

# Sign Language Recognition for Static and Dynamic Gestures

Avi Patel

*Received: 11 June 2021 Accepted: 2 July 2021 Published: 15 July 2021*

---

## Abstract

Humans are called social animals, which makes communication a very important part of humans. Humans use verbal and non-verbal forms of language for communication purposes, but not all humans can give oral speech. Hearing impaired and mute people. Sign language became consequently advanced for them and nevertheless impairs communication. Therefore, this paper proposes a system that uses streams to use CNN networks for the classification of alphabets and numbers. Alphabet and number gestures are static gestures in Indian sign language, and CNN is used because it provides very good results for image classification. Use hand-masked (skin segmented) images for model training. For dynamic hand gestures, the system uses the LSTM network for classification tasks. LSTMs are known for their accurate prediction of time zone distributed data. This paper presents different types of hand gestures, namely two models for static and dynamic prediction, CNN and LSTM.

---

**Index terms**— indian sign language, skin-segmentation, CNN (convolutional neural network), LSTM (long shortterm memory).

## 1 Introduction

Everyone uses language to communicate with others, whether it is English, Spanish, sign language, or touch or smell language. Sign language is the language used by deaf-mute people to talk. It varies from country to country and has its own vocabulary. Indian Sign Language (ISL) is a collection of gestures used by the deaf community in India. These gestures are also different in different parts of India.

It is always a challenge for normal people to communicate with deaf-mute people and vice versa. Sign language translation is the solution to this problem. It provides a bridge of communication between the community at large and the deaf-mute community. There are two main methods for the recognition of sign language, glove-based and computer vision-based [1]. In this article, a computer vision-based approach to interpreting ISL in two different ways is discussed. ISL letter recognition includes camera frame extraction, hand masking, feature extraction and classification recognition. This is to identify the alphabet from a single frame. The second method is to recognize gestures through words. The camera frame sequence is used to recognize gestures. It consists of the same modules as letter recognition, but it uses a series of frames instead of one frame. This article focuses on ISL recognition through deep learning and computer vision. The rest of thesis is organized as follows; the second part presents the related work done in gesture recognition. The third part contains the methodology of the two methods of recognition of the Indian sign language. The first method is suitable for static gestures and the second method is suitable for dynamic gestures. Discussion of the results and conclusions are explained in Sections 4 and 5, respectively.

## 2 II.

## 3 Related Work

Many techniques have been developed to recognize sign language. There are two main approaches that use tracking sensors or computer vision to track various movements. Much research has been done on sensor-based approaches using gloves and wires [1,2,3]. Therefore, it is inconvenient to wear these devices continuously. Additional work will primarily focus on computer vision-based approaches.

## 5 SKINSEGMENTATION:

---

44 A lot of work has been done using a computer vision based approach. The authors have proposed various  
45 methods of recognizing sign language using CNN (Convolution Neural Network), HMM (Hidden Markov Model)  
46 and contour lines [4,5,6,7]. Different methods are used to split images such as HSV and color difference images  
47 [4,5]. The authors proposed an SVM (Support Vector Machine) method for classification [6,8]. Archana and  
48 Gajanan also compared different methods for partitioning and feature extraction [9]. All previous treatises have  
49 successfully recognized the ISL alphabet. But in reality, deaf or mute people use speech gestures to convey  
50 messages. If the word has a static action, you can use these previous methods to check the word.

51 ISL Many words required hand movements. The image classification method is not a simple image classification  
52 technique, but is suitable for identifying these dynamic gestures. Video-based action recognition has already  
53 attracted attention in several studies [10,11,12]. Instead of capturing the color image data for each frame of the  
54 video, some researchers performed differences between successive frames and randomly provided these segments  
55 to the TSN (Temporal Segment Network) [11]. Sun, Wang and Yeh use LSTM (Long Short-Term Memory) [13] to  
56 describe video classification and captions. Juilee, Ankita, Kaustubh and Ruhina used video to suggest a method  
57 of hydration recognition in India [14]. As a result of searching the sign language recognition system, most  
58 studies use static sign language gestures and video recognition techniques to study dynamic gesture identification  
59 and only perform video classification for various actions discovered.

60 III.

### 4 Methodology a) Static gesture classification

62 Experiments were performed on the data set provided by [15]. The dataset contains 36 folders representing 09  
63 and A-Z, each consisting of an image of a hand subdivided by the corresponding alphanumeric skin color. There  
64 are 220 images of 110 x 110 pixels each for each alphanumeric character. Figure 1 shows an image of each label  
65 in the data set. After training the model, predict the output by performing the following steps:

66 Frame Extraction: Uses OpenCV library to capture video from webcam for live prediction. After capturing  
67 the video, take a single frame and define a region of interest (ROI) in that frame. The area of interest is the area  
68 in which a person runs a stream.

### 5 Skinsegmentation:

70 The ROI of the frame is transformed into a hand-masked image to provide to the model for predictive purposes.  
71 First, you need to blur the image to reduce noise. This is done by applying Gaussian Blur. After blurring ROI is  
72 converted to HSV color scale in RGB. Converting an image to the HSV color scale helps detect better skin than  
73 RGB. Next, lower and upper limits are set for skin extraction. Here, (108, 23, 82) was used in the low range and  
74 (179, 255, 255) was used in the high range. This range offers us the best results. After selecting a range, compare  
75 the values of each pixel and if the value is not within the range, it will be converted to black, otherwise it will  
76 be converted to white pixels. This provides us with handmasked images. Still, the hand-masked image is noisy  
77 and the edges are not aligned. To solve this problem, use the Dilate and Erode features available in OpenCV to  
78 smooth the edges. Prediction: The ROI of the frame is converted to a handmasked image. This hand-masked  
79 image is provided as input to the CNN model for prediction. 09 or AZ is provided for the output of the original  
80 frame, but it is a predicted value. But this leads to another problem, the frame output keeps blinking. To solve  
81 this problem, we used a 25-frame forecast and used the maximum forecast class as the output.

82 Figure 2 shows a handmasked image of the alphabet L and the final output. Neural networks can help predict  
83 complex data and values most of the time. Input is not related to time or is not required in chronological order.  
84 This is the case for static gestures in ISL, so a multi-layer CNN architecture is sufficient. However, for dynamic  
85 gestures, you cannot perform CNN since you have to keep the previous state. Therefore, LSTM networks  
86 are useful in this case. LSTM is an RNN (Recurrent Neural Network) type that has a structure similar to a  
87 chain of repeating modules that is useful for learning long-term dependencies from sequential data. 3. The input  
88 continuously delivers a sequence of 8 frames/images extracted from images in the training dataset. Apply an  
89 RGB difference filter before serving these 8 frames as input. The RGB difference subtracts the current frame  
90 from the previous frame. Therefore, only the changed pixels remain in the frame and the remaining still images  
91 are deleted. In this way, it helps to capture time-varying visual features. Here in our case it helps to capture  
92 the gesture pattern and the background is also removed so it becomes independent in a variety of background  
93 scenarios.

94 If so, these frame sizes change to 224 x 224 pixels. This is because the next layer is a MobileNetV2 layer  
95 that only accepts image sizes up to 224 x 224 pixels. As a pre-trained model, we used the weight of 'Imagenet'  
96 and MobileNetV2. MobileNetV2 can be used as a pre-trained model used for image segmentation, eliminating  
97 the task of building CNN models for image segmentation. Separate the Mobile Net V2 layer for passing these 8  
98 frames with Time Distributed layer used. Here I use the Time Distributed Global Average Pooling layer as I need  
99 to flatten the frames to insert a series of frames into the LSTM. Finally, there is a multi-layer LSTM structure  
100 with several dropouts and a fully connected layer to reduce the sum of overcharges. LSTMs help recognize  
101 pattern formation with dynamic / moving hand gestures. Finally, SGD is used in optimization programs because  
102 it provides better results when the available data set is low. Adam provides good results even when the dataset  
103 is large.

## 105 6 Results and Discussion

### 106 7 a) Static gesture classification

107 During checking out of the static hand gesture version and figuring out the greatest architecture, 10epoch models  
108 were each trained using various optimizations such as RMSProp, SGD, Adam, etc. By the way, RMSProp gave  
109 the best results with a precision of 73.6°. A graph of accuracy and time is shown in Figure 4 Skin segmentation is  
110 an integral part of a system for predicting static hand gestures. It was concluded that the lower range (108,23,82)  
111 and the higher range (179,255,255) would give the best results. Figure 5-6 shows the gestures predicted to be skin  
112 segmentation. There are some limitations to using skin segmentation to recognize static hand gestures. Most  
113 importantly, you need a skin-free background. Predictions are wrong because the background contains colors  
114 in the skin color range and it is difficult to hide the skin. For example, if the background is a shade of yellow  
115 that falls within our range, this problem will occur. The second problem is a stream of similar shape. The equal  
116 gestures with alphabets and numbers overlap. For example, the alphabet "V" and the number "2" have the same  
117 gesture and cannot be properly distinguished by the system. There is also a similar hand movement problem  
118 that reduces accuracy. For example, the letters 'M' and 'N' are very similar. Other similar pairs are 'FX' and  
119 'II'. Removed static parts of the frame sequence using RGB differences to overcome the background color issue.  
120 It also leaves the moving hand in the frame, which helps detect hand gesture patterns. The only problem with  
121 this approach is that if the background is moving, the sequence of frames will also have a background, which will  
122 affect the prediction accuracy. You can also add more videos to different backgrounds and people's datasets for  
123 greater accuracy.

124 V.

## 125 8 Conclusion and Destiny Scope

126 The Deafmute community is faced with communication challenges every day. This white paper describes two  
127 methods for recognizing hand gestures: static gestures and dynamic gestures. For static gesture classification, a  
128 CNN model is implemented that classifies the motions alphabetically (AZ) and numerically (09) with a precision  
129 of 73. Use hand mask skin subdivision with the model For dynamic gestures, we trained a model using multi-layer  
130 LSTM using 12word MobileNetV2 and gave very satisfactory results with an accuracy of 85°. For future work  
131 with static gestures, another approach to skin segmentation that does not rely on skin color can be built. For  
132 dynamic gestures, you can increase the size of data sets with different backgrounds.

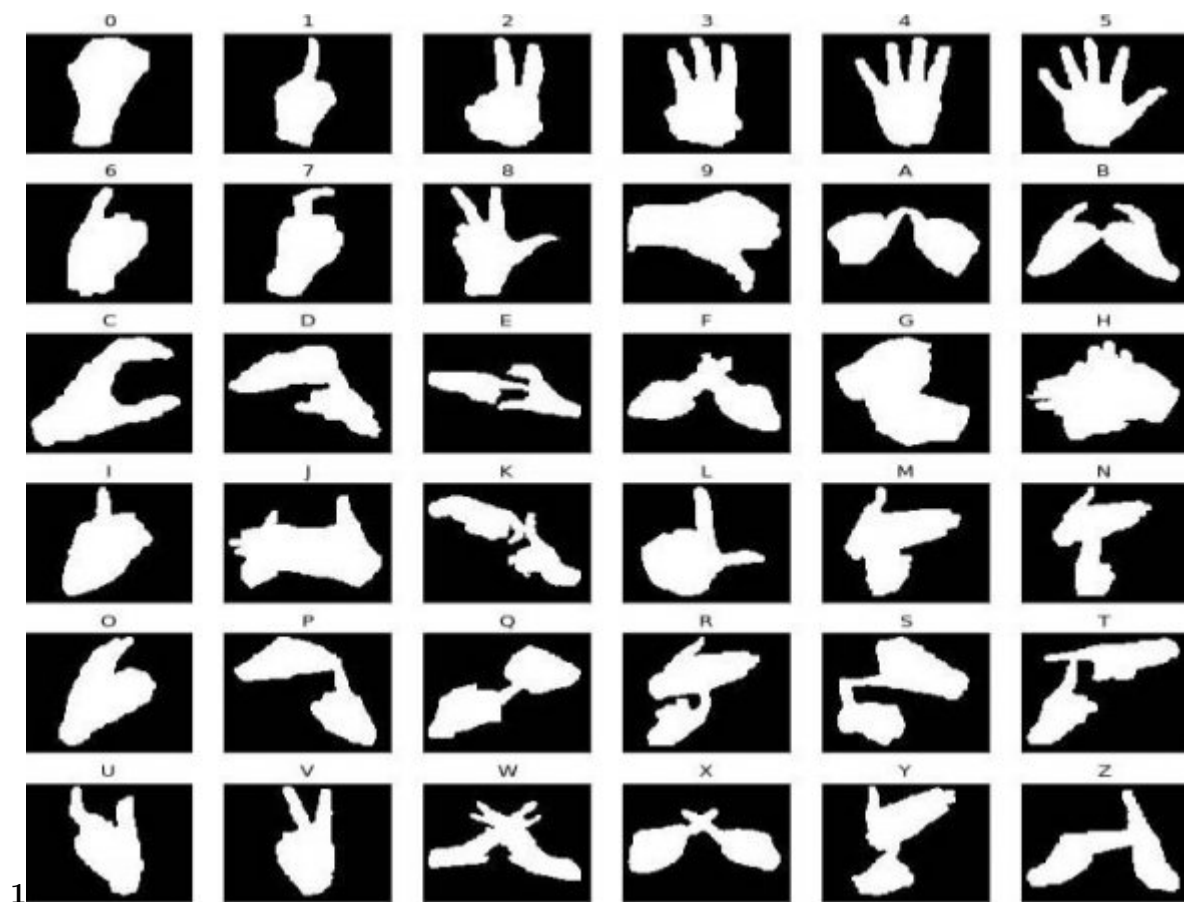


Figure 1: Figure 1 :

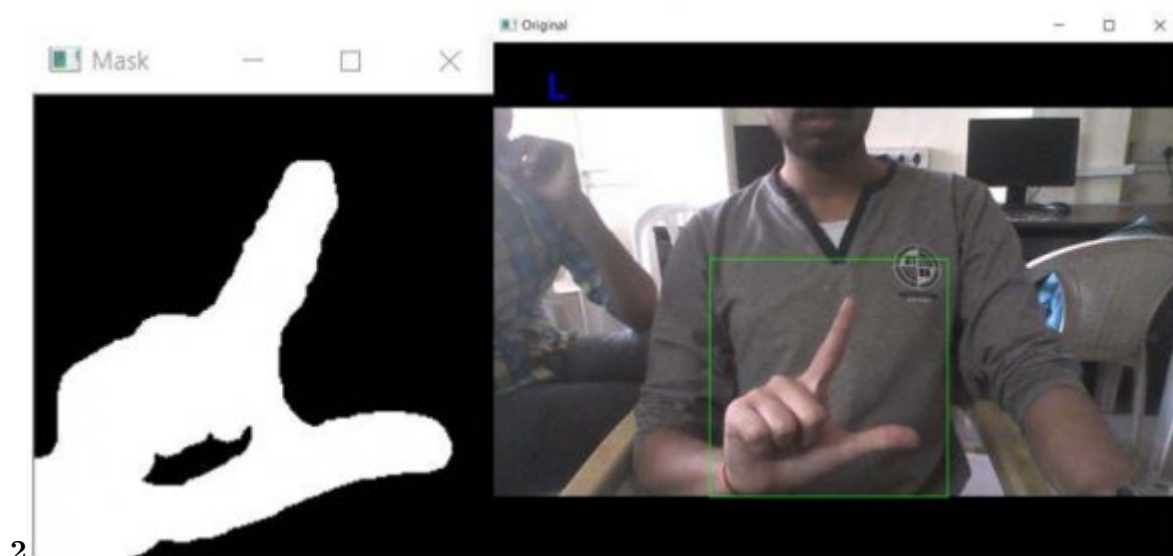


Figure 2: Figure 2 :

Model: "sequential\_1"

Layer (type)	Output Shape	Param #
time_distributed_1 (TimeDist)	(None, 8, 7, 7, 1280)	2257984
time_distributed_2 (TimeDist)	(None, 8, 1280)	0
lstm_1 (LSTM)	(None, 8, 64)	344320
dropout_1 (Dropout)	(None, 8, 64)	0
lstm_2 (LSTM)	(None, 64)	33024
dense_1 (Dense)	(None, 64)	4160
dropout_2 (Dropout)	(None, 64)	0
dense_2 (Dense)	(None, 24)	1560
dropout_3 (Dropout)	(None, 24)	0
dense_3 (Dense)	(None, 5)	125

Total params: 2,641,173  
Trainable params: 2,607,061  
Non-trainable params: 34,112

Figure 3: Figure 3 :

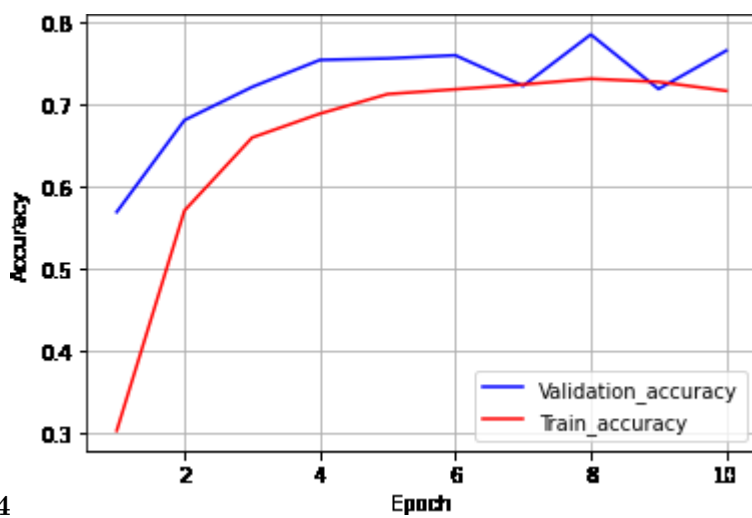


Figure 4: Figure 4 :



5

Figure 5: Figure 5 :



6

Figure 6: Figure 6 :

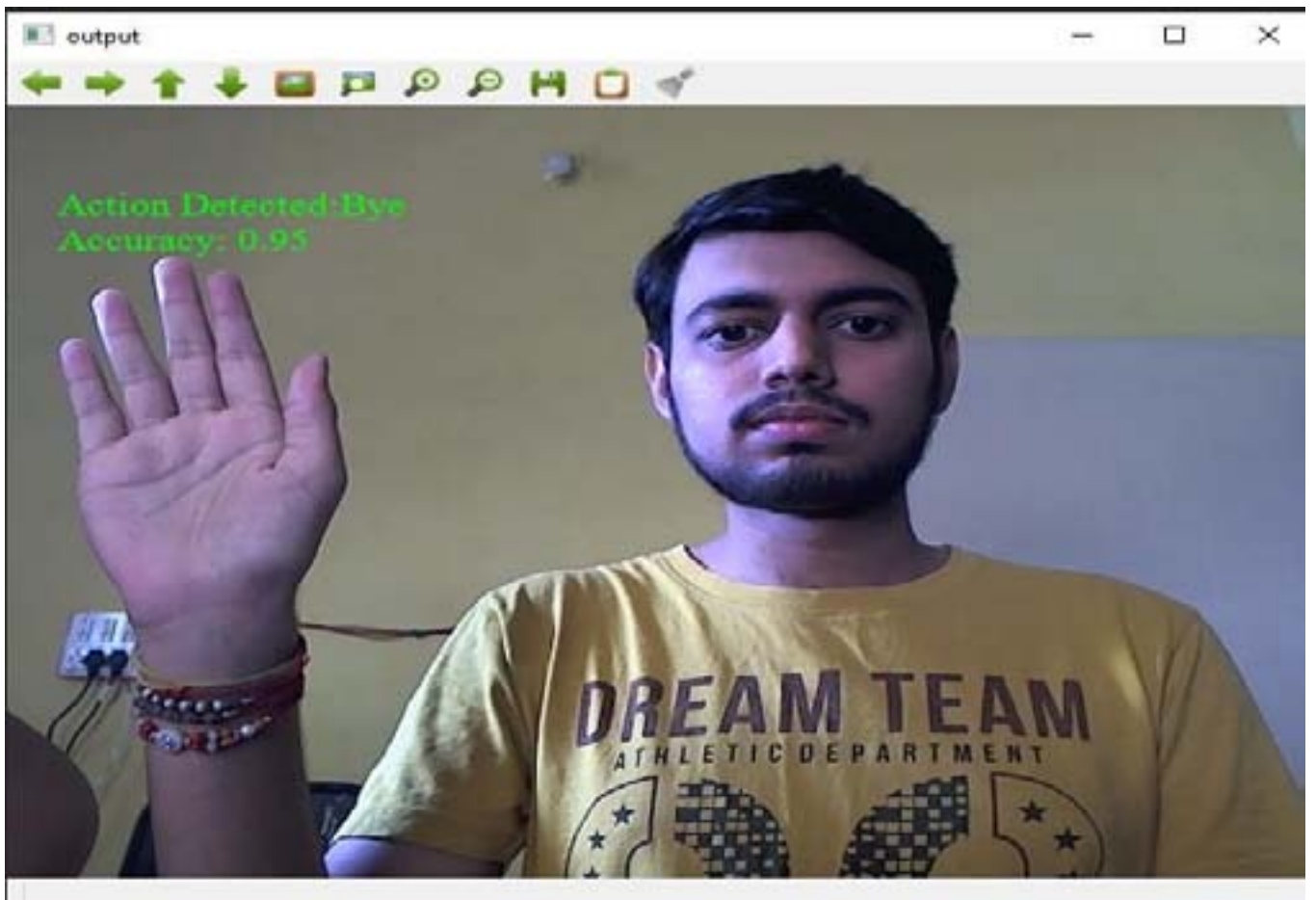


Figure 7:

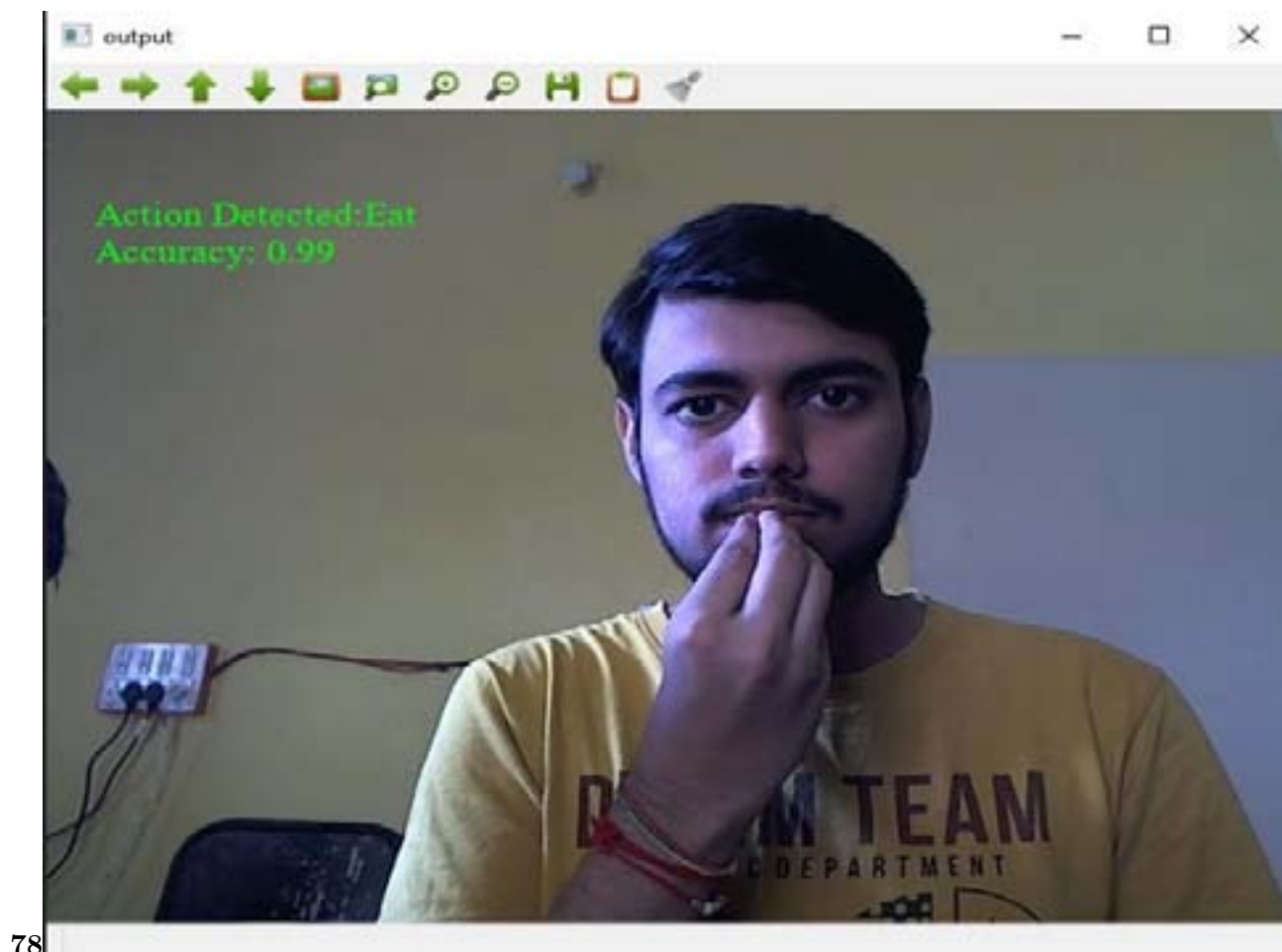


Figure 8: Figure 7 :Figure 8 :

1.1

Property	Value
ConvolutionLayer	3Layers (32,64,128 nodes)
ConvolutionLayer(KernelSize)	3,3, 2
MaxPoolingLayer	3 Layers-(2, 2)
FullyConnected Layer	128nodes
OutputLayer	36nodes
ActivationUsed	Softmax
Optimizer	RMSPProp
Hyperparameters	
Learningrate	0.01
No.ofepochs	10

Figure 9: Table 1 below . Table 1 :



- 
- 133 [Springer ()] , Cham Springer . 2016.
- 134 [Jaya and Rajendran ()] ‘Hand-Talk Assistive Technology for the Dumb’. T Jaya , V Rajendran . *International*  
135 *Journal of Scientific Research in Network Security and Communication (IJSRNSC)* 2018. 6 (5) p. .
- 136 [Lele ()] ‘Image Classification Using Convolutional Neural Network’. N S Lele . *International Journal of Scientific*  
137 *Research in Computer Science and Engineering (IJSRCSE)* 2018. 6 (3) p. .
- 138 [Mohandes et al. ()] ‘Image-Based and Sensor-Based Approaches to Arabic Sign Language Recognition’. M  
139 Mohandes , M Deriche , J Liu . *IEEE Transactions on Human-Machine Systems* 2014. 44 (4) p. .
- 140 [Raheja et al. ()] ‘Indian Sign Language Recognition Using SVM’. J L Raheja , A Mishra , A Chaudhary . *Pattern*  
141 *Recognition and Image Analysis*, 2016. 26 p. .
- 142 [Rege et al. ()] ‘Interpretation of Indian Sign Language through Video Streaming’. J Rege , A Naikdalal , K  
143 Nagar , R Karani . *International Journal of Computer Science and Engineering (IJCSE)* 2015. 3 (11) p. .
- 144 [Karpathy et al. ()] ‘Large-scale Video Classification with Convolutional Neural Networks’. A Karpathy , G  
145 Toderici , S Shetty , T Leung , R Sukthankar , L Fei-Fei . *the Proceedings of the IEEE Conference on*  
146 *Computer Vision and Pattern Recognition (CVPR)*, 2014. p. .
- 147 [Gupta et al. ()] ‘Sign Language Problem and Solutions for Deaf and Dumb People’. P Gupta , A K Agrawal ,  
148 S Fatima . *the Proceedings of the International Conference on System Modeling & Advancement in Research*  
149 *Trends*, (Moradabad, India) 2014.
- 150 [Ramkumar et al. ()] ‘Sign Language Recognition using Depth Data and CNN’. L K Ramkumar , S Premchand ,  
151 G K Vijayakumar . *SSRG International Journal of Computer Sciences and Engineering (SSRG-IJCSE)* 2019.  
152 6 (1) p. .
- 153 [Nikam and Ambekar ()] ‘Sign Language Recognition Using Image Based Hand Gesture Recognition Techniques’.  
154 A S Nikam , A G Ambekar . *the Proceedings of the Online International Conference on Green Engineering*  
155 *and Technologies (IC-GET)*, (Coimbatore, India) 2016. p. .
- 156 [Ghotkar and Kharate ()] ‘Study of Vision Based Hand Gesture Recognition Using Indian Sign Language’. A S  
157 Ghotkar , G K Kharate . *International Journal on Smart Sensing and Intelligent Systems* 2014. 7 (1) .
- 158 [Wang et al.] ‘Temporal Segment Networks: Towards Good Practices for Deep Action Recognition’. L Wang ,  
159 Y Xiong , Z Wang , Y Qiao , X Lin , L Tang , Vangool . *the Proceedings of the European conference on*  
160 *Computer Vision*, p. .
- 161 [Simonyan and Zisserman ()] ‘Two-stream Convolutional Networks for Action Recognition in Videos’. K Si-  
162 monyan , A Zisserman . *Advances in Neural Information Processing Systems*, 2014. p. .
- 163 [Sun et al. ()] ‘Video understanding: from video classification to captioning’. J Sun , J Wang , T C Yeh . *the*  
164 *Proceedings of the Computer Vision and Pattern Recognition*, 2017. p. . Stanford University
- 165 [Patel and Patel (2019)] ‘Vision Based Realtime Recognition of Hand Gestures for Indian Sign Language using  
166 Histogram of Oriented Gradients Features’. Pradip Patel , Narendra Patel . *International Journal of Next-*  
167 *Generation Computing* July 2019. 10 (2) .
- 168 [Zhu and Sheng ()] ‘Wearable Sensor-Based Hand Gesture and Daily Activity Recognition for Robot-Assisted  
169 Living’. C Zhu , W Sheng . *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and*  
170 *Humans*, 2011. 41 p. .