

SIGN LANGUAGE RECOGNITION USING CONVOLUTIONAL NEURAL NETWORKS

¹*Dr.P.Veeresh,*
²*Gowlla Ramireddy,*
³*Kunigiri Hemanth,*
⁴*K.Sathyanarayana,*
⁵*Golla Manohar*

Journal for Educators, Teachers and Trainers, Vol.14 (2)

<https://jett.labosfor.com/>

Date of reception: 18 Jan 2023

Date of revision: 09 Feb 2023

Date of acceptance: 21 Mar 2023

Dr.P.Veeresh, Gowlla Ramireddy, Kunigiri Hemanth, K.Sathyanarayana, Golla Manohar (2023). SIGN LANGUAGE RECOGNITION USING CONVOLUTIONAL NEURAL NETWORKS. Journal for Educators, Teachers and Trainers, Vol. 14(2). 650-657.

SIGN LANGUAGE RECOGNITION USING CONVOLUTIONAL NEURAL NETWORKS

¹Dr.P.Veeresh,²Gowlla Ramireddy,³Kunigiri Hemanth,⁴K.Sathyanarayana,⁵Golla Manohar

¹Professor & HOD,²Students

Department Of CSE

St. Johns College of Engineering & Technology, Errakota, Yemmiganur

ABSTRACT

Sign Language Recognition (SLR) aims to translate sign language into written or spoken language, in order to enhance communication between those who are deaf or mute and those who are not. This work has significant societal implications, but remains very demanding owing to the intricate nature and extensive range of hand movements. Current approaches for SLR use manually designed characteristics to characterize sign language movement and construct classification models using these features. Designing dependable features that can accommodate the wide range of hand movements is a challenge. In order to address this issue, we suggest the use of a groundbreaking convolutional neural network (CNN) that can automatically extract distinctive spatial-temporal characteristics from unprocessed video streams, without the need for any previous information or the creation of specific features. In order to enhance performance, the CNN utilizes several video streams, including color information, depth clues, and body joint locations, as input. This allows for the integration of color, depth, and trajectory information. We assess the suggested model using an actual dataset obtained from Microsoft Kinect and showcase its superior performance compared to conventional methods that rely on manually designed features.

1. INTRODUCTION:

Sign language, a prevalent mode of communication for those with hearing impairments, is conveyed by a combination of hand forms, body gestures, and facial expressions. Sign language identification remains a very tough problem due to the complexity of effectively using information from hand-shapes and body movement trajectory in a collaborative manner. This study presents a proficient recognition model for converting sign language into text or voice, aiming to facilitate communication between those with hearing impairments and those without by means of sign language.

The primary difficulty in sign language identification is in the development of descriptors that accurately represent hand forms and motion trajectories. Specifically, hand-shape description entails the process of identifying hand regions within a video stream, extracting hand-shape pictures from a complex backdrop in each frame, and addressing the challenges associated with recognizing motions. The motion trajectory is also associated with the tracking of critical points and the matching of curves. Despite extensive study on these two challenges, obtaining satisfactory results for SLR remains challenging because to the variability and obstruction of hand and body joints. In addition, the task of combining hand-shape characteristics with trajectory information is a complex one. In order to overcome these challenges, we have developed Convolutional Neural Networks (CNNs) that effectively combine hand forms, movement trajectories, and face expressions. Instead of using typically used color photos as input for networks such as [1, 2], we employ color images, depth images, and body skeleton images concurrently as input, all of which are acquired by Microsoft Kinect.

Kinect is a motion sensor that is capable of providing both color stream and depth feed. The public Windows SDK allows for real-time retrieval of body joint positions, as seen in Figure 1. Thus, we have selected Kinect as the capture device for recording the dataset of sign words. The variation in color and pixel depth provides valuable data for distinguishing between various sign actions. The temporal fluctuation of bodily joints may represent the trajectory of sign activities. CNNs are able to detect changes in color, depth, and trajectory by using various visual sources as input. It is important to note that we may circumvent the challenges of monitoring hands, separating hands from the backdrop, and creating descriptors for hands by using Convolutional Neural Networks (CNNs). CNNs possess the capacity to autonomously learn features from raw data without any previous information [3].

Convolutional Neural Networks (CNNs) have been used for video stream categorization in recent years. An issue that may arise with Convolutional Neural Networks (CNNs) is their time-consuming nature. Training a Convolutional Neural Network (CNN) using millions of films at a million-scale often takes many weeks or months. Thankfully, it is still feasible to get real-time performance by using CUDA for parallel computing. Our proposal involves using Convolutional Neural Networks (CNNs) to extract both spatial and temporal characteristics from a video stream in order to recognize Sign Language (SLR). Current approaches to SLR include manually designed characteristics to characterize the motion of sign language and construct a classification model using these features. On the other hand, Convolutional Neural Networks (CNNs) have the ability to automatically extract motion information from unprocessed video data, eliminating the need for manual feature creation. We use Convolutional Neural Networks (CNNs) that utilize several forms of data as input. This architectural design incorporates color, depth, and trajectory data by using convolution and subsampling techniques to consecutive video frames. The experimental findings indicate that 3D CNNs exhibit superior performance compared to Gaussian mixture model with Hidden Markov model (GMM-HMM) baselines in recognizing certain sign phrases that were recorded by us.

2. SYSTEM ANALYSIS

Existing System

Creating a desktop application that uses a computer's webcam to capture a person signing gestures for American Sign Language (ASL), and translate it into corresponding text and speech in real time. The translated sign language gesture will be acquired in text which is farther converted into audio.

Disadvantage of existing system

1. Less efficiency.

Proposed system

In this manner we are implementing a finger spelling sign language translator. To enable the detection of gestures, we are making use of a Convolutional neural network (CNN). A CNN is highly efficient in tackling computer vision problems and is capable of detecting the desired features with a high degree of accuracy upon sufficient training.

Advantage of Proposed system

1. More efficiency.

3. MODULES DESCRIPTION:

User:

The User can start the project by running run.py file. User has Upload Hand Gesture Dataset, Train CNN with Gesture Images User has to open cv class VideoCapture(0) means primary camera of the system, VideoCapture(1) means secondary camera of the system. VideoCapture(Videofile path) means with out camera we can load the video file from the disk. Vgg16, Vgg19 has programitaically configured. User can change the model selection in the code and can run in multiple ways.

HSR System:

Video-based Hand Sign recognition can be categorized as vision-based according. The vision based method make use of RGB or depth image. It does not require the user to carry any devices or to attach any sensors on the hand. Therefore, this methodology is getting more consideration nowadays, consequently making the HSR framework simple and easy to be deployed in many applications. We first extracted the frames for each activities from the videos. Specifically, we use transfer learning to get deep image features and trained machine learning classifiers.

Hand Gesture Recognition In the past decade, the computational power of computers has doubled, while the human computer interface (HCI) has not changed too much. When we work with a computer, we are constrained by intermediary devices (keyboards and mice). However, these are inconvenient and have become a bottleneck in human-computer interaction. In our daily life, we use speech to communicate with each other, and use gestures to point, emphasize and navigate. They are the more natural and preferable means to interact with computers for

human beings. To make computers understand this however is not an easy task. Gesture recognition is a topic in computer science and language technology with the goal of interpreting human gestures via mathematical algorithms. Gestures can originate from any bodily motion or state but commonly originate from the face or hand. Gesture recognition can be seen as a way for computers to begin to understand human body language, thus building a richer bridge between machines and humans. Gesture recognition enables humans to communicate with the machine and interact naturally without any mechanical devices. Gesture recognition can be conducted with techniques from computer vision and image processing.

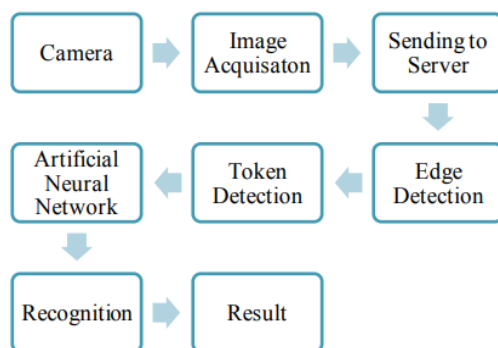


Figure 1: Block diagram of System

Hand and Fingers and Palm Segmentation

The original images used for hand gesture recognition in the work are demonstrated. These images are captured with a normal camera. These hand images are taken under the same condition. The background of these images is identical. So, it is easy and effective to detect the hand region from the original image using the background subtraction method. However, in some cases, there are other moving objects included in the result of background subtraction. The skin color can be used to discriminate the hand region from the other moving objects. The color of the skin is measured with the HSV model. The HSV (hue, saturation, and value) value of the skin color is 315, 94, and 37, respectively. The image of the detected hand is resized to make the gesture recognition invariant to image scale.

The output of the hand detection is a binary image in which the white pixels are the members of the hand region, while the black pixels belong to the background. Then, the following procedure is implemented on the binary hand image to segment the fingers and palm.

Palm Point. The palm point is defined as the center point of the palm. It is found by the method of distance transform. Distance transform also called distance map is a representation of an image. In the distance transform image, each pixel records the distance of it and the nearest boundary pixel.

Inner Circle of the Maximal Radius. When the palm point is found, it can draw a circle with the palm point as the center point inside the palm. The circle is called the inner circle because it is included inside the palm. The radius of the circle gradually increases until it reaches the edge of the palm.

Wrist Points and Palm Mask. When the radius of the maximal inner circle is acquired, a larger circle the radius of which is 1.2 times of that of the maximal inner circle is produced.

Hand Rotation. When the palm point and wrist point are obtained, it can yield an arrow pointing from the palm point to the middle point of the wrist line at the bottom of the hand. Then, the arrow is adjusted to the direction of the north. The hand image rotates synchronously so as to make the hand gesture invariant to the rotation. Meanwhile, the parts below the wrist line in the rotated image are cut to produce an accurate hand image that does not enclose the part of the arm.

Convolutional Neural Network (CNN)

Sign language, as one of the most widely used communication means for hearing-impaired people, is expressed by variations of hand-shapes, body movement, and even facial expression. Since it is difficult to collaboratively exploit the information from hand-shapes and body movement trajectory, sign language recognition is still a very challenging task. This paper proposes an effective recognition model to translate sign language into text or speech in order to help the hearing impaired communicate with normal people through sign language.

Technically speaking, the main challenge of sign language recognition lies in developing descriptors to express hand-shapes and motion trajectory. In particular, hand-shape description involves tracking hand regions in video

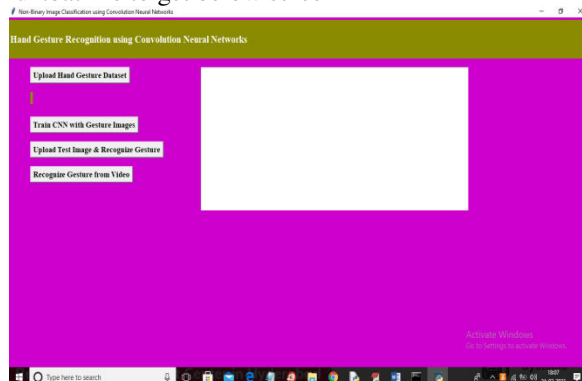
stream, segmenting hand-shape images from complex background in each frame and gestures recognition problems. Motion trajectory is also related to tracking of the key points and curve matching. To address these difficulties, we develop a CNNs to naturally integrate hand-shapes, trajectory of action and facial expression. Instead of using commonly used color images as input to networks like [1, 2], we take color images, depth images and body skeleton images simultaneously as input which are all provided by Microsoft Kinect.

Kinect is a motion sensor which can provide color stream and depth stream. With the public Windows SDK, the body joint locations can be obtained in real-time as shown in Fig.1. Therefore, we choose Kinect as capture device to record sign words dataset. The change of color and depth in pixel level are useful information to discriminate different sign actions. And the variation of body joints in time dimension can depict the trajectory of sign actions. Using multiple types of visual sources as input leads CNNs paying attention to the change not only in color, but also in depth and trajectory. It is worth mentioning that we can avoid the difficulty of tracking hands, segmenting hands from background and designing descriptors for hands because CNNs have the capability to learn features automatically from raw data without any prior knowledge [3].

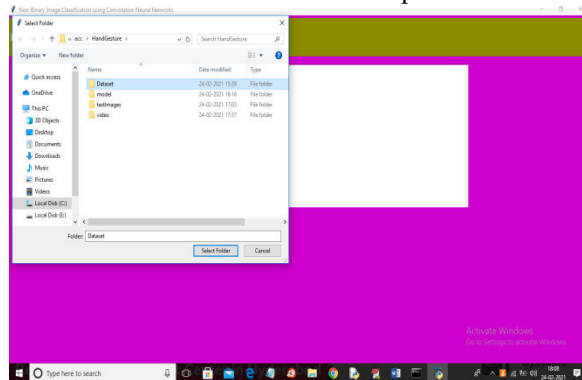
CNNs have been applied in video stream classification recently years. A potential concern of CNNs is time consuming. It costs several weeks or months to train a CNNs with million-scale in million videos. Fortunately, it is still possible to achieve real-time efficiency, with the help of CUDA for parallel processing. We propose to apply CNNs to extract spatial and temporal features from video stream for Sign Language Recognition (SLR). Existing methods for SLR use hand-crafted features to describe sign language motion and build classification model based on these features. In contrast, CNNs can capture motion information from raw video data automatically, avoiding designing features. We develop a CNNs taking multiple types of data as input. This architecture integrates color, depth and trajectory information by performing convolution and subsampling on adjacent video frames. Experimental results demonstrate that 3D CNNs can significantly outperform Gaussian mixture model with Hidden Markov model (GMM-HMM) baselines on some sign words recorded by ourselves.

4. SCREEN SHOTS

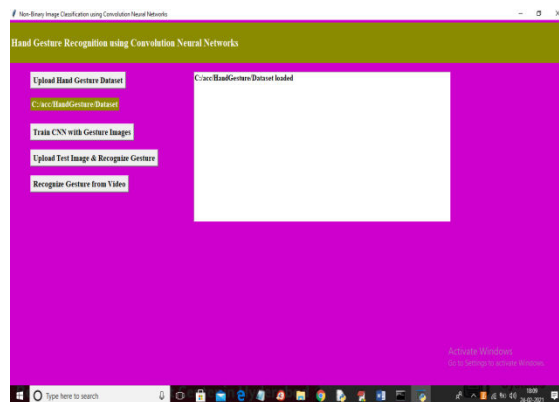
To run project double click on run.bat file to get below screen



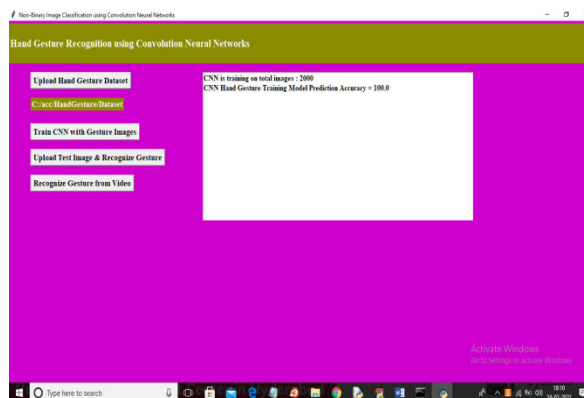
In above screen click on 'Upload Hand Gesture Dataset' button to upload dataset and to get below screen



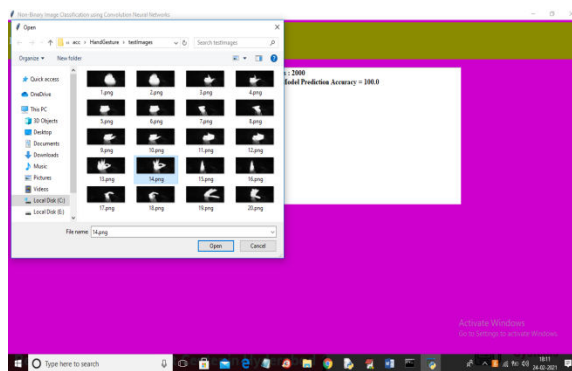
In above screen selecting and uploading 'Dataset' folder and then click on 'Select Folder' button to load dataset and to get below screen



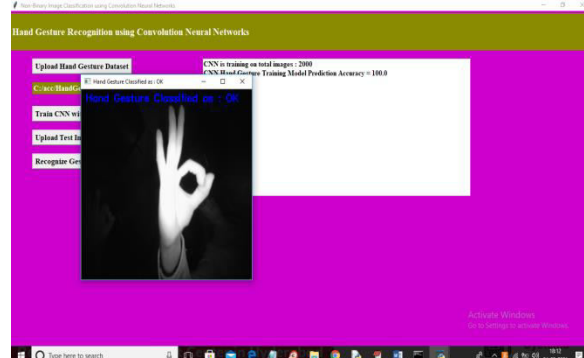
In above screen dataset loaded and now click on 'Train CNN with Gesture Images' button to trained CNN model and to get below screen



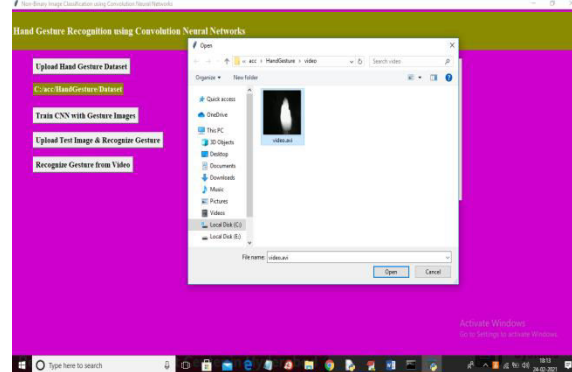
In above screen CNN model trained on 2000 images and its prediction accuracy we got as 100% and now model is ready and now click on 'Upload Test Image & Recognize Gesture' button to upload image and to gesture recognition



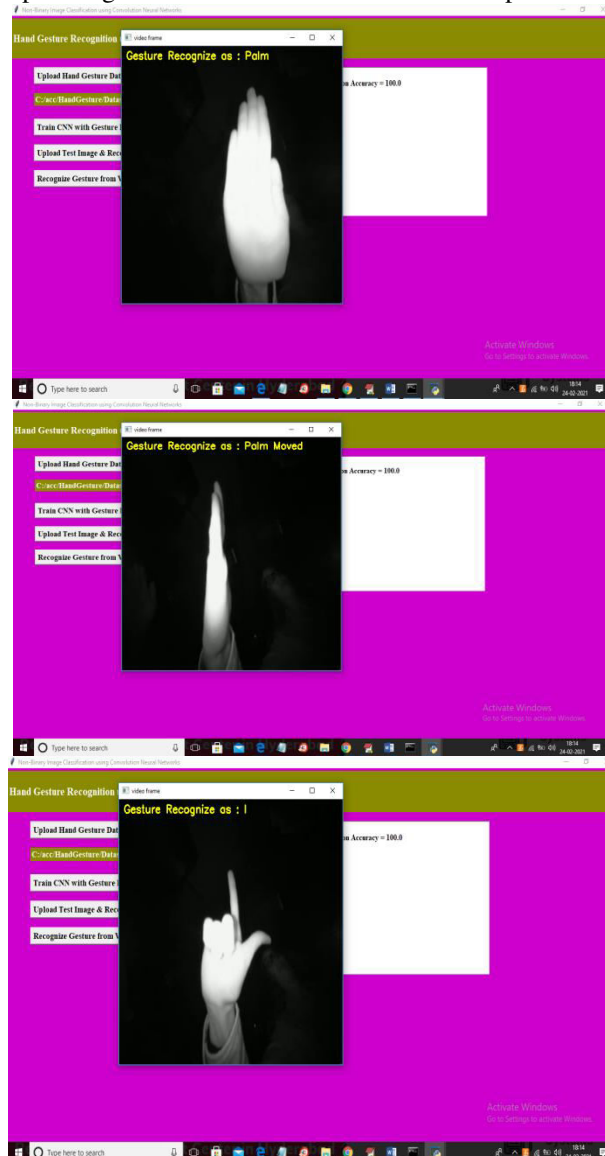
In above screen selecting and uploading '14.png' file and then click Open button to get below result



In above screen gesture recognize as OK and similarly you can upload any image and get result and now click on 'Recognize Gesture from Video' button to upload video and get result.



In above screen selecting and uploading 'video.avi' file and then click on 'Open' button to get below result



In above screen as video play then will get recognition result.

5. CONCLUSION

We created a Convolutional Neural Network (CNN) model specifically designed for the purpose of recognizing sign language. Our model acquires and isolates both spatial and temporal characteristics via the use of 3D convolutions. The advanced deep architecture collects several sorts of information from neighboring input frames and then carries out convolution and subsampling operations independently. The ultimate feature representation amalgamates information from all sources. We use a multilayer perceptron classifier to categorize these feature representations. To provide a comparison, we assess the performance of both Convolutional Neural Network (CNN) and Gaussian Mixture Model-Hidden Markov Model (GMM-HMM) using the same dataset. The experimental findings validate the efficacy of the suggested technique.

REFERENCES

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [2] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei, "Large-scale video classification with convolutional neural networks," in *CVPR*, 2014.
- [3] Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [4] Hueihan Jhuang, Thomas Serre, Lior Wolf, and Tomaso Poggio, "A biologically inspired system for action recognition," in *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*. Ieee, 2007, pp. 1–8.
- [5] Shuiwang Ji, Wei Xu, Ming Yang, and Kai Yu, "3D convolutional neural networks for human action recognition," *IEEE TPAMI*, vol. 35, no. 1, pp. 221–231, 2013.
- [6] Kirsti Grobel and Marcell Assan, "Isolated sign language recognition using hidden markov models," in *Systems, Man, and Cybernetics, 1997. Computational Cybernetics and Simulation., 1997 IEEE International Conference on*. IEEE, 1997, vol. 1, pp. 162–167.
- [7] Thad Starner, Joshua Weaver, and Alex Pentland, "Realtime american sign language recognition using desk and wearable computer based video," *IEEE TPAMI*, vol. 20, no. 12, pp. 1371–1375, 1998.
- [8] Christian Vogler and Dimitris Metaxas, "Parallel hidden markov models for american sign language recognition," in *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*. IEEE, 1999, vol. 1, pp. 116–122.
- [9] Kouichi Murakami and Hitomi Taguchi, "Gesture recognition using recurrent neural networks," in *Proceedings of the SIGCHI conference on Human factors in computing systems*. ACM, 1991, pp. 237–242.
- [10] Chung-Lin Huang and Wen-Yi Huang, "Sign language recognition using model-based tracking and a 3D hopfield neural network," *Machine vision and applications*, vol. 10, no. 5-6, pp. 292–307, 1998.
- [11] Jong-Sung Kim, Won Jang, and Zeungnam Bien, "A dynamic gesture recognition system for the korean sign language (ksl)," *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, vol. 26, no. 2, pp. 354–359, 1996.
- [12] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," *arXiv preprint arXiv:1311.2524*, 2013.
- [13] Ronan Collobert and Jason Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *ICML*. ACM, 2008, pp. 160–167.
- [14] Clement Farabet, Camille Couprie, Laurent Najman, and Yann LeCun, "Learning hierarchical features for scene labeling," *IEEE TPAMI*, vol. 35, no. 8, pp. 1915–1929, 2013.
- [15] Srinivas C Turaga, Joseph F Murray, Viren Jain, Fabian Roth, Moritz Helmstaedter, Kevin Briggman, Winfried Denk, and H Sebastian Seung, "Convolutional networks can learn to generate affinity graphs for image segmentation," *Neural Computation*, vol. 22, no. 2, pp. 511–538, 2010.