

## Hand Sign Recognition using CNN

D. Bhavana<sup>a,\*</sup>, K. Kishore Kumar<sup>b</sup>, Medasani Bipin Chandra<sup>a</sup>, P.V. Sai Krishna Bhargav<sup>a</sup>,  
D. Joy Sanjana<sup>a</sup>, and G. Mohan Gopi<sup>b</sup>

<sup>a</sup>Department of Electronics and Communication Engineering, Koneru Lakshmaiah Education Foundation, Guntur, 522502, India

<sup>b</sup>Department of Mechanical Engineering, Koneru Lakshmaiah Education Foundation, Guntur, 522502, India

---

### Abstract

Our aim is to produce a model that can recognize hand gestures and signs. We will train a model for the purpose of sign language conversion, a simple gesture recognizing model; this will help people converse with people who are innately deaf and mentally disabled. This project can be implemented in several ways such as KNN, Logistic Regression, Naïve Bayes Classification, Support vector machine and can be implemented with CNN. The method we have chosen is CNN as it gives better accuracy compared to the rest of the methods. A computer program is developed using python language which is used to train the model based on the CNN algorithm. The program will be able to recognize hand gestures by comparing the input with preexisting dataset formed using the American sign Language. We will be able to convert Sign Language into text as output for users to recognize the signs presented by the sign language speaker. This model is implemented in Jupyter Lab, an extension to the platform Anaconda documentation. To further improve, we will also add / integrate the inputs into black and white and take input from camera after using the method of Background subtraction. With the mask set to detect the human skin, this model will not require a plain background to function and can be implemented using a basic camera and a computing device.

*Keywords:* American sign recognition; CNN; background subtraction; OpenCV; hand tracking and segmentation; feature extraction

© 2021 Totem Publisher, Inc. All rights reserved.

---

### 1. Introduction

In the modern world, we have computers that work at very high speeds, making significantly huge amount of calculations in a split second. Now we humans are trying to achieve a goal where we want the computer to start thinking [1] and working like a human being. This requires the most basic property 'learning'. This takes us to artificial intelligence. According to AI, the computer starts or begins performing tasks on its own without human intervention. For this to happen, the computer needs to learn how to react to certain inputs and situations. The computer needs to be trained with huge amounts of data; the data that is used to train depends upon the preferred outcome and the working of the machine. We develop a computer model that can detect hand gestures made by human beings; there are so many applications that we see in our day to day life that work on hand gestures. Take a look at the console in our living room; hook it up with a sensor and we can start playing tennis with our hand. We realize a gesture recognition model that converts sign language into speech [2]. There are so many devices that exist that depend on gesture recognition, such as for security or entertainment purposes. Sign language may be a vision-based language that uses an amalgamation of various types of visuals, gestures, and the shape of the hands, fingers and their alignment, orientation, and movement of hand and body, eyes, lip, and complete facial movements and expressions. Like the spoken language, regional variants of signing also exist, e.g., Indian signing (ISL), American Sign Language (ASL), and Portuguese Sign Language, spelling each alphabet with our fingers, particular vocabulary is maintained for words and sentences, with the help of hands and body movement, facial expressions, and lip movement. Sign language can be isolated also as continuous. In isolated sign language, people communicate with gestures that represent a single word, while continuous sign language may be a sequence of gestures that generate a meaningful sentence [3]. All the methods for recognizing hand gestures are often broadly classified as vision-based and based on measurements picked up by sensors inside gloves. This vision-based method involves interaction between human and computer for gesture recognition, while the latter (glove based) method depends on the use of programmed external hardware for gesture recognition [4]. We will not be using a glove for our project; our model can recognize gestures made by bare hands. In this project, sign language recognition

---

\* Corresponding author.

E-mail address: [bhavanaee@kluniversity.in](mailto:bhavanaee@kluniversity.in)

is done using OpenCV. It recognizes the hand gestures made by the user via a webcam; the output will be the text displayed on the screen. Our project aims to help people who are not aware of sign language by detecting the symbols made by a person and convert them into readable texts. We can achieve this with the help of Machine learning, specifically CNN [5]. We aim to train a model that predicts the text by taking an image as an input from the web camera/phone camera (using IP camera software and OpenCV). This model will be able to function at an accuracy of at least 75%, as the model is trained using a renowned dataset (ASL), enough data is available for the training algorithm to produce a precise and accurate model. The code used to develop the model is written in Python language, a simple computer can be used to realize the project without the need for high processing units and GPU's [6]. A platform called Anaconda Documentation is used as a base for the software Jupyter Lab, where this model is trained and implemented [7]. The various concepts involved and the procedure on how this project is carried out is discussed in a detailed fashion in the upcoming sections of the paper.

### 1.1. Design and Description

**OpenCV-** We use OpenCV in Artificial Intelligence, face recognition, machine learning, etc. Computer Vision (CV) is an open-source Library in Python [8]. Computer Vision, helps the computer study objects through its eyes, sensors that allow the computer to take in visual data. Content of digital images such as videos and photographs are taken in by the machine through this platform. The objective of computer vision is to understand and extract data from the images [9]. It extracts the description from the pictures. It can be an object, a text description, and may also be a three-dimension model, and so on. For example, gaming consoles can be facilitated with computer vision, which will be able to identify different gestures made by the player with his arms and legs. And depending upon the movements made by the user, the computer performs the designated task in the game. The computer performance keeps on improving based on the time the data model is set to use on a single subject (one frequent user) signs, and so on, and acts accordingly [10]. Now machines with the help of computer vision has more of an edge to start thinking and acting like humans to some extent and can perform tasks like humans or similar efficiency. There are two tasks involved, object classification and object identification. In object classification, we train a model on a dataset of certain pre-defined objects; the model classifies these new objects arranging them to one or more of your training categories. In object identification, the object is identified based on its unique features which are in turn belonging to a certain class. OpenCV is used for computer Vision for the following reasons:

- 1) OpenCV is available free of cost.
- 2) Since the OpenCV library is written in C/C++, it is quite fast and can be used with Python.
- 3) It requires less RAM to usage, around 60-70 MB.
- 4) Computer Vision is as portable as OpenCV and can run on any device that C can run on.

**Convolutional Neural Networks-** Our entire model cannot be realized without this concept. CNN comes under deep neural networks belonging to one of its classes. We can see CNN applied in many areas but most of these areas including visual images. CNN is made up of several neurons whose value i.e., weights can be altered to serve the required purpose. These weights and biases are learnable [11]. Dot products are performed between the neurons resulting in non-linearity. The entire network expresses a single function; the main difference between the ordinary neural networks and convolutional neural networks is the assumption made by CNN, that all the inputs are explicit images. Thus, this will be best for training the model since our project revolves around images. Due to this, CNN will be able to take advantage the architecture can be constrained in a sensible way with images as input explicitly; an arrangement of neurons is done in 3 dimensions: width, depth and height. Depth means the volume required for activation.

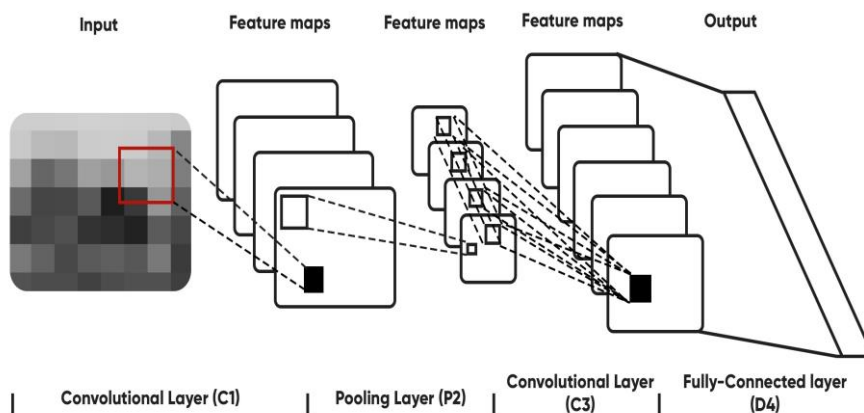


Figure 1. Structure of Convolutional Neural Network

As shown in Figure 1, we know that a convolutional network consists of a number of layers and the volume existing is transferred to another form with the help of a function(differentiable). The layers used to build CNN architecture are: pooling layer, convolutional layer and fully-connected layer. Stacking these in proper order will result in a convolutional network architecture [12]. Now these layers extract features that are unique to a property from the image, the loss function needs to be minimized. With this equation. it can be done so positive classes of the sample are denoted by M. Each positive class's CNN score is denoted by  $S_p$ ; the scaling factor is denoted by  $1/M$

$$CE = \frac{1}{M} \sum_p^M -\log \left( \frac{e^{S_p}}{\sum_j^C e^{S_j}} \right) \quad (1)$$

**Background Subtraction-** This means subtracting the background, i.e., taking it out of the picture. This technique involves separating the foreground with background elements. This is achieved with the help of a mask that is generated as per the user's desire [13]. We use this procedure to detect through static cameras. Background subtraction is vital for tracking an object, this can be realized in many ways.

### 1.2. How Does the Computer Recognize the Image

There is a lot of information around us and this is picked up by our eyes selectively, which is different for each person depending upon their preference. But the machines, on the other hand, see everything and take in every image and then convert the data into 1's and 0's which the computer can understand. How does this conversion take place? Pixel. This is the smallest unit of a digital image, which can be displayed or projected onto a display device. The image has various intensity levels at various positions and these levels are represented by numbers, for our images we have shown values consisting of only one value( grayscale), the value is assigned automatically by the computer based on the intensity of the darkness or level [14]. The two common ways to identify images are by Grayscale or RGB. In grayscale, picture a black and white image; these images consist of only two colors. The black color is assumed to be or treated as a measurement as the weakest intensity meaning white is the strongest intensity. The values are allotted by the computer based on the darkness levels [15]. In RGB, Red, green, blue are termed as RGG. All these colors together make an actual color. Any color that exists on this planet can be defined only using these three colors. The computer checks and extracts value from each pixel and assigns the results in an array.

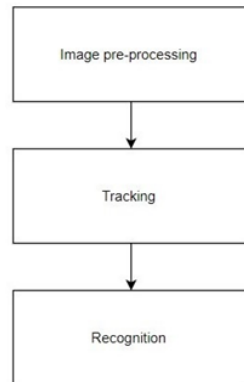


Figure 2. Gesture Recognition System

## 2. PROPOSED METHODOLOGY

As shown in Figure 2, we follow a process to devise a model that uses CNN to predict the text of a hand gesture/symbol. We also use methods like background subtraction (see Figure 3) to improve the efficiency of the model by eliminating the light/ambience of the background.

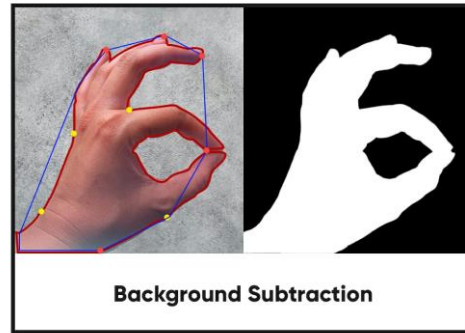


Figure 3. Background Subtraction

To build a Sign Language Recognition system, three things are needed:

- 1) Dataset as shown in Figure 4
- 2) Model (In this case we will use a CNN)
- 3) Platform to apply our model (We are going to use OpenCV)

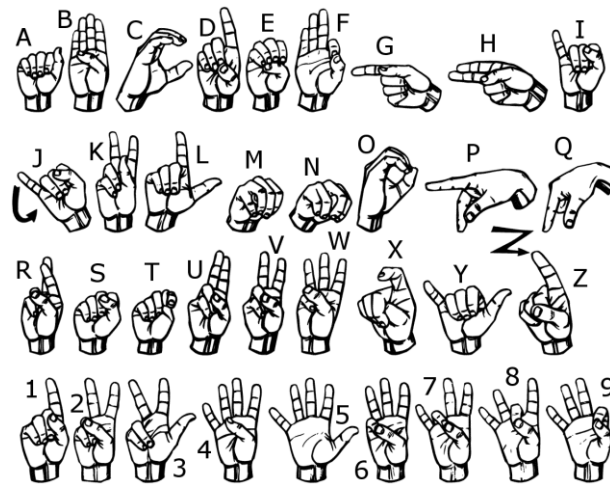


Figure 4. American Sign Language [ASL benchmark dataset]

To train a deep neural network, proper hardware is required (e.g., a powerful GPU). A powerful GPU will not be needed for this project. But still, the best way to go through with this is by using online platforms like Google Collab. We will use the Modified National Institute of Standards and Technology dataset [16]. Our dataset consists of many images of 24 (except J and Z) American Sign Language alphabet; there is a total of 784 pixels per image, meaning each image has a size of 28x28.

Our dataset is in CSV (Comma-separated values) format; train\_X and test\_Y contain the pixel values. train\_Y and test\_Y contains the label of image. Moving on to pre-processing, both the trained x and y consists an array of all the pixel values. An image is to be created from these values. Our image size is 28x28; we divide the array into 28x28 pixel groups. We use this dataset to coach our model. We use CNN (Convolutional Neural Network) to recognize the alphabets. We use keras. Like any other CNN, our model consists of a couple of layers like Conv2D and the MaxPooling that is followed by fully connected layers. The first Conv2D layer takes the shape of input image and the last fully connected layer provides us with the output for 26 alphabets [17]. We are using a Dropout after 2nd Conv2D layer to regularize our training. Soft Max is used as the activation function in the final layer [18], which will give us the probability for each alphabet as an output. The proposed model is given Table 1.

Table 1. Model Overview

| Layer (Type)                 | Output Shape    | Param# |
|------------------------------|-----------------|--------|
| Conv2d_1(Conv2D)             | (None,28,8,8)   | 80     |
| Max_pooling_1(MaxPooling2)   | (None,14,14,8)  | 0      |
| Conv2d_2(Conv2D)             | (None,14,14,16) | 1168   |
| Dropout_1(Dropout)           | (None,14,14,16) | 0      |
| Max_pooling2d_2(MaxPooling2) | (None,3,3,16)   | 0      |
| Dense_1(dense)               | (None,3,3,128)  | 2176   |
| Flatten_1(Flatten)           | (None,1152)     | 0      |
| Dense_2(dense)               | (none,26)       | 29978  |

The SGD optimizer is utilized for compiling our model. You may decrease the epochs to 25. Check the accuracy of the trained model. Open CV: We must create a window to take the input from our webcam. The images that are taken as input should be a grayscale image(28x28). Because we trained our model on 28x28 size image, the alphabet from the input image is to be predicted. Our model will give outputs as integers rather than alphabets; that is because the labels are given as integers (1-A, 2-B, 3-C, etc.). With maximum value of accuracy, our model should be able to recognize the alphabet, with the help plain background and descent lights, without any hindrance.

### 2.1. Steps Involved

Hand Recognition: I window is required to take the input image from our camera (web cam) so we must first create it. As the requirements of the image dimension mentioned earlier(28x28), a greyscale image of these dimensions must be taken. Because the model has been trained for those dimensions only, following the same dimensions will result in accurate results. Segmentation: After tracking the hands, the next step is to segment the hands from the background. The HSV color model is used for this purpose. Feature Extraction: Several general features are extracted; the relationship between features and classes is inferred by an appropriate classifier. Gesture Recognition: After extracting features of the input character, we search its features in the database and consider the most similar features as the result.

### 2.2. Execution of the Model

The execution includes the following steps:

- Capture the image through the camera.
- Before sending the captured image to the server, store some of the test images of the person.
- In the server, give it to the learning API.
- Now send this captured image through the API to the server.
- To send this image to the server we use REST API
- REST stands for Representational State Transfer, meaning when a REST API is called.
- The server will transfer the representation of the state of the requested source to the client.
- After sending this image to the server, the server will finalize the results i.e., gives back the name and employee id of the person in that image using learning API and data sets to the REST API.
- REST API gives back the information passed by the server to the clients.

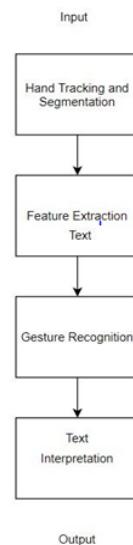


Figure 5. Execution Process



### 3. Outputs

After opening the platform to perform the execution of the model as shown in Figure 5, we used jupyter lab to go through with the execution of our model launched by anaconda navigator as the base software when the code is run. We are presented with three windows: mask window, the camera window and the trackbar window. It is advised to have a plain background before going through with the recognition process, but if that seems to be an issue we can always change the HSV colour scheme in the trackbar to match that ratio to that of the human skin. So even if the background is not plain or a solid color, we will be able to make the model work at good efficiency. I am using a basic dell model laptop and working from home because of the pandemic. Due to this, I do not have many resources at my expense. I decided to go through the model using a computer software which enables the built-in webcam to act as the input device. To resolve the background issue, I set the trackbar levels to upper HSV- 0,58,50 and lower HSV- 30,255,255. With this setting we will be able to deliver a black background where our skin will be detected in white. We move on to checking the accuracy and the precision of the model. As mentioned in Figure 3, our model can detect American sign language and display the respective alphabet or character in text on screen. I have tested the model and am displaying four random alphabets as detected by the computer. Parts of the source code are also visible in the background.

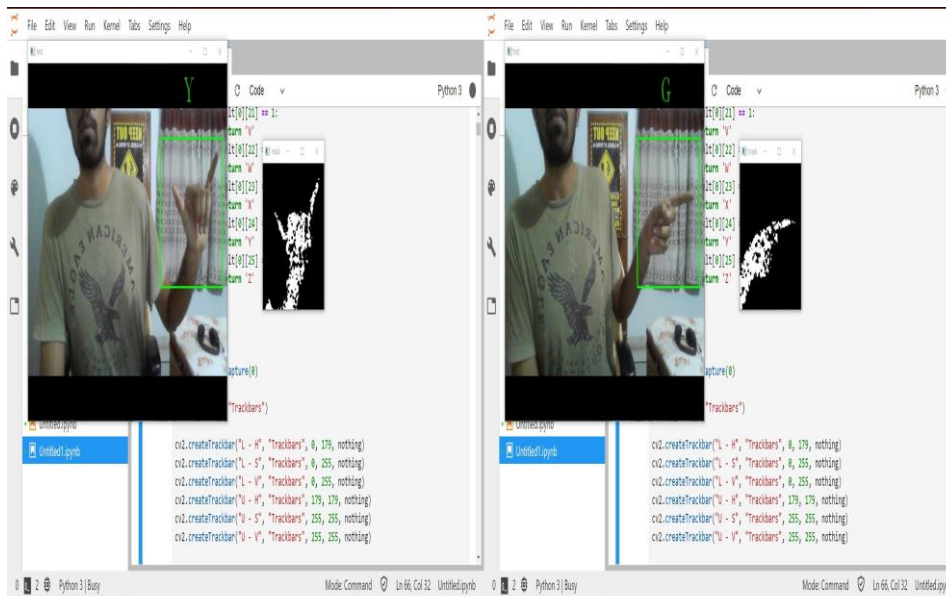


Figure 6. Alphabets 'Y' and 'G' as detected by the model

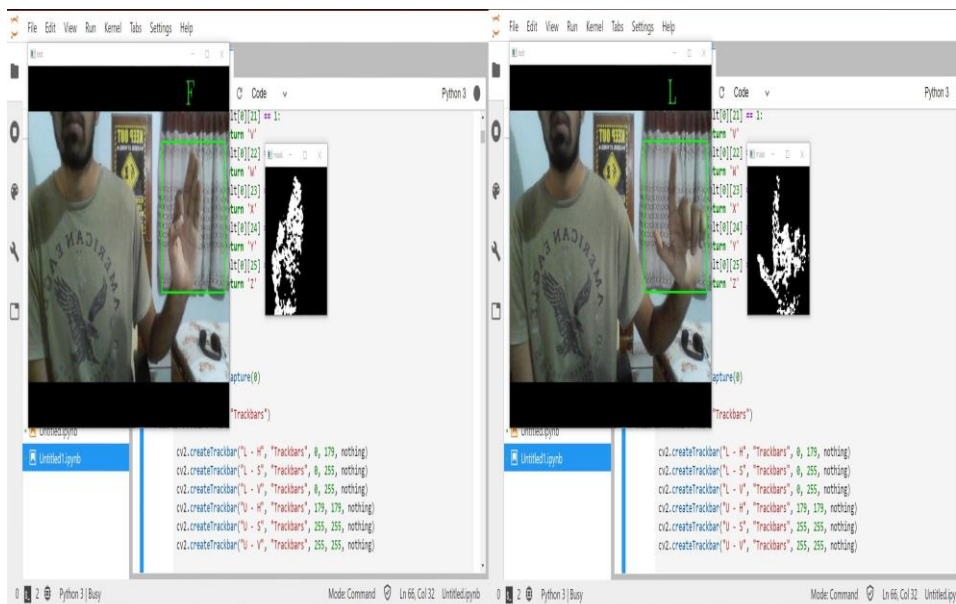


Figure 7. Alphabets 'F' and 'L' as detected by the model

As seen in the output screenshots shown in Figure 6 and Figure 7, I did not use a plain background but due to the process of background subtraction and selecting the appropriate HSV values in the trackbar, we were able to secure a mask that can easily detected and distinguished human skin from other elements existing in the camera vision. This feature is achieved due to selection of proper algorithm to model the background and foreground separator.

#### 4. Applications

- For ssecurity purposes, companies are using face recognition to secure their premises.
- Immigration checkpoints use facial recognition to enforce smarter border control.
- Vehicle security management companies can use facial recognition to give access to the engine and secure their vehicles.
- Ride-sharing companies can use facial recognition to ensure proper passengers are picked up by the proper drivers [19].
- IoT benefits from facial recognition allows enhanced security measures and automatic access control reception.
- Enforcement can use facial recognition technologies together as a part of AI-driven surveillance systems.
- Retailers can use facial recognition to customize offline offerings and to theoretically map online purchasing habits with their online ones [20].

#### 5. Conclusion

The dimensionality of representations is often exploited to move representations to a higher level by discovering the spatial or spatiotemporal regions of interest or selecting/extracting features that enhance the discrimination of similar looking expressions of different emotions. To these ends, most existing systems rely on generic dimensionality reduction techniques. The optimality of such techniques, however, is being questioned in the scope of affect recognition, and new trends address the importance of making use of domain knowledge explicitly when developing dimensionality reduction technique.

In this survey, we analyzed facial recognition systems by breaking them down into their fundamental components and we analyzed their potentials and limitations. In this section, we summarize the progress in the literature and highlight future directions. Advances in the field and the transition from controlled to naturalistic settings have been the focus of several survey papers. Zeng et al. focused on automatic affect recognition using visual and auditory modalities. Gunes and Schuller highlighted the continuity aspect for affect recognition both in terms of input and system output. Yet no survey has analyzed systems by isolating their fundamental components and discussing how each component addresses the above-mentioned challenges in facial affect recognition. Furthermore, some new trends and developments are not discussed in previous survey papers. Novel classification techniques that aim at capturing affect specific dynamics are proposed; validation protocols with evaluation metrics tailored for affect analysis are presented and affect recognition competitions are organized. Our in-depth analysis of these developments will expose open issues and useful practices and facilitate the design of real-world affect recognition systems.

The dimensionality of representations is often exploited to move representations to a higher level by discovering the spatial or spatiotemporal regions of interest, or selecting/extracting features that enhance the discrimination of similar looking expressions of different emotions. To these ends, most existing systems rely on generic dimensionality reduction techniques. The optimality of such techniques, however, is being questioned in the scope of affect recognition, and new trends address the importance of making use of domain knowledge explicitly when developing dimensionality reduction technique

#### References

1. Eason, G., Noble, B. and Sneddon, I.N. On certain integrals of Lipschitz-Hankel type involving products of Bessel functions. *Philosophical Transactions of the Royal Society of London. Series A, Mathematical and Physical Sciences*, 247(935), pp.529-551, 1955.
2. Maxwell, J.C. A treatise on electricity and magnetism (Vol. 1). *Oxford: Clarendon Press*, 1873.
3. Jacobs, I.S. Fine particles, thin films and exchange anisotropy. *Magnetism*, pp.271-350, 1963.
4. Ahammad, S.H., Rajesh, V. and Rahman, M.Z.U. Fast and accurate feature extraction-based segmentation framework for spinal cord injury severity classification. *IEEE Access*, 7, pp.46092-46103, 2019.
5. Yoroizu, T., Hirano, M., Oka, K. and Tagawa, Y. Electron spectroscopy studies on magneto-optical media and plastic substrate interface. *IEEE translation journal on magnetics in Japan*, 2(8), pp.740-741, 1987.
6. Sudhamani, M.J., Venkatesha, M.K. and Radhika, K.R. Fusion at decision level in multimodal biometric authentication system using Iris and Finger Vein with novel feature extraction. In *2014 Annual IEEE India Conference (INDICON)*, pp. 1-6, 2014.
7. Rao, G.A. and Kishore, P.V.V. Selfie sign language recognition with multiple features on adaboost multilabel multiclass

- classifier. *Journal of Engineering Science and Technology*, 13(8), pp.2352-2368, 2018.
8. Sunitha, R.A.V.I., Suman, M., Kishore, P.V.V. and Eepuri, K.K. Sign language recognition with multi feature fusion and ANN classifier. *Turkish Journal of Electrical Engineering and Computer Science*, 26(6), pp.2871-2885, 2018.
  9. Kumar, K.V.V., Kishore, P.V.V. and Anil Kumar, D. Indian classical dance classification with adaboost multiclass classifier on multifeature fusion. *Mathematical Problems in Engineering*, 2017.
  10. Bhavana, D., Kumar, K.K., Rajesh, V., Saketha, Y.S.S.S. and Bhargav, T., Deep Learning for Pixel-Level Image Fusion using CNN. *International Journal of Innovative Technology and Exploring Engineering*, 8, pp. 49-56, 2019
  11. Siva, D. and Bojja, P. MLC based Classification of Satellite Images for Damage Assessment Index in Disaster Management. *International Journal of Advanced Trends in Computer Science and Engineering*, 8(3), pp.10-13, 2019.
  12. Pardhasaradhi, P., Madhav, B.T.P., Sindhuja, G.L., Sreeram, K.S., Parvathi, M. and Lokesh, B. Image enhancement with contrast coefficients using wavelet based image fusion. *International Journal of Engineering and Technology (UAE)*, 7(2.8), pp.432-435, 2018.
  13. Suryanarayana, G. and Dhuli, R. Super-resolution image reconstruction using dual-mode complex diffusion-based shock filter and singular value decomposition. *Circuits, Systems, and Signal Processing*, 36(8), pp.3409-3425, 2017.
  14. Suresh, B., Subhani, S., Ghali, V.S. and Mulaveesala, R. Subsurface detail fusion for anomaly detection in non-stationary thermal wave imaging. *Insight-Non-Destructive Testing and Condition Monitoring*, 59(10), pp.553-558, 2017.
  15. Rahman, M.Z.U. and Reddy, B.M.K. Efficient SAR Image Segmentation Using Bias Field Estimation. *Journal of Scientific and Industrial Research (JSIR)*, 76(6), June 2017.
  16. Vallabhaneni, R.B. and Rajesh, V. On the performance characteristics of embedded techniques for medical image compression. *Journal of Scientific and Industrial Research (JSIR)*, 76(10), October 2017.
  17. Gattim, N.K., Rajesh, V., Partheepan, R., Karunakaran, S. and Reddy, K.N. Multimodal image fusion using Curvelet and Genetic Algorithm. *Journal of Scientific and Industrial Research (JSIR)*, 76(11), November 2017.
  18. Inthiyaz, S., Madhav, B.T.P. and Kishore, P.V.V. Flower image segmentation with PCA fused colored covariance and gabor texture features based level sets. *Ain Shams Engineering Journal*, 9(4), pp.3277-3291, 2018.
  19. Ahammad, S.K. and Rajesh, V. Image processing based segmentation techniques for spinal cord in MRI. *Indian Journal of Public Health Research and Development*, 9(6), 2018.
  20. Ahammad, S.H., Rajesh, V., Neetha, A., Sai Jeemitha, B. and Srikanth, A. Automatic segmentation of spinal cord diffusion MR images for disease location finding. *Indonesian Journal of Electrical Engineering and Computer Science (IJECS)*, 15(3), pp.1313-1321, 2019.