hand gesture recognition can be found in a variety of fields, including education, healthcare, and entertainment. It enhances the sense of presence and immersion in virtual reality and augmented reality apps by allowing users to manipulate virtual items and traverse immersive settings with natural hand movements. Gesture-based interfaces in healthcare enable hands-free engagement with assistive technologies and medical equipment, enabling

people with disabilities to communicate and navigate their surroundings more successfully.

There are still obstacles in the way of creating powerful and trustworthy hand gesture recognition systems, despite tremendous advancements. Accurate gesture interpretation is hindered by variations in hand sizes, shapes, and orientations as well as surrounding elements like illumination and occlusions. Furthermore, meticulous hardware and software component integration and optimization are necessary to provide real-time performance and scalability across a variety of applications.

## II. LITERATURE SURVEY

Recent years have seen notable progress in hand gesture detection, mostly due to the widespread use of deep learning models in computer vision applications. In order to develop reliable and precise gesture detection systems, researchers have looked into a variety of architectures, datasets, and approaches. This section provides an overview of important research and methods for deep learning models-based hand gesture recognition.

This research work aims to address by putting out a real-time hand-gesture recognition system based on American Sign Language (ASL), this research endeavor seeks to bridge the communication gap between the general public and those with speech and hearing impairments. The researchers used the Vision Transformer Model (ViT)[1], a cutting-edge deep learning architecture renowned for its efficiency in processing visual data, to train the model. The dataset used to train the ViT model included 87,000 RGB samples that showed different combinations of the 29 ASL motions. During one epoch, or training phase, the model's parameters were optimized. There were 2719 batches, or 32 photos per batch, in each epoch.

This paper's research work[2] consists of preprocessing techniques are applied to the MNIST dataset in the first phase. The preprocessed hand gesture image's many crucial properties are then calculated. 24 classes are used in the ASL to represent the alphabetic letters A through Y (excluding J and Z). Nevertheless, as the letters 9 and 25 involve gesture motions related to these letters, there are no cases accessible for J or Z, respectively. The MNIST dataset indicates that there are 34,627 total training instances, of which 80% are for training and 20% are for testing. All of the images have a uniform dimension of $28 \times 28$ pixels and range in grayscale from 0 to 255. The findings show that the recommended methods produced positive outcomes, with a 98.75% categorization recognition accuracy.

This paper[3] is the fact on how cnn can be applied on different datasets because deep learning is currently experiencing a boom, there is a lot of active practical research being done in the field of computer vision. Even though multiple image processing algorithms have employed color and depth cameras to distinguish hand gestures up to this point, it is still challenging to reliably classify movements from distinct subjects. This study's goal is to recognize hand motions by swiftly following the region of interest—in this

example, the hand region—in the visual range with a camera. In this project, a real-time hand gesture detection system using convolutional neural networks has been proposed in this study.

This method[4] by Bader Alsharif entailed extracting features from several input photos that represented ASL motions using ResNet. After that, these characteristics were combined and supplied to the ViT model for additional processing. By using this knowledge transfer approach, which is an excellent example of both multi-image-focused fusion and transfer learning, the ResNet acquired representations can be utilized to the advantage of the ViT model.Our efforts paid off, as the accuracy of the ViT model increased significantly, with performance rising from 88.59% to an astounding 97.09%. This accomplishment highlights how well our suggested system works to improve ASL alphabet categorization problems, especially for the benefit of hearing impaired people. Our method shows great promise for improving speech and hearing aid accessibility and inclusion by utilizing the complimentary qualities of ViT and ResNet.

The goal of this paper by K. Bantupalli[5] is to create a vision-based application that will help sign language users and non-signers communicate more effectively. The model takes advantage of current developments in computer vision and deep learning to extract spatial as well as temporal characteristics from video sequences containing sign language motions. With Inception, a Convolutional Neural Network (CNN), spatial features—static hand locations and configurations—are identified. Meanwhile, a Recurrent Neural Network (RNN) processes temporal characteristics, which indicate dynamic changes across time. The American Sign Language Dataset is used to train the model, giving it the ability to recognize and understand a wide range of sign language motions. By using this method, the program converts gestures from sign language into text, allowing people who don't use sign language to communicate with each other easily.

In this paper by A. Mohan[6] Before creating the sign language gesture prediction system, they worked on improving the overall accuracy and the individual accuracies of the given gestures in this paper. To do this, they took a dataset of previously processed images from Kaggle and applied a CNN model with four layers of relu activation function, three layers of pooling, and a layer of softmax activation. After execution, they obtained an overall accuracy of 99.95%.

The work by Muang[7] a picture is transformed into a feature vector by a transforrn, which is then compared to the feature vectors of a training set of motions by a pattern recognition system. The Perceptron implementation in MATLAB will be the final system. This report presents the findings of studies using thirty-three hand postures. Tests indicate that the system is appropriate for real-time applications and can reach an average recognition rate of 90%.

Another study in this field is by Shivashankar S 2018 [8], who proposed a use of the HSV model, YCbCr model and grayscale model to pre-process the images before feeding

them into the predicting system. The system then computes various properties such as the centroid, gesture area and boundary and uses the results for recognizing the gesture being performed. However, this system didn't account for video processing and failed to recognize alphabets such as J and Z which involved hand movements. The pre-processing stage involved utilizing three different models: the HSV model, the YCbCr model, and the grayscale model. By employing these models, the images were transformed into different color spaces, allowing for better analysis and feature extraction. After the pre-processing stage, the system computed various properties of the gesture, including the centroid, gesture area, and boundary. These properties were then utilized to recognize the specific gesture being performed. The proposed system achieved an accuracy rate of 92.88%, indicating its effectiveness in identifying and classifying gestures accurately. However, it is important to note that this system had certain limitations. As a result, the system failed to recognize gestures that involved hand movements, such as the gestures for the alphabets J and Z. This limitation suggests that the system was not suitable for real-time applications or scenarios where video input was essential for accurate gesture recognition. Furthermore, the study employed a deep Convolutional Neural Network (CNN) model for gesture recognition. CNNs are a type of neural network commonly used for image processing tasks due to their ability to extract hierarchical features from images. Overall, the proposed system showed promising results in gesture recognition, achieving a high accuracy rate. However, its limitations in video processing and failure to recognize certain hand movements indicate the need for further improvements and advancements in the field of gesture recognition systems.

## III. METHODOLOGY

The research methods employed in this project for plant leaf disease detection using deep learning encompass several key steps, each critical to achieving the project's objectives. Here is a detailed description of the research methods

### A. Data Collection and Pre-processing:

Data Acquisition: The project begins by collecting a diverse dataset of plant leaf images. These images should represent various plant species and encompass healthy and diseased leaves. Field surveys, collaboration with agricultural institutions, and publicly available datasets can all contribute to data acquisition.

Data Annotation: Each image in the dataset must be labeled with the corresponding plant species and disease type. Accurate labeling is essential for model training and evaluation.

### B. Data Pre-processing:

Data pre-processing techniques are applied to ensure the dataset's quality and consistency. This involves resizing images to a uniform resolution, removing artifacts, normalising lighting conditions, and augmenting the dataset with techniques like rotation and cropping. Pre-processing helps reduce noise in the data and enhances the model's ability to learn relevant features.

### C. DeepLearning Model Selection:

Model Architecture: Several pre-trained deep learning models, such as VGG16, ResNet50, and Inception, are considered. The choice of the model architecture depends on factors like computational resources, model size, and the specific characteristics of the plant leaf dataset. Fine-tuning a model that has been pre-trained on a large-scale dataset like ImageNet is a common practice.

### D. Model Training and Optimisation:

Loss Function: A suitable loss function, often cross-entropy, is chosen for the classification task. The loss function quantifies the error between predicted and actual class labels.

Validation and Testing: The dataset is typically divided into training, validation, and test sets. The model's performance is monitored on the validation set during training, and it is rigorously evaluated on the test set to assess its generalisation capabilities.

### E. Evaluation and Generalisation:

Performance Metrics: The model's performance is evaluated using accuracy, precision, recall, and F1 score metrics. These metrics provide insights into the model's ability to correctly identify healthy and diseased leaves and minimize false positives and negatives.

The dataset collected forms the training data, which is trained against our image processing model. This model is then saved and used to test the images taken by the camera. The Web interface comes into a requirement when we need a user interface wherein a user needs to upload the captured picture into the front end, and the model is pre-trained by the dataset of images in the back end. The result is generated on the user interface without needing the user to navigate between the training and testing phases of the system. This interface renders an easy and smooth control flow, and the user does not need to know the entire mechanism behind the HGR Hand Gesture Recognition System.

The architecture and design of a project for hand gesture recognition using Deep learning are built around a systematic and efficient framework. The project comprises several essential components, both hardware and software. On the hardware side, server/cloud infrastructure and user devices form the basis for data processing and user interaction. Software components include data pre-processing to ensure data quality and compatibility, model selection and configuration to adapt a pre-trained deep learning model for the specific task, model training and optimisation for fine-tuning and hyperparameter tuning, and model evaluation to gauge the model's performance using various metrics. Additionally, a component dedicated to assessing the model's generalisation abilities ensures its adaptability to different scenarios. Figure 1 shows the modules of the proposed model. The process of the proposed model is shown in Figure 2.
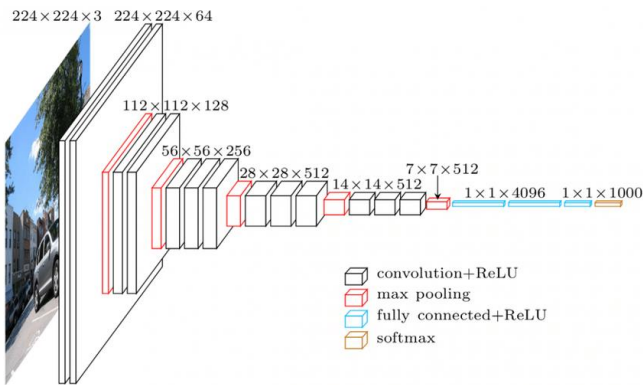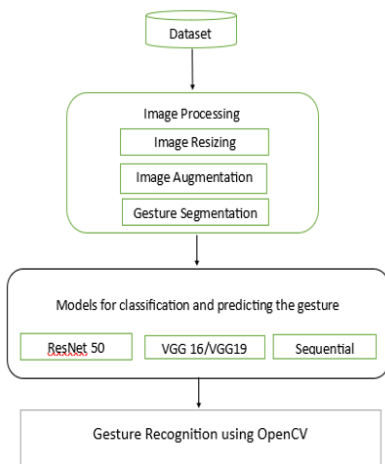
Fig. 1.  The proposed architecture of the model



Fig. 2.  The process of the proposed model

## IV.  RESULTS AND DISCUSSION

This dataset consists of images pertaining to Indian Sign Language, which deaf and hard-of-hearing people in India utilize for communication, which is probably included in the dataset "Indian Sign Language (ISL)". The dataset comprises 42,745 photos included in the collection are 300 x 300 pixel images with three color channels (300 x 300 x 3). These pictures, which are divided into 35 classes, most likely correspond to different Indian Sign Language (ISL) movements or signs. Using the ImageDataGenerator, augmentation techniques are employed to improve the dataset's robustness and diversity. By transforming preexisting photos using techniques including rotation, scaling, shearing, and flipping, augmentation creates new images. By increasing the dataset's variability, this technique helps the model better adapt to the various orientations, locations, and illumination conditions that may arise in real-world situations. In computer vision problems, augmentation is frequently used to reduce overfitting and enhance the model's capacity for generalization. The model gains improved recognition and classification skills even when it is exposed to slightly different photos from the original dataset by being exposed to a greater variety of changes in the training data.

The VGG16 model is explicitly selected as the pre-trained model for the recognition with the help of tailored and optimized learning rates for gesture Recognition. The training process involved fine-tuning of the model using a batch size of 128 and the complete model undergoes a total

of 10 epochs for each each epoch having 201 steps_per_epoch of iterative learning. Using keras. optimizers.Adam (Adam Optimizer) with an enhanced learning rate set to 0.001, The model built for the recognition has exhibited a commendable accuracy of 98.17 %. This is the major contribution and achievement scored under its ability to effectively recognise the hand gestures.
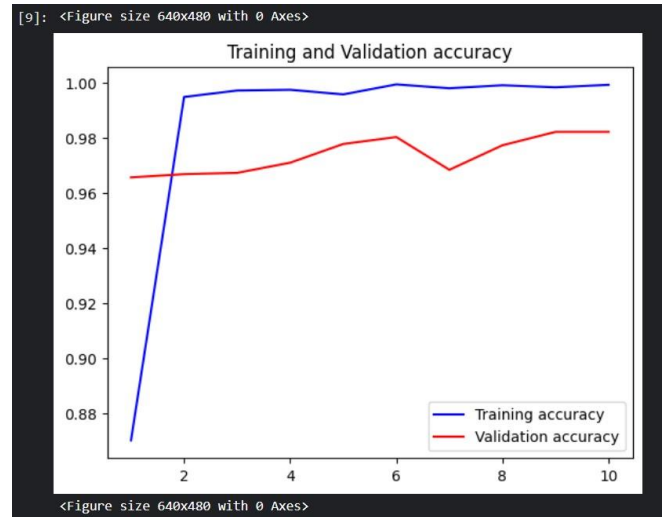


Fig.3 Training and Validation Accuracy of the Model

Figure 3 portrays a line plot graph of both the training and validation accuracies across epochs during the model training phase. The x-axis represents the no of epochs and the y-axis represents the corresponding accuracy values.
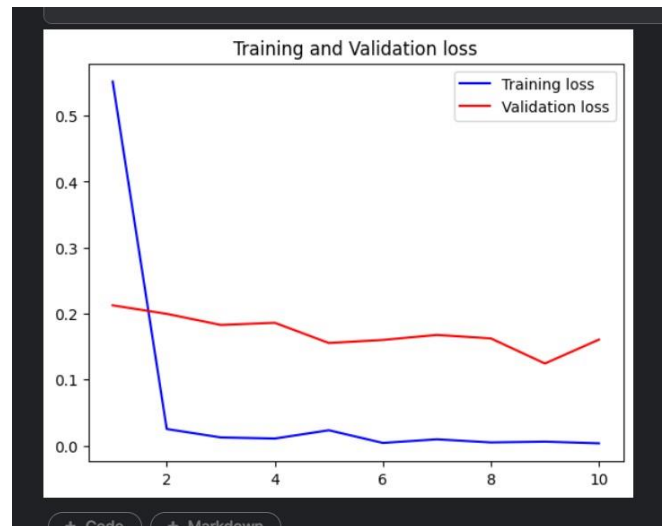


Fig4 Training and Validation Loss of the Model

Figure 4 is the visualization is the loss function that describes the exact model's losses under different phases of running the model.

We used a deep convolutional neural network(VGG16) architecture and hand gesture images to train a model to identify gestures present in new photos. Hand gesture

photographs and a deep convolutional neural network architecture were used to train a model that could recognize gestures in new images. A 99.92% maximum accuracy within the 42745 photos with 35 classes of images in the Indian Sign Language data set indicates that this goal has been achieved. Therefore, using the model without feature engineering, 993 out of 1000 pictures (or 38%) were correctly classified as gestures. Although training the model takes a long time (several hours on a high-performance GPU cluster computer), the classification step is very fast (less than a second on a CPU), which makes it suitable for deployment on mobile devices.

## V. Conclusion

In our study, the model was successful in getting the highest accuracy using the VGG16 model. Our project leveraged achieved impressive results in hand gesture detection by utilizing the VGG16 architecture in conjunction with the capability of deep learning. Our model illustrates the efficiency of deep convolutional neural networks for complicated picture classification tasks, with peak training accuracy reaching an astonishing 99.92% and peak validation accuracy at 98.75%. Our work was well-founded by the VGG16 architecture, which is well-known for its simplicity and efficacy. This allowed us to concentrate on optimizing the model for hand gesture identification using the Indian Sign Language dataset. The tremendous success made during the study, in spite of the high processing demands during model training, highlights the promise of deep learning approaches in promoting inclusion and accessibility for people who use sign languages.

We were able to create a reliable and effective real-time gesture detection system by utilizing deep learning architectures such as VGG16, which opens up new avenues for future study and useful applications in assistive technology. This study emphasizes the significance of ongoing research and innovation in this quickly developing field and the transformative power of deep learning in addressing difficult societal concerns. Many challenges were faced during developing the model some of them are **Variability in motions**: There is a great deal of variation in hand motions with regard to appearance, orientation, lighting, and backdrop clutter. It is difficult for models to generalize well across various instances of the same gesture because of this heterogeneity.
**The intricacy of Motions**: To handle temporal information effectively, dynamic gestures like finger spelling or continuous hand movements need specific strategies.

**Gesture Ambiguity**: It can be difficult for a model to reliably discern between some motions due to their small visual differences or similar look.[p0kb[

REFERENCES

[1] S. N. Reddy Karna, J. S. Kode, S. Nadipalli and S. Yadav, "American Sign Language Static Gesture Recognition using Deep Learning and Computer Vision," *2021 2nd International Conference on Smart Electronics and Communication (ICOSEC)*, Trichy, India, 2021, pp. 1432-1437, doi: 10.1109/ICOSEC51865.2021.9591845.

[2] T. D. Gunvantray and T. Ananthan, "Sign Language to Text Translation Using Convolutional Neural Network," *2024 International Conference on Emerging Smart Computing and Informatics (ESCI)*, Pune, India, 2024, pp. 1-5, doi: 10.1109/ESCI59607.2024.10497209

[3] S. Meshram, R. Singh, P. Pal and S. K. Singh, "Convolution Neural Network based Hand Gesture Recognition System," *2023 Third International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT)*, Bhilai, India, 2023, pp. 1-5, doi: 10.1109/ICAECT57570.2023.10118267

[4] B. Alsharif, M. Alanazi, A. S. Altaher, A. Altaher and M. Ilyas, "Deep Learning Technology to Recognize American Sign Language Alphabet Using Mulit-Focus Image Fusion Technique," *2023 IEEE 20th International Conference on Smart Communities: Improving Quality of Life using AI, Robotics and IoT (HONET)*, Boca Raton, FL, USA, 2023, pp. 1-6, doi: 10.1109/HONET59747.2023.10374775.

[5] K. Bantupalli and Y. Xie, "American Sign Language Recognition using Deep Learning and Computer Vision," *2018 IEEE International Conference on Big Data (Big Data)*, Seattle, WA, USA, 2018, pp. 4896-4899, doi: 10.1109/BigData.2018.8622141

[6] A. Mohan, D. Mohan, S. Vats, V. Sharma and V. Kukreja, "Classification of Sign Language Gestures using CNN with Adam Optimizer," *2024 2nd International Conference on Disruptive Technologies (ICDT)*, Greater Noida, India, 2024, pp. 430-433, doi: 10.1109/ICDT61202.2024.10489158.

[7] Maung, T.H.H., 2009. Real-time hand tracking and gesture recognition system using neural networks. International Journal of Computer and Information Engineering, 3(2), pp.315-319.

[8] Shivashankar S, Srinath S, 'American Sign Language Recognition System: An Optimal Approach', August, I.J. Image, Graphics, and Signal Processing, 2018.