

SIGN LANGUAGE TRANSLATION BASED ON HAND GESTURES CONVERTING INTO TEXT AND SPEECH IN ENGLISH AS WELL AS TAMIL USING RANDOM FOREST ALGORITHM

Gayathri M^{#1}, Kaviya K^{#2}, Priyanka S^{#3}, Dr. Madhurikha S^{#4}

^{#1,2,3}Department Of Computer Science And Engineering,

^{#4}Associate Professor / AI&DS

Jaya Sakthi Engineering College, Chennai, 602024, Tamil Nadu, India

gayathri30kavi@gmail.com^{#1}, kaviyak98240@gmail.com^{#2}, js.priyanka06@gmail.com^{#3},

madhurikha@gmail.com^{#4}

Abstract— A real-time system that can decipher sign language from a live webcam stream is presented by the Sign Language Conversion Project. Using the powerful landmark identification features of the Mediapipe library, the project takes important data out of every frame, including hand landmarks. Following their detection, these landmark coordinates are collated and saved in a CSV file for later examination. This landmark data is used to train a Random Forest Classifier, which uses machine learning techniques to identify different sign language patterns. The trained model predicts the sign language class and its probability in real-time as it processes the camera data. In order to improve accessibility, the system offers English translation as for the identified signs, which are accompanied by voice translations into Tamil and English. The video feed has these translations superimposed on it, providing real-time translations and visual cues to enhance accessibility for individuals who are deaf or hard of hearing, allowing them to communicate more effectively with others whom they may not understand sign language.

Keywords— Random forest algorithm, Mediapipe, sign language Machine learning.

I. INTRODUCTION

The sign Language Conversion Project is to meet the critical requirement of providing deaf and mute people with standardised and accessible communication, especially in India with its diverse linguistic population. Sign language is an essential communication tool for thoughts, feelings, and experiences. It is used in both manual and visual modalities. Though sign languages differ throughout Indian regions, there is an urgent need for an Indian Sign Language standard that the deaf community and the general public can both understand. In the same way that most Indians find it easier to communicate effectively when they speak English, the

creation of an official Indian Sign Language might help reduce obstacles to communication and promote inclusivity across the country.

The primary objective of the Sign Language conversion project is to develop a comprehensive system that leverages cutting-edge technologies to facilitate real-time interpretation of sign language gestures. The project aims to integrate the Mediapipe library for landmark detection, enabling the extraction of vital hand landmark information from live webcam feeds. These landmark coordinates will serve as input data for a Random Forest Classifier, trained to recognize and classify various sign language patterns with high accuracy and efficiency. Additionally, the project seeks to incorporate natural language processing (NLP) techniques to provide English suggestions for detected signs and to translate them into Tamil, enhancing accessibility and linguistic diversity. By combining computer vision, machine learning, and NLP methodologies, the project endeavors to create a robust and versatile tool that empowers individuals with hearing and speech impairment to communicate effectively and seamlessly with the broader community.

The project's ability to offer English suggestions for detected signs, followed by Tamil and English voice translations, enhances accessibility and promotes linguistic diversity. Beyond these domains, the project has the potential for broader societal impact, fostering greater understanding and integration of individuals with hearing and speech impairments into mainstream social and professional spheres. As such, the scope of the project extends far beyond mere technological innovation, encompassing the promotion of inclusivity, accessibility, and effective communication across

various sectors of society. In this paper, the system utilizes the MediaPipe library, leveraging its landmark identification features to extract essential data from each frame of the live webcam stream. Specifically, the coordinates of hand landmarks are detected and extracted, providing crucial information about the hand gestures being performed. These landmark coordinates are then processed and used as input for the subsequent steps in the translation pipeline. Once the hand landmarks are detected, their coordinates are collated and stored in a structured format, such as a CSV file, for further analysis and training. This dataset serves as the foundation for training the Random Forest Classifier, a supervised learning algorithm that excels in classification tasks. Before training the model, the data may undergo preprocessing steps such as normalization and feature engineering to ensure optimal performance.

II. LITERATURE SURVEY

1. Gestures in American Sign Language (ASL) are characterized by fast, highly articulate motion of upper body, including arm movements with complex hand shapes and facial expressions. In this work, we propose a new method for word-level sign recognition from American Sign Language (ASL) using video. Our method uses both motion and hand shape cues while being robust to variations of execution. We exploit the knowledge of the body pose, estimated from an off-the-shelf pose estimator. Using the pose as a guide, we pool spatio-temporal feature maps from different layers of a 3D convolutional neural network. We train separate classifiers using pose-guided pooled features from different resolutions and fuse their prediction scores during test time. This leads to a significant improvement in performance on the WLASL benchmark dataset [25]. The proposed approach achieves 10%, 12%, 9.5% and 6.5% performance gain on WLASL100, WLASL300, WLASL1000, WLASL2000 subsets respectively. To demonstrate the robustness of the pose-guided pooling and proposed fusion mechanism, we also evaluate our method by finetuning the model on another dataset. This yields 10%

performance improvement for the proposed method using only 0.4% training data during the finetuning stage.

2. Sign language recognition (SLR) plays a crucial role in bridging the communication gap between the hearing and vocally impaired community and the rest of the society. Word-level sign language recognition (WSLR) is the first important step towards understanding and interpreting sign language. However, recognizing signs from videos is a challenging task as the meaning of a word depends on a combination of subtle body motions, hand configurations and other movements. Recent pose-based architectures of WSLR either model both the spatial and temporal dependencies among the poses in different frames simultaneously or only model the temporal information without fully utilizing the spatial information. We tackle the problem of WSLR using a novel pose-based approach, which captures spatial and temporal information separately and performs late fusion. Our proposed architecture explicitly captures the spatial interactions in the video using a Graph Convolutional Network (GCN). The temporal dependencies between the frames are captured using Bidirectional Encoder Representations from Transformers (BERT). Experimental results on WLASL, a standard word-level sign language recognition dataset show that our model significantly outperforms the state-of-the-art on pose-based methods by achieving an improvement in the prediction accuracy by up to 5%.

3. Understanding complex hand actions, such as assembly tasks or kitchen activities, from hand skeleton data is an important yet challenging task. In this paper, we analyze hand skeleton-based complex activities by modelling dynamic hand skeletons through a spatiotemporal graph convolutional neural network (ST-GCN). This model jointly learns and extracts spatio-temporal features for activity recognition. Our proposed technique, Symmetric Sub-graph spatio-temporal graph convolutional neural network (S2-ST-GCN), exploits the symmetric nature of hand graphs to decompose them into smaller sub-graphs, which allow us to build a separate temporal model for the rela-

tivemotionofthefingers.This

subgraph approach can be implemented efficiently by preprocessing input data using a Haar unit based orthogonal matrix. Then, in addition to spatial filters, separate temporal filters can be learned for each sub-graph. We evaluate the performance of the proposed method on the First-Person Hand Action dataset. While the proposed method shows comparable performance with the state of the art methods in train: test=1:1 setting, it achieves this with greater stability. Furthermore, we demonstrate significant performance improvement in comparison to state of the art methods in the cross-person setting, where the model did not come across a test subject's data while learning. S2-ST-GCN also shows superior performance than a finger-based decomposition of the hand graph where no preprocessing is applied.

III. EXPERIMENTAL DETAILS

A confusion matrix for American Sign Language (ASL) recognition represents the performance of a classification model by showing the predicted labels against the true labels for each class (ASL letter gesture). It provides detailed information about the Accuracy and misclassifications of the system. Moreover, Naive Bayes, which performs best with independent features, has the lowest accuracy rates. Yet, the traits are interdependent in the suggested strategy.

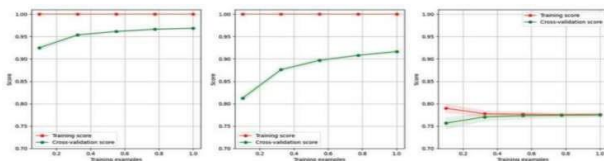


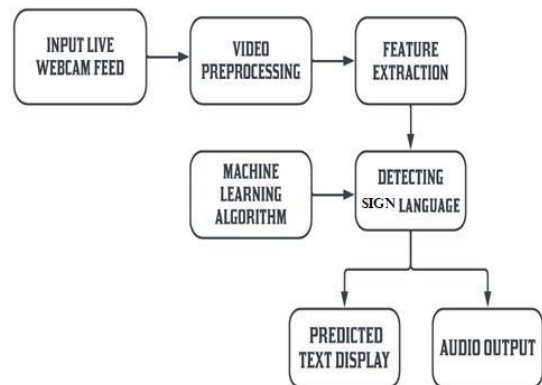
Fig- Accuracy score on the y-axis and training examples on the x-axis

The angle between landmark 12 and landmark 4, for instance, changes when the angle between landmark 0 and landmark

12 (12 is located on the tip of the middle finger), which influences the other characteristics as well. From the confusion matrix, we can calculate various performance metrics such as accuracy, precision, recall, and F1-score. It provides insights into which ASL gestures are more

challenging for the system to recognize and helps identify potential areas for improvement in the recognition model or training process.

SYSTEM ARCHITECTURE



DATA ACQUISITION:

The initial or the first step of this system is vision based i.e. to acquire data at runtime via the camera. Then, these data will be stored in a csv file format directory, in which all the images will be stored and trained by the user and the saved & trained data will be used to compare the recently captured image to the stored data's specific words.

FEATURE EXTRACTION:

The palm is extracted from the data's via image segmentation. This procedure revolves around converting raw data, such as images, into a meaningful set of features that can be effectively utilized for analysis and machine learning algorithms. In the context of sign language recognition, these extracted features hold vital information encompassing distinct patterns, and gestures that are indicative of various emotional states or behaviors. To begin, the dataset undergoes essential data preprocessing steps. This involves handling any missing data points, normalizing the data if necessary, and ensuring the overall cleanliness and preparedness of the dataset for subsequent phases. Upon loading the CSV file using relevant programming libraries, the data reveals itself as

rows, each representing a sample of sign language data, while columns correspond to specific attributes.

GESTURE RECOGNITION:

Gesture recognition within the realm of sign language recognition is a critical process that involves the identification and interpretation of various hand movements to deduce meaningful insights about a person's intentions, emotions, and communication cues. This sophisticated technology leverages advancements in computer vision and machine learning to translate physical gestures into actionable information. Through a analysis of posture, motion, and the spatial relationships of hand sign, gesture recognition systems can discern intricate details such as handshakes, nods, thumbs-up, and more complex gestures like pointing or even specific cultural gestures.

TEXT TO SPEECH:

Once the character is successfully recognized, the resulting output undergoes an additional transformation from text to speech. This conversion process is facilitated through the utilization of the English language process and GTTS library processing, a powerful text-to-speech conversion tool in Python. Unlike some other alternatives, this library operates offline, which ensures its compatibility and efficiency. This integration enables users to observe and simultaneously hear the translated sign language within our system, enhancing the overall convenience and usability of the application.

SUGGESTION OF WORDS:

The module description for English words suggestion and translation to Tamil incorporates sophisticated natural language processing (NLP) techniques to enhance the accessibility and comprehension of sign language interpretation. This module first utilizes NLP algorithms to suggest English words corresponding to detected sign language gestures, aiding both deaf users and non-signers in understanding the intended message. Leveraging advanced language models, the system generates contextually relevant word suggestions in real-time, ensuring accurate interpretation of sign language cues. Additionally, the module integrates

machinetranslationalgorithmstoseamlesslytranslatethesuggested English words into Tamil, catering to the linguisticdiversity of users. By providing instantaneous translations, thesystemfacilitateseffectivecommunicationandcomprehension across languagebarriers. Through thefusionof NLP technologies, this module enhances the usability andinclusivityofthesignlanguageconversionssystem,empoweringuserstoengageinseamlessandmeaningfulinteractions.

EXISTINGWORK:

BeforethedevelopmentoftheSignLanguageConversionproject,therewasno real-timesystemcapableofautomatically interpreting and classifying sign language cuesfrom a live webcam feed. Traditionally, understanding signlanguage requiredhumanobservationandanalysis,which could be subjective and time-consuming. Existing computervision systems focused on basic gesture detection but lackedcomprehensive sign language interpretation. Moreover, real-time analysis of sign language using landmark detection andmachine learning was not readily available. As a result, therewas a need for an innovative system that could efficientlydetect and analyzelandmarks from live video streams, andthen classify various sign language patterns in real-time. TheSign Language Conversion project addresses these limitationsandprovidesaneffectivesolutionfor non-verbalcommunication analysis, offering significant advancements inthefieldofhuman-computerinteraction.

PROPOSEDSYSTEM:

TheproposedsystemfortheSignLanguageconversionproject integrates advanced technologies to enable real-timeinterpretationofsignlanguagegestures.UtilizingtheMedia pipe library, the system detects hand landmarks fromlivewebcamfeeds,providingcrucialdataforanalysis.A

Random Forest Classifier is then employed to recognize and classify these landmarks into various sign language patterns, ensuring accurate interpretation. Moreover, natural language processing techniques are utilized to generate English suggestions for detected signs and translate them into Tamil, enhancing accessibility and comprehension. The system overlays the interpreted sign language cues, including detected signs and translations, onto the live video stream in real-time, facilitating immediate understanding for users. Additionally, a user-friendly interface allows for seamless interaction and feedback, ensuring ongoing refinement and improvement. Through rigorous testing and validation procedures, the system aims to deliver a reliable and effective solution for bridging communication barriers and promoting inclusivity for individuals with hearing and speech impairments.

TRAINED MODULE

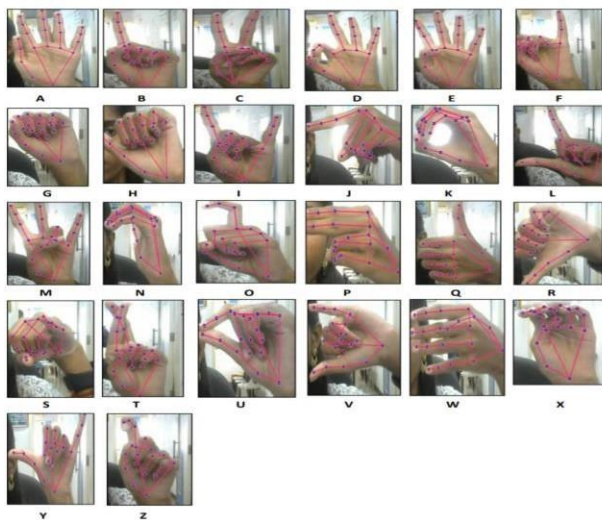


Fig-modules

III.RESULT AND DISCUSSION

The comprehensive examination of the simulation results using the "Machine Learning" is provided in this part.

Additionally, employing the dataset, the suggested method's performance is contrasted with that of current approaches.

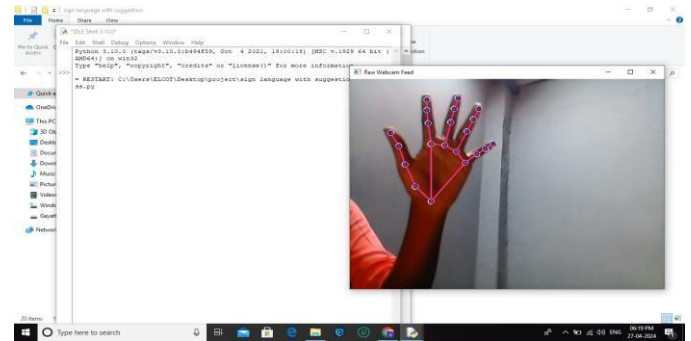


Fig.Hand Gesture Recognition

This method of hand identification uses the media pipe library, which is used for image processing, to first identify a hand from an image captured by a webcam. Thus, once the hand is located in the image, we obtain the region of interest (Roi), crop the image, and use the OpenCV library to transform it to a grayscale image after applying gaussian blur. The OpenCV library, commonly referred to as the Open Computer Vision Library, makes it simple to apply the filter. Next, we used the threshold and adaptive threshold methods to convert the grayscale image to a binary image. For the sign letters A through Z, we have a set of pictures of various signs at various perspectives.

Mediapipe Landmark System:



Fig.Landmark System

Now that we have these landmark points, we use the OpenCV library to draw them on a simple white background. We address the background and lighting situations by doing this since the mediapipe library will provide us with landmark locations in any backdrop and primarily under all lighting circumstances.



Fig.LandmarkProcessingssystemusingOpenCV**V.CONCLUSION**

In conclusion, the Sign Language Conversion project represents a significant milestone in the development of automated, real-time systems for interpreting and classifying sign language cues. By leveraging the power of computer vision and machine learning, the project achieves remarkable accuracy and efficiency in detecting hand landmarks and interpreting non-verbal communication from live webcam feeds. The integration of a Random Forest Classifier ensures reliable and objective sign language classification, contributing to the system's consistency and effectiveness. Moreover, the user-friendly frontend enhances the interactive experience, providing users with real-time analysis results and immediate feedback. With diverse applications in human-computer interaction and user behavior analysis, the project marks a substantial advancement in the field of non-verbal communication analysis. Its success offers valuable insights and opportunities for further research and development in this burgeoning domain, promising to revolutionize accessibility and inclusivity for individuals with hearing and speech impairments.

REFERENCES:

- [1] Y. Saleh and G. F. Issa, "Arabic sign language recognition through deep neural networks fine-tuning," *Int.J.Online Biomed.Eng.*, vol.16, no. 5, pp.71–83, 2020.
- [2] X. Jiang, M. Lu, and S.-H. Wang, "A eight-layer convolutional neural network with stochastic pooling, batch normalization and dropout for fingerspelling recognition of Chinese sign language," *Multimedia Tools Appl.*, vol.79, nos.21–22, pp.15697–15715, Jun. 2020.
- [3] O. Sevli and N. Kemaloglu, "Turkish sign language digits classification with CNN using different optimizers," *Int. Adv. Researches Eng. J.*, vol.4, no.3, pp.200–207, Dec.2020
- [3] A. Tyagi and S. Bansal, "Feature extraction technique for vision-based Indian sign language recognition system: A review," in *Computational Methods and Data Engineering*. Singapore: Springer, 2021, pp.39–53.
- [4] M. Mukushev, "Evaluation of manual and non-manual components for sign language recognition," in *Proc.12th Lang. Resour. Eval. Conf., Eur. Lang. Resour. Assoc. (ELRA)*, 2020, pp. 1–6.
- [5] A. Tunga, S. V. Nuthalapati, and J. Wachs, "Pose-based sign language recognition using GCN and BERT," in *Proc. IEEE Winter Conf. Appl. Comput. Vis. Workshops (WACV W)*, Jan. 2021, pp.31–40.
- [6] N. C. Camgoz, S. Hadfield, O. Koller, H. Ney, and R. Bowden, "Neural sign language translation," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Jun. 2018, pp.7784–7793.
- [7] S. Aly and W. Aly, "Deep ArSLR: A novel signer-independent deep learning framework for isolated Arabic sign language gestures recognition," *IEEE Access*, vol.8, pp.8319–83212, 2020.
- [8] D. Li, X. Yu, C. Xu, L. Petersson, and H. Li, "Transferring cross-domain knowledge for video sign language recognition," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 202

[9] R.Rastgoo,K.Kiani,andS.Escalera,“Video-basedisolated hand sign language recognition using a deep cascadedmodel,” MultimediaTools Appl.,vol.79, nos.31–32, pp.22965–22987, Aug. 2020.