IARJSET



International Advanced Research Journal in Science, Engineering and Technology Impact Factor 8.066 ∺ Peer-reviewed & Refereed journal ∺ Vol. 12, Issue 4, April 2025 DOI: 10.17148/IARJSET.2025.12467

AI-driven hand signs and face feel recognition system

B. Haritha¹, D. Lavanya², K. Jayasree Nagamani³, B. Himaja⁴, A. Vasantha⁵

M.Tech, Asst.Professor, Computer Science & Engineering, Bapatla Women's Engineering College,

Bapatla,India¹

B.Tech, Computer Science & Engineering, Bapatla Women's Engineering College, Bapatla, India²⁻⁵

Abstract: With the growing need for intelligent human-computer interaction systems, recognizing human emotions and interpreting sign language have become essential components in bridging communication gaps. This research presents a unified deep learning-based system that integrates both facial emotion recognition and sign language translation. The proposed model utilizes pre- trained VGG16 and VGG19 architectures to extract high-level spatial features from facial images and sign language gestures. For facial emotion recognition, the FER2013 dataset is used, and real- time emotion prediction is achieved using live webcam input. In parallel, sign language gestures are interpreted using the American Sign Language (ASL) dataset, where the temporal dynamics are captured and processed. The extracted features are used to train classifiers to enhance recognition accuracy. Experimental evaluations demonstrate the effectiveness of the combined approach, showing promising performance in accurately detecting emotions and translating sign gestures. This integrated system offers a valuable tool for enhancing communication, especially for individuals with speech impairments and in emotionally aware interactive systems.

Keywords: Facial Emotion Recognition, Sign Language Translation, VGG16 and VGG19, Deep Learning, Human-Computer Interaction and Real-Time Gesture Recognition

I. INTRODUCTION

In human interaction, emotions serve as the primary medium for expressing internal states. These emotional cues whether conveyed through facial expressions, voice modulation, or physiological signals—play a critical role in how individuals connect with one another and with technology. Facial expressions alone contribute significantly to nonverbal communication, accounting for approximately 55% to 93% of the emotional context in daily interactions. Consequently, analyzing facial expressions provides a rich source of emotional information, making automated Facial Emotion Recognition (FER) a prominent research focus in the field of computer vision.

Automated FER has found widespread applications across various domains, including human-computer interaction, mobile technology, security systems, psychological assessments, healthcare, driver fatigue monitoring, and even robotics, where it supports emotional intelligence in machines. Deep Convolutional Neural Networks (DCNs) have long been the preferred models for FER due to their robustness against minor image transformations such as scaling and translation.

However, recent studies have challenged this assumption, revealing that DCNs are not entirely invariant to such changes and are, in fact, shift-variant.

A primary contributor to this limitation is the down-sampling technique, particularly strided convolutions, which often violate the sampling theorem and lead to aliasing. Aliasing occurs when high-frequency image details are misrepresented as lower frequencies during the down-sampling process, resulting in information loss and distorted features. This can severely affect FER accuracy, causing misclassification of emotional states.

To address this issue, signal processing principles advocate for blurring (low-pass filtering) prior to subsampling to prevent aliasing. Despite this, many contemporary CNN architectures overlook this crucial step. Unlike traditional approaches, the current work introduces an anti- aliasing strategy within a CNN framework specifically tailored to minimize the aliasing effect in FER systems, thereby improving recognition performance and feature preservation.

Sign language serves as a vital mode of communication for individuals who are hearing or speech impaired. However, understanding sign language is typically limited to those closely associated with the deaf community, such as family members or specially trained instructors. Sign language can take the form of both informal gestures and structured, rule-

IARJSET



International Advanced Research Journal in Science, Engineering and Technology Impact Factor 8.066 ∺ Peer-reviewed & Refereed journal ∺ Vol. 12, Issue 4, April 2025 DOI: 10.17148/IARJSET.2025.12467

based expressions that reflect the grammar of spoken languages. Through specific hand movements and facial expressions, signers convey complex thoughts and emotions without speaking, bridging a crucial gap in communication. In many social and professional settings, effective communication is essential to avoid misunderstandings and promote inclusion. One of the major barriers to communication between hearing-impaired individuals and the wider population is the general lack of awareness or understanding of sign language. This communication divide can create challenges in everyday activities—from contacting authorities to participating in economic sectors such as agriculture— simply due to the inability to express themselves clearly.

Traditional approaches to sign language recognition have often relied on sensor-based systems, such as gloves equipped with motion and position detectors. While these systems can offer precision, they are typically expensive and require complex hardware setups, making them less practical for widespread adoption. As a result, vision-based approaches using computer vision and deep learning techniques have gained more popularity due to their non-intrusive nature and scalability.

Despite extensive research, developing a reliable sign language recognition system remains a complex challenge. Sign language is highly dynamic and varies significantly across individuals in terms of speed, hand orientation, duration, and expression. Building a model that can generalize well across different users while maintaining real-time responsiveness requires a careful balance between accuracy and computational efficiency.

This research focuses on developing a vision-based sign language recognition system using Convolutional Neural Networks (CNNs), specifically leveraging pre-trained VGG16 and VGG19 models. These architectures are chosen for their ability to extract robust and high-level spatial features from video frames of hand gestures. The system processes both static and dynamic signs captured via a standard webcam. Each frame undergoes preprocessing and feature extraction using the CNN backbones, followed by classification into the corresponding sign language category.

The objectives of the proposed system are as follows:

- To enhance the precision of sign language classification using deep CNN models;
- To enable real-time detection and recognition of both static and sequential dynamic signs;
- To provide visual feedback of recognized signs along with their corresponding text and confidence levels.

Experimental evaluations show that the CNN-based framework achieves strong classification performance for a variety of sign gestures, demonstrating its potential to serve as an assistive communication tool for the speech- and hearing-impaired community. By integrating VGG16 and VGG19 into a unified model, this approach contributes toward bridging the gap between signers and non-signers in everyday interactions, thus supporting inclusivity and digital accessibility.

Related work: -

C. Raghavachari and colleagues focused their study on enhancing communication tools for individuals with hearing and speech impairments through the interpretation of hand, facial, and body gestures. They proposed a vision-based finger-spelling system using Convolutional Neural Networks (CNNs). Although their model surpassed the performance of architectures like ResNet and VGG16 in some respects, the overall accuracy remained modest. Interestingly, it performed comparably to InceptionV3 in certain cases but still highlighted the need for further optimization.

A. M. Arjun and team explored the significance of non-verbal cues, particularly hand gestures, in conveying meaning without relying on facial expressions. Their work emphasized that hand gestures alone can be powerful communication tools, and when integrated with other non-verbal elements, they can convey complex messages. They implemented an image classification system using CNNs along with TensorFlow and OpenCV, but their work was constrained to a relatively small dataset. The system was designed to classify emotional states purely based on hand movements, without incorporating facial features.

The research conducted by Sundar B et al. spanned multiple domains including communication for the deaf-mute community, healthcare applications, home automation, and human-robot interaction. They utilized Google's MediaPipe framework to detect and extract hand landmarks, and applied it to American Sign Language (ASL) recognition. Their model demonstrated high accuracy in recognizing all 26 alphabet gestures, reaching up to 99% using deep learning techniques. However, while their method accurately translated gestures into text, integrating real-time message display remains an area for improvement.

L. V. B., S. K. B., P. H., and S. Abhishek developed a deep learning model that translates ASL signs into either text or



International Advanced Research Journal in Science, Engineering and Technology

Impact Factor 8.066 $\,\,st\,$ Peer-reviewed & Refereed journal $\,\,st\,$ Vol. 12, Issue 4, April 2025

DOI: 10.17148/IARJSET.2025.12467

speech. The researchers evaluated the performance of multiple CNN architectures including VGG16, ResNet, and AlexNet. The results showed that ResNet achieved the highest accuracy of 96.9%, followed by AlexNet at 95.87%, and VGG16 at 93.2%. Despite these strong results, some underfitting issues were noted, highlighting the need for model tuning and dataset expansion.

The paper by Simonyan et al. introduced an innovative approach to image classification using deep CNNs. While the method showed promising results, the evaluation was limited to a subset of the ImageNet dataset. The study did not assess the impact of dataset diversity or generalizability across other datasets, leaving room for further investigation. Aloysius et al. proposed a vision-based system for continuous sign language recognition, which largely depended on accurate hand detection and tracking. However, their method struggled with consistency under varied lighting and background conditions. Moreover, the use of a K-Nearest Neighbors (KNN) classifier may not have been optimal, as it showed limitations in both speed and accuracy compared to deep learning approaches.

Recent advancements in artificial intelligence have significantly contributed to the development of systems aimed at assisting individuals with hearing impairments through sign language recognition (SLR). Many researchers have explored gesture prediction techniques to improve communication accessibility for this community. Despite the notable progress in this area, challenges remain that must be addressed to enhance the accuracy, speed, and robustness of these systems.

This section reviews contemporary research efforts focusing on SLR using both sensor-based and vision-based deep learning methods. The literature reveals a wide variety of techniques employed to recognize gestures from video inputs.

For instance, one study utilized Hidden Markov Models (HMMs) in combination with Bayesian Network Classifiers and Gaussian Tree-Augmented Naive Bayes Classifiers to recognize facial expressions from video frames. These hybrid probabilistic models aimed to increase accuracy by considering the temporal and statistical dependencies between gestures.

Francois et al. proposed a method for human posture recognition using a combination of 2D and 3D appearance-based models. They applied Principal Component Analysis (PCA) to extract silhouettes from video captured by a static camera and then constructed 3D models for posture analysis.

However, this method encountered limitations, particularly with ambiguous transitional gestures, which reduced the model's overall prediction accuracy.

In another approach, Neural Networks were applied to analyze video sequences by extracting visual features and classifying them. Although effective in learning complex patterns, neural networks face challenges such as accurate hand tracking, subject-background segmentation, occlusions, illumination changes, and variability in movement and positioning.

Nandy et al. addressed these challenges by segmenting video datasets, extracting visual features from each segment, and classifying them using Euclidean Distance and K-Nearest Neighbors (KNN). Their method focused on pattern recognition through feature distance measures.

Similarly, Kumud et al. proposed a framework for continuous Indian Sign Language recognition. Their process involved extracting individual frames from video sequences, performing pre- processing steps, and identifying key frames for gesture analysis. Videos were first converted into RGB frames of uniform dimensions, followed by skin color segmentation using the HSV color space. The segmented images were then binarized, and key frames were selected by analyzing the gradient differences between successive frames. From these key frames, orientation histograms were computed as feature descriptors. The recognition stage involved comparing features using multiple distance metrics, including Euclidean, Manhattan, and Chessboard distances.

Overall, while numerous techniques exist for sign language gesture recognition, ongoing research is focused on improving system performance in terms of real-time processing, user independence, and recognition under diverse environmental conditions.

II. METHODOLOGY

The proposed system for vision-based sign language recognition involves a comprehensive pipeline starting with data collection, where American Sign Language (ASL) gestures are captured in the form of images and video sequences representing both static signs (e.g., letters) and dynamic signs (e.g., words or short phrases). The dataset is either sourced

446



International Advanced Research Journal in Science, Engineering and Technology Impact Factor 8.066 ∺ Peer-reviewed & Refereed journal ∺ Vol. 12, Issue 4, April 2025 DOI: 10.17148/IARJSET.2025.12467

from public ASL repositories or custom-built using recorded video frames via a webcam. After data acquisition, the image preprocessing stage is applied to enhance the quality and consistency of the inputs—this includes resizing images to a standard dimension, converting them to grayscale or RGB format, normalizing pixel values, removing background noise, and applying histogram equalization to enhance contrast. In cases of hand segmentation, skin color detection or MediaPipe hand landmark extraction is employed to isolate the hand region. The preprocessed data is then split into training, validation, and testing subsets, typically in a 70:15:15 ratio, to ensure balanced and fair model evaluation across unseen data.



Input Image

Figure 1 Face emotion detection model architecture



Figure 2 Sign Language model architecture

Following preprocessing, the system proceeds to model building using deep learning architectures such as Convolutional Neural Networks (CNNs) and pretrained models like VGG16 and VGG19. These models are fine-tuned for sign recognition by modifying the top classification layers to match the number of ASL classes. Model training is carried out using the training dataset, with optimization techniques such as Adam optimizer and categorical cross-entropy loss function.

Techniques like data augmentation and early stopping are used to prevent overfitting and improve generalization. Once trained, the models are evaluated on the test set using performance metrics such as accuracy, precision, recall, F1-score, and confusion matrix analysis. After achieving satisfactory results, the best-performing model is integrated with OpenCV to perform real-time gesture recognition via webcam. The system captures live frames, applies the same preprocessing pipeline, and feeds the processed image into the trained model to predict the gesture along with its predicted label and confidence score, displaying the results in real time on the screen. This practical evaluation helps validate the model's effectiveness in real-world scenarios and ensures its usability by the hearing- and speech-impaired community.

Datasets: - For facial emotion recognition, the CK+ (Cohn-Kanade Plus) dataset has been utilized. Although it contains only 981 labeled images, it is well-regarded for its high-quality annotations and clear facial expressions, making it a strong candidate for deep learning tasks despite its small size. The dataset includes a diverse range of facial emotions, enabling effective training and evaluation of emotion classification models and same use sign language dataset also from Kaggle. **Image Preprocessing**:



International Advanced Research Journal in Science, Engineering and Technology Impact Factor 8.066 ∺ Peer-reviewed & Refereed journal ∺ Vol. 12, Issue 4, April 2025 DOI: 10.17148/IARJSET.2025.12467

To ensure uniform and high-quality input for the system, the preprocessing phase begins by resizing all images and frames to a standard dimension for consistency. Depending on the data type, images are either converted to grayscale or retained in RGB format to maintain essential features.

Brightness and contrast are adjusted through techniques like histogram equalization to improve visibility of key patterns. To reduce variations and speed up processing, pixel values are normalized to a common scale, typically between 0 and 1. For hand gesture recognition, methods like background filtering, skin tone segmentation in the HSV color space, or region isolation are applied to focus on the hand area. In the case of facial emotion data, the face is detected and cropped to eliminate unnecessary background details. These steps ensure that each input is cleaned, focused, and standardized before entering the classification pipeline, thus supporting accurate and reliable recognition.

Data splitting: - To evaluate the model's performance effectively, the dataset is divided into two main subsets: 80% is allocated for training, while the remaining 20% is reserved for testing. The training set is used to teach the model how to recognize patterns, such as facial expressions or hand gestures, by learning from labeled examples. This phase helps the model build an understanding of the features that define each class. On the other hand, the testing set contains unseen data that the model has never encountered before. It is used to assess how well the trained model performs on new inputs and to ensure that it can generalize its learning beyond the training data. This 80-20 split strikes a balance between providing the model with enough information to learn while still preserving a portion of the data for an unbiased evaluation of its real-world effectiveness.

Model Training: - CNN

model:-

In this project, two separate Convolutional Neural Network (CNN) models are trained—one for recognizing facial emotions and another for interpreting sign language gestures. Each CNN model follows a structured pipeline to classify the input images. First, the input image (either a facial expression or a hand sign) is passed through multiple convolutional layers. These layers apply filters to detect essential features such as edges, shapes, and textures. As the image moves deeper into the network, more complex patterns and high-level features are extracted, which help distinguish between different emotions or hand signs.

After the convolutional operations, pooling layers are used to reduce the spatial dimensions of the feature maps, which not only decreases computation time but also helps the model focus on the most relevant features. These condensed features are then passed through fully connected layers that act as a decision-making network. These layers analyze the extracted patterns and assign a probability score to each class. The final output layer uses a softmax function to determine the most likely class label—for example, predicting whether a face shows happiness or sadness, or whether a hand sign represents the letter 'A' or 'B'.

By learning from a large number of labeled examples during training, the CNN models become capable of generalizing well to new, unseen images. The result is accurate and efficient classification for both facial emotion recognition and sign language translation, making the system practical for real-time applications.

VGG- 16 Model:-

In this project, VGG16 plays a key role in identifying facial emotions and hand gestures by analyzing image features at different levels. VGG16 is a deep convolutional neural network consisting of 16 layers, mainly using 3x3 convolutional filters that are stacked together to extract detailed spatial features from the input images.

For both emotion and sign classification tasks, the input image is first resized and passed through multiple convolution layers. These layers help in detecting features starting from basic edges and textures to more complex patterns like facial muscle positions or hand shapes. After every few convolutional layers, max-pooling is applied to reduce the spatial dimensions and retain the most important information.

As the image goes deeper into the network, VGG16 captures high-level abstractions that are crucial for distinguishing between different classes. For example, it can tell the difference between a smile and a frown, or between a thumbs-up and a fist, based on the learned patterns.

Once all feature maps are generated, they are flattened and sent to dense layers which perform the classification. A softmax activation is used in the final layer to output the probabilities for each emotion or sign class. The class with the highest probability is considered the predicted label.

VGG16's architecture helps the model learn more robust and general features, especially when pre- trained weights are used. This results in better performance and faster convergence, making it a reliable choice for both facial emotion recognition and hand sign detection tasks.

Results and Analysis: -

© <u>IARJSET</u> This work



International Advanced Research Journal in Science, Engineering and Technology Impact Factor 8.066 ∺ Peer-reviewed & Refereed journal ∺ Vol. 12, Issue 4, April 2025 DOI: 10.17148/IARJSET.2025.12467

IARJSET

The performance of both the emotion detection and sign language recognition models was evaluated using two different architectures: a custom-built CNN and the pre-trained VGG16 model. For the sign language recognition task, the VGG16 model achieved a remarkable accuracy of **95%**, significantly outperforming the custom CNN model which reached **80%**. Similarly, in the case of emotion recognition, VGG16 once again delivered superior results with an accuracy of **96%**, compared to **82%** obtained using the CNN model. These results clearly demonstrate the advantage of using deeper, pre-trained networks like VGG16, which are capable of learning more abstract and high-level features from the image data.



In addition to the accuracy metrics, the training and validation loss plots further validate the robustness of the VGG16 models. **Figure 3** shows the loss curve for the sign language recognition model using VGG16, and **Figure 4** displays the loss trend for the emotion recognition model. Both plots reveal a smooth decrease in training and validation loss, indicating good convergence and minimal overfitting. The consistency in performance across training and validation sets confirms that the VGG16 models not only learn effectively but also generalize well to unseen data. These results underline the importance of using deep pre-trained models, especially when working with complex tasks like facial emotion and hand gesture classification.



International Advanced Research Journal in Science, Engineering and Technology Impact Factor 8.066 ∺ Peer-reviewed & Refereed journal ∺ Vol. 12, Issue 4, April 2025 DOI: 10.17148/IARJSET.2025.12467

IARJSET

Model	Task	Accuracy
VGG16	Sign Language	95%
CNN	Sign Language	80%
VGG16	Emotion Recognition	96%
CNN	Emotion Recognition	82%
70.1.1.1		

Table 1

The accuracy comparison and loss plots confirm that the VGG16 model consistently outperforms the CNN model, demonstrating its capability to better generalize from complex datasets. The accuracy and loss trends further support that deeper, pre-trained models like VGG16 are highly beneficial for tasks such as sign language recognition and emotion detection. **Table 1** illustrates the comparison between the accuracy of the two models, highlighting the significant performance improvement with VGG16. This reinforces the effectiveness of using well-established deep learning models in real-world applications for gesture and emotion recognition.



Model Accuracy Comparison for Emotion and Sign Language Recognition

Here is the count plot comparing the accuracy of the models used for both sign language and emotion recognition. The plot visually represents the performance of each model in terms of accuracy, with VGG16 outperforming the CNN model in both cases. Specifically, the VGG16 model achieved 95% accuracy for sign language recognition and 96% for emotion recognition, while the CNN model achieved 80% accuracy for sign language and 82% for emotion recognition.

III. CONCLUSION

This research successfully demonstrates the development of an integrated deep learning-based system that performs both facial emotion recognition and American Sign Language (ASL) translation. Leveraging pre-trained VGG16 models, the system achieved impressive accuracy rates of 96% for emotion recognition and 95% for sign language recognition, significantly outperforming traditional CNN models which attained 82% and 80% respectively. The use of advanced image preprocessing techniques, combined with the strength of transfer learning, enabled the models to learn complex patterns efficiently and generalize well to unseen data. Real-time implementation using webcam input further validated the system's practical applicability, offering a robust and effective solution for enhancing human-computer

IARJSET



International Advanced Research Journal in Science, Engineering and Technology

Impact Factor 8.066 $\,\,st\,$ Peer-reviewed & Refereed journal $\,\,st\,$ Vol. 12, Issue 4, April 2025

DOI: 10.17148/IARJSET.2025.12467

interaction, particularly benefiting individuals with hearing and speech impairments.

Moving forward, the system can be expanded to include a wider variety of sign languages beyond ASL, as well as additional emotions to make the interaction more comprehensive. Incorporating temporal models such as LSTM or 3D-CNNs could further improve recognition of dynamic gestures and emotional transitions in videos. Moreover, integrating voice-to-sign and sign-to-voice capabilities could create a fully bidirectional communication platform. Deployment on mobile or edge devices can also make the solution more accessible for real-world use, particularly in assistive technology for special needs communities.

REFERENCES

- [1]. R.C. Gonzalez, R.E. Woods, Digital Image Processing, 2nd edn. (Prentice Hall, New Jersey, 2008), p. 693
- [2]. M.J.Z. Omar, M.H. Jaward, A review of hand gesture and sign language recognition techniques. Int. J. Mach. Learn. Cyber. 10, 131–153 (2019)
- [3]. C. Raghavachari and G. Shanmugha Sundaram, "Deep Learning Framework for Fingerspelling System using CNN," 2020 International Conference on Communication and Signal Processing (ICCSP), Chennai, India, 2020, pp. 469-473, doi: 10.1109/ICCSP48568.2020.9182155.
- [4]. M. Arjun, S. Sreehari and R. Nandakumar, "The Interplay Of Hand Gestures And Facial Expressions In Conveying Emotions A CNN–BASED APPROACH," 2020 Fourth International Conference on Computing Methodologies and Communication (ICCMC), Erode, India, 2020, pp. 833-837, doi: 10.1109/ICCMC48092.2020.ICCMC-000154.
- [5]. Sundar, B., and Bagyammal, T. (2022). American Sign Language Recognition for Alphabets Using MediaPipe and LSTM.Procedia Computer Science, 215, 642–651. <u>https://doi.org/10.1016/j.procs.2022.12.066</u>
- [6]. L. VB, S. KB, P. H, S. Abhishek and A. T, "An Empirical Analysis of CNN for American Sign Language Recognition," 2023 5th International Conference on Inventive Research in Computing Applications (ICIRCA), Coimbatore, India, 2023, pp. 421-428, doi: 10.1109/ICIRCA57980.2023.10220822.
- [7]. Simonyan, K. and Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556.
- [8]. N. Aloysius and M. Geetha, "Understanding vision-based continuous sign language recognition", Multimed Tools[9] S. Akshay, T. K. Mytravarun, N. Manohar and M. A. Pranav, "Satellite Image Classification for Detecting Unused Landscape using CNN," 2020 International Conference on Electronics and Sustainable Communication Systems (ICESC), Coimbatore, India, 2020, pp. 215-222, doi: 10.1109/ICESC48915.2020.9155859.
- [9]. Mrinal Waila, "Object Detection vs. Image Classification vs. Keypoint Detection"
- [10]. Niklas Donges,"A Guide to Recurrent Neural Networks: Understanding RNN and LSTM Networks" https://www.turing.com/kb/recurrent-neural-networks-and-lstm

BIOGRAPHY



B. HARITHA, M.Tech, Asst. Professor, Department of Computer Science & Engineering, BWEC, Andhra Pradesh, INDIA





451





K. JAYASREE NAGAMANI [B.Tech], Student, Department of Computer Science & Engineering, BWEC, Andhra Pradesh, INDIA



B. HIMAJA [B.Tech], Student, Department of Computer Science & Engineering, BWEC, Andhra Pradesh, INDIA



A.VASANTHA [B.Tech], Student, Department of Computer Science & Engineering, BWEC, Andhra Pradesh, INDIA