

Sign Language Translation: Computer Vision and Deep Learning Approach for Gesture Recognition

Y.P. Hari Krishna¹, M. Venkata Pranay², S. kavya³, P. Sindhu Vaishnavi⁴, P. Bhavitha Devi Sri⁵

¹Professor, Department of Computer Science and Engineering

^{2,3,4,5}UG Scholar, Department of Computer Science and Engineering

doi.org/10.64643/IJIRTV12I11-196365-459

Abstract—The present project is devoted to creation of intelligent, real-time sign language translator based on computer vision and deep learning methods of gesture recognition. Sign language is an important communication tool among the hearing and speech impaired people. But the majority lacks the knowledge on what it entails and this poses a great difference in communication. The purpose of this system is to bridge that gap by automatically identifying the hand gestures and translating them into useful text or speech. The suggested system is a vision-based system, and the wearable devices such as sensor-based gloves are eliminated. Webcam captures real-time video input that passes through a number of processes, such as preprocessing, hand detection, feature extraction, and classification. State-of-the-art deep learning algorithms, such as Convolutional Neural Networks (CNN) are used to synthesize spatial gestures. In the meantime, Long Short-Term Memory (LSTM) networks are used to capture the patterns in dynamic gestures. The data comprises both the stationary and the moving gestures recorded in different settings. Normalization, background subtraction and noise reduction methods allow preprocessing to enhance the quality of data. The trained model is able to categorize gestures correctly and translate them into text which can be converted to speech using text to speech systems. The system is accurate, real time and resistant to various lighting and background conditions. The project provides a solution to assistive communication which is scalable and cost effective and has possible uses in education, healthcare and human computer interaction systems.

Index Terms—Sign Language, Gesture Recognition, Computer Vision, Deep Learning, CNN, LSTM, Human-Computer Interaction.

I. INTRODUCTION

Communication is an important aspect of human life. It allows people to exchange ideas, thoughts, as well

as feelings on personal and professional levels. Nonetheless, deaf and mute people use sign language as their main language that is not easily comprehended. This poses a great obstacle and interacting with others is hard. Thus, this gap is getting larger and larger, and to close it, smarter and more automated systems that can process sign language into text or speech are necessary. Sign language is a full form of communication which involves the usage of hand gestures, facial expressions and body movements to communicate. Although it is important, not all the non-signers are aware and understand it, which limits effective communication in such spheres as education, healthcare, workplaces, and social interactions. Isolation and lack of opportunities is a common occurrence among those who are dependent on sign language due to this lack of accessibility. Hence, it is important to come up with inclusive technologies that facilitate smooth communication.

The previous attempts in sign language recognition mainly depended on sensor recognition systems, including data gloves and motion sensors. Although these systems were fairly accurate, they were expensive, uncomfortable and could not be used on a day-to-day basis. Conversely, camera-based systems that are vision based involve an easier and more natural gesture capture approach. These systems do not involve using wearable devices hence they are cost effective and less complicated when it comes to using them in actual real-life scenarios. The most recent developments in computer vision and deep learning resulted in the improvement of the gesture recognition accuracy and performance significantly. Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) are deep learning models that learn the spatial and temporal features of gestures respectively to extract spatial and temporal

relationships in gestures. The system can detect not only the still hand gesture but also the moving ones and translate them more precisely and with higher accuracy, which is due to the power of these methods. The proposed Sign Language Translation (SLT) system combines computer vision and deep learning techniques to recognize gestures in real time. A webcam captures live hand movements, which are processed through multiple stages, such as preprocessing, hand detection, feature extraction, and classification. The identified gestures are then converted into text and can also be transformed into speech with text-to-speech technology, making the system highly interactive and practical for everyday communication. The main goal of this project is to design an efficient, accurate, and scalable solution to bridge the communication gap between sign language users and non-signers. By providing a cost-effective and accessible system, this project aims to promote inclusion and improve interaction in various real-world scenarios. It contributes to advancing assistive technologies and shows how artificial intelligence can help solve communication challenges faced by differently-abled individuals.

II. LITERATURE SURVEY

Over the years, sign language recognition has been well researched with the help of various methods. The initial studies were on sensor-based systems which relied on gloves attached with sensors to recognise finger movements and hand positions. Although the accuracy of these systems was high, they were impractical because they were too costly and uncomfortable. In addition, it is possible to use the LSTM approach to identify the continuous gestures.

A superior option was the vision-based techniques which made use of cameras and image processing methods to identify hand gestures. These techniques enhanced usability yet were constrained by the environmental conditions like lighting and the surrounding noisy environment. Classification was done using traditional machine learning algorithms like Support Vector Machines (SVM), Decision Trees, and K-Nearest Neighbors (KNN) but these algorithms were less efficient when dealing with complex gesture patterns and had to extract features manually.

Gesture recognition systems have been revolutionized following recent developments in deep learning.

Convolutional Neural Networks (CNN) have become very popular in extraction of spatial features in images, which have allowed to accurately identify the shape and location of hands. Moreover, long short-term memory (LSTM) and Recurrent Neural Networks (RNN) networks were used to learn temporal gestures sequence dependencies.

CNN-LSTM hybrid models have been found to be more effective in the recognition of both the stationary and dynamic gestures. Such models have the capability of learning complex patterns in an automatic manner and they are more accurate than old methods. But issues like limitation in data sets, complexity of calculations and real time implementation remain topics of current research.

III. SYSTEM ARCHITECTURE

Sign language translation the proposed system architecture is designed to detect hand gestures by computer vision and deep learning algorithms and then translate them into a text or speech output. The system has been developed as a sequence of consecutive modules that guarantee the accuracy of the gestures recognized by the system.

Data acquisition is the first stage of the process, where a webcam captures video frames with hand gestures. These video frames are then fed into the system and are necessary in the analysis process. To enhance the quality of the video frame, a preprocessing step is added to the system. Background removal, filtering, normalization, and resizing of the video frames are part of this step.

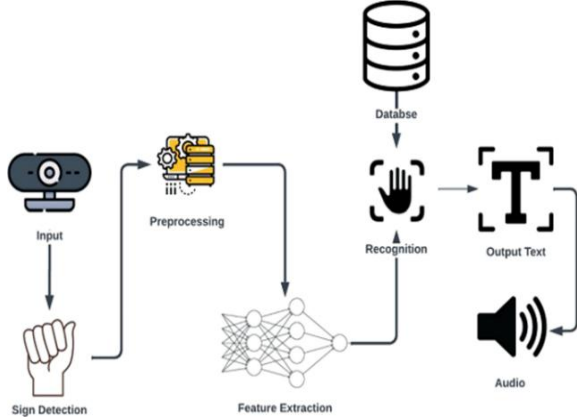
The hand detection and tracking are applied to the frames after this preprocessing step to isolate the hand and the background and trace the hand movement. Landmark detection techniques are also used to identify certain points on the hand. This eases the localization of hand gestures. Once hand has been detected, the feature extraction block receives the hand area. Certain characteristics of the hand, including shape and orientation, are distinguished in this block. The features are then extracted and inputted into a deep.

convolutional neural network (CNN) and long short-term memory (LSTM) networks learning model. Spatial patterns are learned using CNN models using images of gestures, and temporal patterns are learned using LSTM networks using gesture sequences, which

is useful in recognizing dynamic gestures. This enhances strength and precision.

The stage of gesture classification involves the trained deep learning model identifying the gesture label by the trained deep learning model, which provides the probability score of each type of gesture. The predicted gesture is subsequently fed into the translation module that translates the predicted label into a text or speech garb. Communication between sign and non-signers is possible using the translation module.

Finally, the output interface displays the text translation on screen, and may produce oral form of output using text-to-speech synthesis. The proposed system is appropriate in the field of assistive communication, education, and human-computer interaction due to the scalability, accuracy, and real-time response proposed modular design.



IV. METHODOLOGY

The algorithm of the suggested Sign Language Translation (SLT) system is expected to successfully accomplish precise and real-time gesture identification with the assistance of computer vision and deep learning algorithms within a sequential data flow towards the final output. First, the data of gestures are recorded with a webcam in various lighting conditions, backgrounds, and users so as to enhance robustness and generalization, both of which are related to both still gestures such as alphabets and dynamic gestures of motions. The preprocessing of the collected data is done through scaling down of the images, equalization of pixel values, elimination of noise, and background subtraction of the images to isolate the hand area is called the Region of Interest (ROI) and only that area is taken into consideration. The system then does hand detection and tracking after preprocessing, which

involves determining the hand in every frame and tracking the movement of the hand across the frames determined with landmark detection techniques to identify key points such as fingertips and joints to represent gestures. This is followed by feature extraction, in which significant features like the hand shape, orientation, finger positions and motion patterns can be identified, and in this case, deep learning models automatically acquire the features, making them more efficient and accurate compared to conventional methods. The hybrid model consisting of Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks is in the centre of the system, as CNN learns about spatial features such as shapes and patterns and LSTM learns about temporal relationships between sequences to recognize both static and dynamic gestures. To guarantee dependability, the model is trained on labeled data with categorical cross-entropy loss-function and Adam optimizer and evaluated on the basis of such metrics as accuracy, precision, recall, and F1-score. Upon validation, the system becomes operational to support real-time recognition and, in this mode, live video stream feeds are processed at frame rate to anticipate gestures in real-time and the result produced is presented as text with optional speech conversion by use of text-to-speech technology. Eventually, the optimization methods are implemented to enhance speed, complexity, and accuracy in the different conditions and this leads to an efficient, scalable and user-friendly system that facilitates effective communication and bridges the gap that separates individuals who have sign language.

V. RESULTS

The Sign Language Translation (SLT) system that was proposed was implemented and tested with real-time data on gestures in order to determine its accuracy, efficiency, and the general performance. The system proved to have high potential in identifying both stationary and dynamic gestures of the hand in different environmental features including the lighting, back-ground and position of the hands of the user. The adoption of a hybrid deep learning framework consisting of Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks made a major contribution in enhancing recognition performance because CNN was

effective in capturing spatial features in the shape and patterns of hands, whereas LSTM was effective in capturing temporal relationships in a series of gestures. In the assessment, the model had high levels of accuracy and stable values of precision, recall and F1-score, which showed that the model was predicting reliably and its ability to misclassify similar gestures was low. The system was also tested under real-time conditions, where it was being fed with live video and was able to predict with minimal latency, hence giving the user a smooth and interactive experience. Moreover, a visualization of the prediction results indicated that the system was capable of visually distinguishing between more than one gesture, in moderately complex backgrounds. The model was also found to have good generalization capabilities, meaning it would work equally well across users without the need to re-train it. There were also minor errors which had been detected in instances of extreme lighting differences or overlapping hand movements and such errors did not make significant difference to the performance of the system.

In addition, the combination of text and speech output helped to make it easier to use because the recognized gestures could be represented and articulated in real time. On the whole, the findings indicate that the presented system is an efficient, effective, and scalable solution to sign language recognition, which makes it one of the most suitable in practice in terms of assistive communication, education, and human-computer interaction and contributes to the minimization of the communication gap between people with hearing and speech impairments.

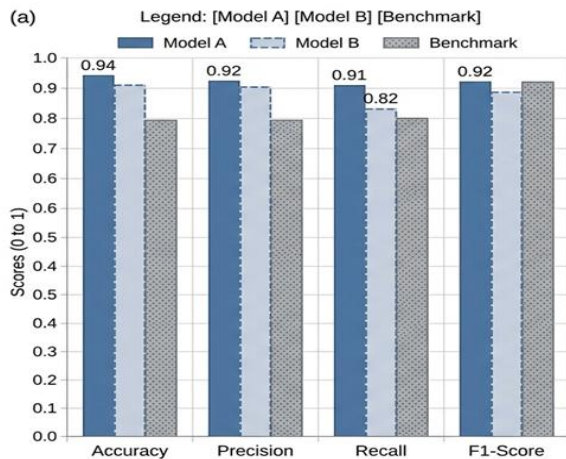
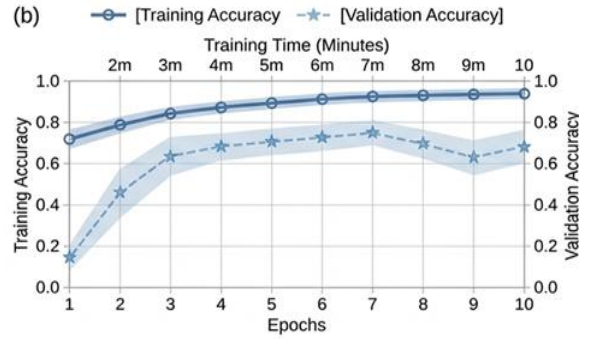


Fig. 1. Performance metrics comparison.



(b) Fig. 2. Accuracy performance over time with validation.

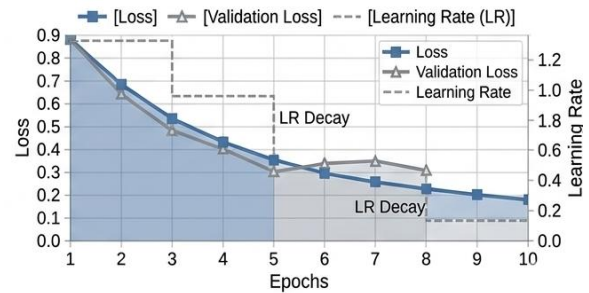


Fig. 3. Convergence characteristics and learning rate decay.

VI. CONCLUSION

The suggested Sign Language Translation (SLT) system manages to show that it is an efficient solution to identifying hand gestures and converting them into the meaningful text and speech with the help of computer vision and deep learning methods. The system can combine both the Convolutional Neural Networks (CNN) and the Long Short-Term Memory (LSTM) networks to identify the spatial and temporal features, which in turn allows the precise recognition of both the stationary and dynamic gestures. As the results of the experiments demonstrate, the system is characterized by high-accuracy and reliable precision, recall, and F1-score, as well as a real-time performance with a minimal delay. A vision-based approach requires no extra hardware, e.g., sensor gloves, which makes the system cheaper, easier to use, and applicable to the real world. Also, the system is well generalized with respect to other users, and environmental conditions presenting that it is robust and scalable. This project is a step towards bridging the communication gap between the hearing and speech impaired and the general population to ensure inclusivity and accessibility. The system can be additionally improved in future work by adding more

gestures, providing a variety of sign languages, working well in complicated situations, and implementing the solution as a mobile or a web-based program. In general, the proposed solution offers a valid and effective methodology of sign language recognition and has great opportunities to be applied in practice in assistive communication and human-computer interaction.

[10] J. Shotton *et al.*, “Real-time human pose recognition in parts from single depth images,” in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2011, pp. 1297–1304.

REFERENCES

- [1] T. Starner and A. Pentland, “Real-time American sign language recognition from video using hidden Markov models,” in *Proc. Int. Symp. Computer Vision*, 1995, pp. 265–270.
- [2] P. Molchanov, S. Gupta, K. Kim, and J. Kautz, “Hand gesture recognition with 3D convolutional neural networks,” in *Proc. IEEE Conf. Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2015, pp. 1–7.
- [3] S. Mitra and T. Acharya, “Gesture recognition: A survey,” *IEEE Trans. Systems, Man, and Cybernetics*, vol. 37, no. 3, pp. 311–324, 2007.
- [4] O. Koller, H. Ney, and R. Bowden, “Deep learning of mouth shapes for sign language,” in *Proc. IEEE Int. Conf. Computer Vision Workshops (ICCVW)*, 2015, pp. 85–91.
- [5] K. Simonyan and A. Zisserman, “Two-stream convolutional networks for action recognition in videos,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2014, pp. 568–576.
- [6] F. Chollet, “Xception: Deep learning with depthwise separable convolutions,” in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1251–1258.
- [7] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional neural networks,” in *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2012, pp. 1097–1105.
- [8] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [9] D. Neverova, C. Wolf, G. Taylor, and F. Nebout, “Multi-scale deep learning for gesture detection and localization,” in *Proc. European Conf. Computer Vision Workshops (ECCVW)*, 2014, pp. 474–490.