

Bridging the Communication Gap - Sign Language Communication System

Varsha Malla¹, Avadhesh Vora², Caleb D'mello³, Shubham Kardel⁴, Prof. Vardha Gotmare⁵

Students, Department of Computer Science^{1,2,3,4}

Assistant Professor, Department of Computer Science⁵

Dr. D. Y. Patil Institute of Technology Pimpri, Pune, India

varshamalla020@gmail.com, avadheshvora2002@gmail.com, calebdmello2@gmail.com

contact.shubham.kardel@gmail.com, vardha.gotmare@gmail.com

Abstract: Now, imagine a world where words are not spoken, and emotions are not expressed. For the Deaf community, this is not a mere thought experiment but a daily reality. Indeed, they are forced to communicate primarily through sign language, therefore, regularly feel alienated in a society where words are the primary tool for translating thoughts and feelings. Under such circumstances, most Deaf and hard-of-hearing people are struggling to communicate the most effectively. They usually opt for interpreters or their own signaling, but these alternatives are not the most efficient or effective choices. While signaling is the most intuitive and meaningful alternative, its grammar and semantic variations make it impossible to comprehend for those people who are not within the culture. Therefore, we have also developed a software prototype that could interpret sign language automatically. This paper presents a new method to recognize Indian sign language alphabets A-Z and digits (0-9) in a real-time video feed by employing the Bag of Visual Words model (BOVW). Furthermore, it not only predicts the labels of the signs but also gives the output as text and speech. The segmentation process here comprehends skin color detection and background subtraction. Speeded Up Robust Features (SURF) features have been extracted again from the images, and histograms are built to map the signs to text labels. Furthermore, we have utilized algorithms such as Support Vector Machine (SVM) and Convolutional Neural Networks (CNN) for binary classification. Finally, we have established an interactive Graphical User Interface (GUI) to make this process more user-friendly.

Keywords: Indian Sign Language (ISL), Bag of Visual Words (BOVW), Convolutional Neural Networks (CNN), Support Vector Machine (SVM), Speeded Up Robust Features (SURF), Speech Recognition, Google Speech

I. INTRODUCTION

There are different types of gestures: static, dynamic, or a combination of the two that serve as modes of non-verbal communication where movements of the body convey information. Communication is an essential part of people's lives that enables ideas to be shared and emotions expressed, thus creating a bond between people through mutual understanding. Regardless of traditional verbal communication is out of reach to certain groups of individuals with disabilities, especially those who are deaf and muted. They rely on sign language, which becomes last mode for them to communicate with the world at large. Despite this fact, us normal human beings don't really struggle much when we connect with each other. We can easily express ourselves through talking, gestures, body language, reading, and writing, with speech being the most common way we do it. But for folks dealing with speech impairment, they rely only on sign language, making it a lot harder for them to communicate with the rest of us. This means we need sign language recognizers that can understand and change sign language into spoken or written language, and vice versa. But these recognizers are limited, expensive, and not very easy to use. Now, researchers from different parts of the world are getting into this sign language recognizer program, and so we are able to see that the world is taking steps forward in this development too.

Even though India is a super diverse country with almost 17.7 percent of the world's population living here, there hasn't been a whole lot of work done in this research area, which is kind of surprising compared to other countries [3], [15], [16]]. This delay in standardization could be the reason for this. Indian Sign Language studies started in India back in 1978, but since there was no standard type of ISL, it was only used in short-term courses. Plus, the gestures used in most deaf schools were really different from each other, and only about 5 percent of all deaf people went to these schools. It wasn't until 2003 that ISL got standardized and caught the attention of researchers [18].

Indian Sign Language (ISL) includes both still and moving signs, single and double-handed signs, and there are lots of different signs for the same letter in different parts of India. This makes it super hard to come up with a system for it. Plus, there's no standard set of data available. All of this shows just how complex Indian sign language really is.

Lately, researchers have started looking into this area. There are mainly two different approaches widely used in sign language recognition: the Sensor-based approach and the Vision-based approach [11]. The sensor-based approach uses gloves or other tools to recognize finger gestures and turn them into electrical signals for sign determination, while web cameras are used to capture video or images in a vision-based approach. Because it doesn't need any special equipment, the vision-based gesture recognition is more natural and preferred by signers [1]. But hand segmentation in a complex setting is a big deal in identification. That's why we need a framework that can handle this problem.

The progress in machine learning and deep learning technology is coming up with new methods and algorithms for recognizing Indian sign language letters in a more efficient, accurate, and affordable way. These very precise and end-to-end models are getting rid of the prior limitations of traditional methods, making the results more accurate and efficient.

In this work, we're laying out a way to build a big, varied, and strong real-time alphabet (A-Z) and number (0-9) recognition system for Indian Sign Language. Instead of using fancy tech like gloves or the Kinect, we're recognizing signs from images taken from a webcam. We're also talking about the accuracy we got in the results. We need a real-time, accurate, and efficient way to recognize ISL signs to bridge the gap between people who can hear and speak and those who can't.

II. LITERATURE SURVEY

Depending on the nature of sign language and the signs, different authors have used different methods.

J. Singha et al. [17] came up with a method for recognizing signs in real time by using Eigen value-weighted Euclidean distance for classification. P. Kishore et al. [9] proposed a system that finds active contours from boundary edge maps using Artificial Neural Network (ANN) to classify signs. Another approach used the Viola Jones algorithm with LBP functions for hand gesture recognition in real-time. It had the advantage of needing less processing power to detect movements. Segmentation is a crucial step in hand processing, and in general, Otsu's algorithm showed a fairly high rate of accuracy. In a separate attempt, the moving block distance parameterization method was used to skip the initialization and segmentation steps. High precision static symbols and 33 basic word units were used.

Most of these works were based on pattern recognition, feature extraction, and so on [13]. However, in most cases, a system with a single feature isn't enough. That's why hybrid approaches were introduced to tackle this issue. For example,

A. Nandy et al. [12] used hybrid approaches with K-Nearest Neighbor (KNN) and Euclidean distance to classify gestures from orientated histogram features. The drawback of this approach was its poor performance with similar gestures. K. Manjushree et al. [6] used single-handed sign classification with histogram of oriented gradients and feature matching. S. Kanade et al. [5] designed a system with a custom dataset using PCA features and SVM, and achieved good accuracy.

A. Sahoo [14] proposed ISL recognition for both single and double-handed character signs. Geetha. M et al used B-Spline approximation for the shape matching of static gestures of ISL alphabets and numerals. In Ref. [4], a method was proposed to classify word symbols using the Neuro-Fuzzy approach and natural language processing (NLP) technology to display the final word. The combination of PCA and the local coordinate system produced high calculation accuracy and was found to be superior to the method based on the condensation algorithm.

However, for real-time systems, researchers needed a faster way to solve this problem. The advancements in Deep Learning technologies have enabled automation of image recognition using various image recognition models. For

example, Convolutional neural networks have made great strides in the field of deep learning in recent years [2]. G. Jayadeep et al. [7] used a CNN (Convolutional Neural Network) to extract image features, LSTM (Long Short Term Memory) to classify these gestures and translate them into text. Binet et al. [10] proposed the InceptionV3 model to use depth sensors to identify static signs. It eliminated the steps of gesture segmentation and feature extraction. In Ref. [8], Vivek Bheda et al. proposed a methodology for using a mini-batch supervised learning method of stochastic gradient descent to classify images for each digit (0–9) and American Sign Language letter using deep convolutional neural networks.

After reviewing these works, the authors were motivated to create a custom dataset and an algorithm that would work exclusively on that dataset without compromising the accuracy of video detection. They decided to use SURF features because it would reduce the measurement time and make the system invariant to rotation. The authors of the paper also addressed the issue of background dependency so that the system can be used anywhere, not just in controlled environments.

III. PROPOSED WORK

Sign language recognition needs reliable and efficient data to create a highly accurate system that would be beneficial for real-time users. The data flow at various stages for sign language recognition, including Dataset, Image Acquisition, Data Pre-processing, Feature Extraction, and Sign Classification, is illustrated in Fig. 1.

- Dataset Creation
- Data Pre-processing
- Feature Extraction
- Sign Classification

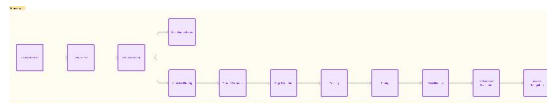


Fig. 1. Flow diagram of Proposed work.

Dataset Creation

In the pursuit of this project, the task of obtaining a comprehensive dataset for hand signs presented a significant challenge, especially in the case of alphabets J and Z, which involve dynamic motion. Our initial exploration led us to a dataset available on Kaggle, which initially showed promise but unfortunately lacked images specifically for these letters. The dataset's format closely resembled the classic MNIST, where each training and test case represented a label (0- 25) corresponding to alphabetic letters A-Z. However, it was notable that the dataset did not include cases for letters J and Z due to their gestural nature. Despite this limitation, the dataset still provided valuable training examples, comprising 27,455 training cases and 7,172 test cases, with each case consisting of a 28x28 pixel image with grayscale values ranging from 0 to 255.

To address the absence of data for letters J and Z, we made the strategic decision to create our own dataset. This involved capturing videos of the signs, ensuring the acquisition of multiple images for each letter to effectively capture the nuanced motion involved. Great care was taken to maintain a clear and consistent background, aligning with the standards set by the existing dataset.

Moreover, in addition to our custom dataset creation, we also sourced a separate Kaggle dataset specifically focused on numerical signs, covering digits 0-9. This comprehensive approach not only addressed the limitations of the existing dataset but also ensured the development of a diverse and all-encompassing training set for our model. This combination of leveraging existing datasets and creating our own lays a strong foundation for the accurate and robust recognition of sign language, aligning with our commitment to excellence in this vital area of research.

The dataset images can be seen below fig2 and fig3

Data Preprocessing

In the preparation of the dataset for the sign language recognition project, the initial challenge lay in cleaning the image data to ensure its suitability for training the model. Various tools such as Python's OpenCV library were used to address issues such as noise and unwanted artifacts in the images. However, one significant hurdle encountered was inconsistent lighting conditions across different image samples, leading to variations in brightness and contrast. To mitigate this issue, histogram equalization techniques were implemented to normalize the brightness levels across all images. This preprocessing step not only improved the quality of the dataset but also enhanced the robustness of the model to varying lighting conditions.



Fig 2. Example of images of alphabets



Fig. 3. Example of images of numbers

Furthermore, resizing the images to a uniform dimension posed another critical aspect of data preprocessing. Despite efforts to maintain consistency in image sizes, challenges were encountered in preserving the aspect ratio while resizing. This resulted in distorted images that could adversely affect the performance of the recognition system. To address this problem, a two-step approach was adopted, first resizing the images to a common width while preserving the aspect ratio and then cropping or padding them to achieve the desired height. This strategy ensured that the hand gestures remained accurately represented in the resized images, thereby maintaining the integrity of the dataset.

Additionally, data augmentation emerged as a crucial step in augmenting the size and diversity of the dataset. However, augmenting images posed its own set of challenges, particularly in maintaining the semantic meaning of the hand gestures. Issues were encountered where augmented images deviated significantly from the original gestures, leading to confusion during training. To overcome this challenge, augmentation techniques such as rotation, translation, and flipping were implemented, while ensuring that the semantic integrity of the gestures was preserved. By carefully selecting augmentation parameters and validating the augmented data, the dataset was successfully expanded without compromising the quality or semantic coherence of the hand gestures.

Feature Extraction

In this phase, our focus is on the construction of a Bag of Visual Words (BOVW), which entails a series of crucial steps including feature extraction, clustering of features, codebook construction for the model, and the generation of histograms. The Bag of Visual Words (BOVW) model is widely utilized in image classification and is adapted from the concept of Bag of Words (BOW) in Natural Language Processing (NLP) and data retrieval. Similar to BOW's approach of counting the frequency of each word in a text to derive keywords and produce a frequency histogram, the BOVW model uses image features as "words" to construct a vocabulary. This involves representing each

image as a frequency histogram of characteristics obtained from its image descriptors and key points, allowing for the prediction of the category of comparable images based on this frequency histogram.

To begin building a bag of visual words (BOVW), the first step is to extract descriptors from each image in the dataset. These descriptors, which consist of 64-member vectors for each interest point, define the distribution of intensity material within the neighborhood of the interest point. For this purpose, we utilize SURF (Speeded Up Robust Features), a local feature detector and descriptor known for its robustness against rotation, variance, point of view occlusion, and its use of box filters for efficient computation.

Each image is represented as a set of image descriptors provided by SURF, as shown in Eq (1):

$$I_m = \{d_1, d_2, d_3, \dots, d_n\} \quad (1)$$

where d_i represents the color, shape, etc. of the hands, and n denotes the total image descriptors. Figure 4 illustrates the extracted SURF features when a binary image representing sign A is passed to the SURF algorithm. This approach allows for the effective extraction and representation of image features crucial for the subsequent steps in the construction of the Bag of Visual Words model

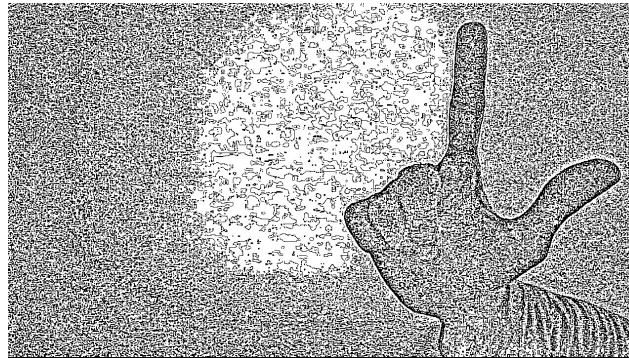


Fig. 4. Surf feature Extraction

In the extraction of features, the next step involves clustering all the features obtained after applying SURF. This process aims to group similar features, allowing us to use the core and cluster them as the dictionary's visual keyword. While the clustering can be performed using the K-means algorithm, we opted for mini batch K-Means due to the large volume of data. This approach, similar to K-means, offers improved processing time and memory utilization by employing small random batches of fixed-size data at a time, thereby reducing the need to hold all the data in memory simultaneously. With a value of k set at 180, this method involves obtaining a new random sample from the dataset in each iteration to update the clusters, which is repeated until convergence.

$$v = \{w_1, w_2, w_3, \dots, w_k\} \quad (2)$$

Following the clustering process, the resulting cluster centers (centroids) are treated as our code vectors for codebook generation. The codebook is instrumental in quantizing features, mapping a feature vector to the index of the nearest code vector. This constructed vocabulary, represented as k (the total number of clusters, i.e., 180), facilitates the mapping of each descriptor to the nearest visual word, as per Eq (3), where $w(d_i)$ represents the visual word assigned to the i th descriptor, and $\text{Dist}(w, d_i)$ denotes the distance between the visual word w and descriptor d_i .

$$w(d) = \text{argmin} \text{Dist}(w, d) \quad (3)$$

The final step involves generating histograms for all the images by calculating the frequency of occurrence of each visual word in an image. The count of bins in the histogram, equal to the total number of visual words in the dictionary (k), is represented by Eq (4). Here, D_i signifies the set of all the descriptors corresponding to a particular visual word w_i in the image, and $C(D_i)$ represents the cardinality, indicating the count of elements in set D_i . This process is repeated for every visual word in the image to obtain final histograms, which are then passed for recognition to the classifier along with their respective labels.

$$\text{bin}_i = C(D_i) \quad (4)$$

where

$$D_i = \{d_j, j \in 1, \dots, n \mid w(d_j) = w_i\} \quad (5)$$

Classification

1) *CNN*: Convolutional Neural Networks (CNNs) is indeed a powerful models, it is made as a replication from the human brain’s visual cortex. It excel in comparing images piece by piece using filter maps that slide over local image patches, enabling the identification and comparison of similar features at corresponding locations in different images. This unique quality provides CNNs an upper hand in taking and naming images compared to other neural networks. Our CNN architecture adopts a standard structure, incor- porating multiple convolutional and dense layers, with each CNN being 3 layers deep. The architecture begins with a pair of 2 convolutional layers containing 32 filters with a window size of 3×3 , followed by a max-pool layer and a dropout layer. This is followed by another set of 2 convolutional layers with 64 filters, along with a max pooling layer and dropout layer. Additionally, there are 2 more convolutional layers with 64 filters and a max pooling layer, leading to a fully connected hidden layer with 512 neurons of the ReLU activation function and an output layer of the softmax activation function. Notably, the first convolution layer handles an input image of size (100,100), while the final output layer comprises 36 neurons, each corresponding to a category of the ISL signs. The architectural diagram for this model is depicted in Fig. 5, providing a visual representation of its structure and components.

In summary, the CNN architecture offers a robust and versatile model for image classification, leveraging its multi-layered design to effectively process and classify visual data. The use of convolutional and dense layers, in conjunction with specific filter sizes and pooling layers, contributes to the network’s ability to extract and process features, ultimately leading to accurate and efficient classification of ISL signs.



Fig. 5. CNN

Output Sign

Once the classifier identifies the predicted class labels, they are automatically translated by the system into text and speech formats, enhancing communication and user convenience. The identified label is then utilized as a key in a dictionary to retrieve the corresponding sign as its value, which is then presented to the user. To make the text-to-speech conversion, the Python text-to-speech module, Pyttsx3, is applied. To avoid the delays in the live video stream caused by the slow pro- cessing of frames, threading was implemented. This approach enables simultaneous prediction of signs and translation of text to speech, ensuring continuous and uninterrupted playback of the audio

IV. EXPERIMENT AND RESULT

The dataset has been divided into two sets to maximize effectiveness as per the all the standred Techniques the data is divided into: 80 percent for training and the rest 20 percent for testing. It’s was really fascinating to see that both the SVM and CNN classifiers have achieved impressive accuracy while processing the images. Moreover the CNN has truly gave us the amazing result which were really impressive in the performance with a really very feable amount of features The system is trained to understand and recongnise total of 36 signs, which includes the 26 alphabets and 10 numerals. The current results are undeniably promising, and it’s exciting to anticipate even greater achievements with a few strategic refinements.

CNN Performance

By using CNN, the model which was tested then we achieved an accuracy of 94 percent on the training set and it gave us more fruitful results over the testing set to get more remarkable testing accuracy of over 99 percent in the final epoch. We achieved this after near about 50 epochs of training. For this we used a categorical cross entropy loss function and the softmax function as the activation function, which resulted in a training loss of 0.1748 and a testing loss of 0.0184 in the last epoch. The class-wise accuracy is provided for further analysis of the model’s performance. Additionally, the accuracy graph for our experiment is shown in Fig. 6, giving a clear visual representation of the model’s accuracy trends. And table. 1 shows the percentage of accuracy which was achieved in indentifying the Labels.

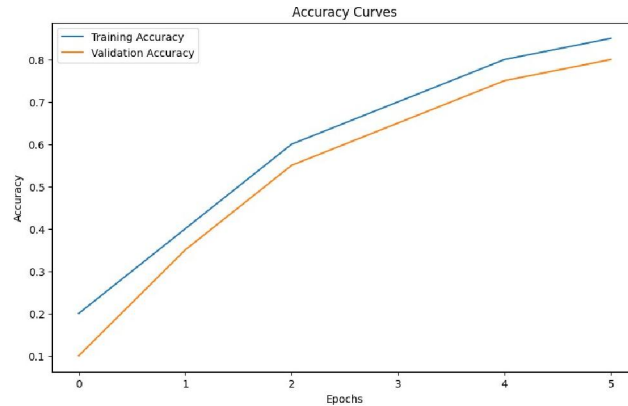


Fig. 6. Accuracy Graph

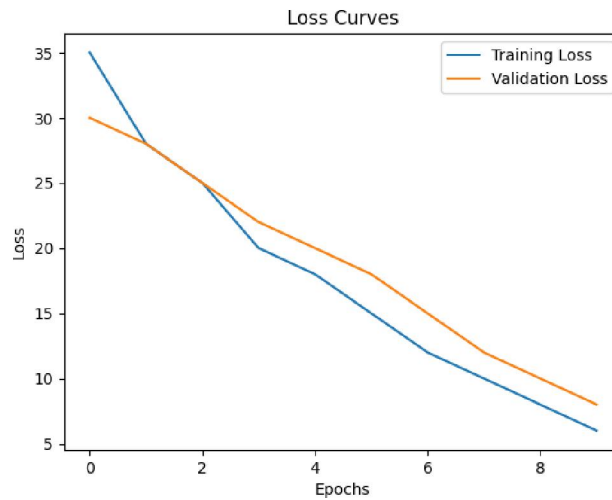


Fig. 7. Loss function graph

Real time testing

An interactive GUI has been meticulously crafted for our system, catering to users with a fully operational sign-in and sign-up system using Tkinter. With this interface, users have the ability to predict signs based on the model trained using our dataset simply by clicking on the predict sign button.

TABLE I: CLASS-WISE ACCURACY TABLE

Label	CNN (%)	Label	CNN (%)	Label	CNN (%)
0	100	C	100	O	99
1	100	D	100	P	100
2	100	E	97	Q	100
3	98	F	100	R	98
4	100	G	100	S	100
5	99	H	100	T	100
6	100	I	100	U	100
7	100	J	100	V	100
8	100	K	100	W	100
9	100	L	100	X	99
A	100	M	99	Y	100
B	100	N	100	Z	100

Additionally, they can create their database using the create signs button, offering a seamless and personalized experience. Furthermore, an option for speech-to-sign conversion has been thoughtfully incorporated to enhance user interaction. For a glimpse into the real-time video testing, screenshots are available for reference, providing valuable insights into the system's functionality and user experience.

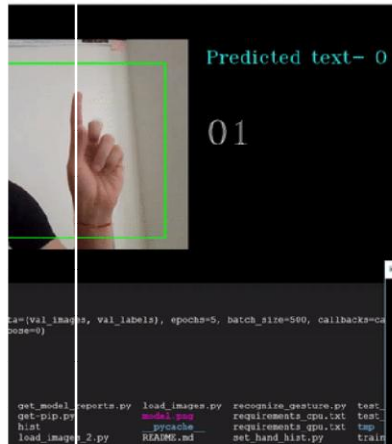


Fig. 8. GUI image

V. CONCLUSION

In this paper, a new way is shown for understanding and naming Indian hand signs. It covers letters (A-Z) and numbers (0-9), using Convolutional Neural Network (CNN). The main aim is to make the system quicker and usable in many situations. This is done by making a special set of data, to make sure the system is strong when things move and to stop problems with the background. Wow! The system is 99 percent right in training, with 36 Indian hand signs! In the future, the plan is to add more signs from other languages. This will make the system better and help with real-time use. Also, the plan is to include whole words and signs, for both fast and separate tasks. To do this, we need to make the system faster—a hard job that we must keep working on.

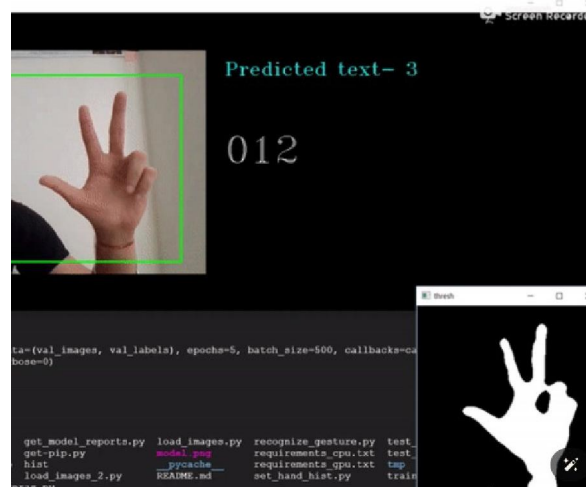


Fig. 9. GUI Image

REFERENCES

- [1]. PK Athira, CJ Sruthi, and A Lijiya. A signer independent sign language recognition with co-articulation elimination from live videos: an indian scenario. Journal of King Saud University-Computer and Information Sciences, 34(3):771–781, 2022.

- [2]. Shailesh Bachani, Shubham Dixit, Rohin Chadha, and A Bagul. Sign language recognition using neural network. *International Research Journal of Engineering and Technology (IRJE)*, 2020.
- [3]. Kshitij Bantupalli and Ying Xie. American sign language recognition using machine learning and computer vision. 2019.
- [4]. Hemina Bhavsar and Jeegar Trivedi. Indian sign language recognition using framework of skin color detection, viola-jones algorithm, correlation-coefficient technique and distance based neuro-fuzzy classification approach. In *Emerging Technology Trends in Electronics, Communication and Networking: Third International Conference, ET2ECN 2020, Surat, India, February 7–8, 2020, Revised Selected Papers 3*, pages 235–243. Springer, 2020.
- [5]. Padmanabh D Deshpande and Sudhir S Kanade. Recognition of indian sign language using svm classifier. *International Journal of Trend in Scientific Research and Development (IJTSRD)*, 2(3):1053–1058, 2018.
- [6]. Manjushree K Divyashree. Gesture recognition for indian sign language using hog and svm. *International Research Journal of Engineering and Technology*, 6(7), 2019.
- [7]. Gautham Jayadeep, NV Vishnupriya, Vyshnavi Venugopal, S Vishnu, and M Geetha. Mudra: convolutional neural network based indian sign language translator for banks. In *2020 4th International Conference on Intelligent Computing and Control Systems (ICICCS)*, pages 1228–1232. IEEE, 2020.
- [8]. Shagun Katoch, Varsha Singh, and Uma Shanker Tiwary. Indian sign language recognition system using surf with svm and cnn. *Array*, 14:100141, 2022.
- [9]. PVV Kishore, MVD Prasad, D Anil Kumar, and ASCS Sastry. Optical flow hand tracking and active contour hand shape features for continuous sign language recognition with artificial neural networks. In *2016 IEEE 6th international conference on advanced computing (IACC)*, pages 346–351. IEEE, 2016.
- [10]. Gongfa Li, Heng Tang, Ying Sun, Jianyi Kong, Guozhang Jiang, Du Jiang, Bo Tao, Shuang Xu, and Honghai Liu. Hand gesture recognition based on convolution neural network. *Cluster Computing*, 22:2719–2729, 2019.
- [11]. Anuja V Nair and V Bindu. A review on indian sign language recognition. *International journal of computer applications*, 73(22), 2013.
- [12]. Anup Nandy, Jay Shankar Prasad, Soumik Mondal, Pavan Chakraborty, and Gora Chand Nandi. Recognition of isolated indian sign language gesture in real time. In *Information Processing and Management: International Conference on Recent Trends in Business Administration and Information Processing, BAIP 2010, Trivandrum, Kerala, India, March 26-27, 2010. Proceedings*, pages 102–107. Springer, 2010.
- [13]. Yogeshwar I Rokade and Prashant M Jadav. Indian sign language recognition system. *International Journal of engineering and Technology*, 9(3):189–196, 2017.
- [14]. Ashok Kumar Sahoo and Kiran Kumar Ravulakollu. Vision based indian sign language character recognition. *Journal of Theoretical & Applied Information Technology*, 67(3), 2014.
- [15]. Shadman Shahriar, Ashraf Siddiquee, Tanveerul Islam, Abesh Ghosh, Rajat Chakraborty, Asir Intisar Khan, Celia Shahnaz, and Shaikh Anowarul Fattah. Real-time american sign language recognition using skin segmentation and image category classification with convolutional neural network and deep learning. In *TENCON 2018-2018 IEEE Region 10 Conference*, pages 1168–1171. IEEE, 2018.
- [16]. S Shivashankara and S Srinath. A comparative study of various techniques and outcomes of recognizing american sign language: a review. *International Journal of Scientific Research Engineering & Technology (IJSRET)*, 6(9):1013–1023, 2017.
- [17]. Joyeeta Singha and Karen Das. Recognition of indian sign language in live video. arXiv preprint arXiv:1306.1301, 2013.
- [18]. Daleesha M Viswanathan and Sumam Mary Idicula. Recent developments in indian sign language recognition: an analysis. *International Journal of Computer Science and Information Technologies*, 6(1):289–293, 2015