
Hand-Gesture Recognition

Nitya Mula
nmula3@uic.edu

Ramana Rao Akula
rakula3@uic.edu

Krishna Garg
kgarg8@uic.edu

1 Introduction

Gesture recognition is the process of understanding and interpreting meaningful movements of the hands, arms, face, or sometimes head. It can help users to control or interact with devices without physically touching them. Gesture recognition has been an interesting challenge in computer vision because of the difficulty involved in segmenting the foreground from the background effectively. It is obvious that there is a semantic gap between how humans look at a picture and how a computer does it. Yet, it has many widely used applications in computer vision and human-computer interaction. Some of these applications include gesture based audio/video control, robotic arm control, sign language system for impaired people as shown in Figure 1. In this project, we tried to develop an application that can recognize the hand gestures from a live video sequence / images.



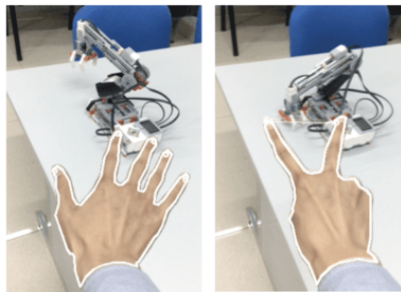
(a) Audio Video Control [Alam, 2020]



(b) Action Control [Heintz, 2018]



(c) Sign Language [Dias, 2019]



(d) Robotic Arm Control [Song *et al.*, 2016]

Figure 1: Applications of hand gesture recognition

2 Related work

The problem of hand gesture recognition has been tackled previously by various methods as shown in Figure 2. One of them includes using a special kind of glove called "the monochrome glove" [Ishiyama and Kurabayashi, 2016]. This glove has markers made of white tape on it that helps in identifying hand gestures. The glove has markers at the corners of the palm and base of fingers,

and white lines through the length of the finger. Based on the bending states and finger angles, the gestures are identified.

The second approach is to use orientation histograms [Lee and Chung, 1999]. In this method, from an image input or live sequence of video feature vectors are extracted. It does not use any gesture glove for this purpose. Firstly, the hand area is detected using any edge detection algorithms and separated from the input image or video. The feature vectors representing the hand are extracted using motion estimation. These are later analysed using orientation histogram scheme based on which the gesture is identified. This approach was used to recognize hand gesture for sign language. These vectors are fed to a Hidden Markov Model (HMM) which is a probabilistic graphical model used to predict a sequence of unknown gestures from a set of observed ones.

Another approach uses Fourier descriptors [Gamal *et al.*, 2013b] for gesture recognition. It has an edge over other methods mentioned above as it gives features that are invariant to rotation, translation and scaling. They represent the shape of a closed curve at varying form of detail. Also the time taken to process is less as only a less number of points are needed for entire image. Initially hand segmentation is performed using hue, HSV models. Later the Fourier descriptors are extracted and used as features for hand signs. Prediction is done using nearest neighbour or other classification techniques.

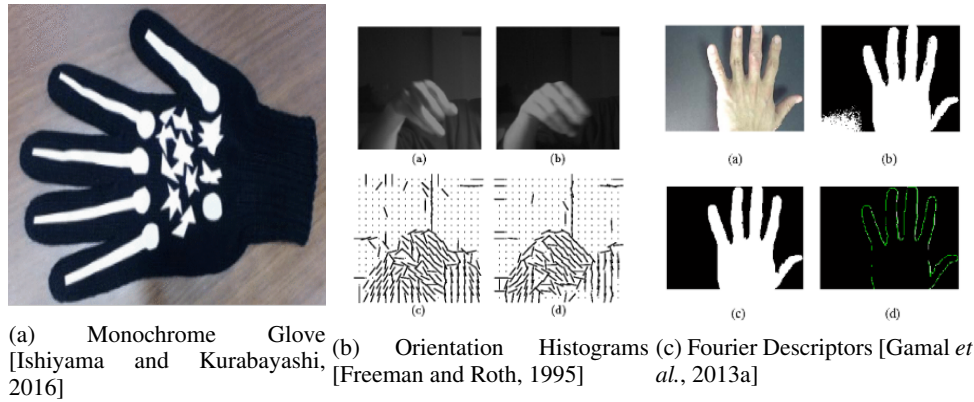


Figure 2: Existing approaches for hand gesture recognition

3 Problem setup

In this project, we tried to come up with a system that can identify gestures from images/live video sequence. Initially, we made use of contours in the image to identify hand region and then counted the number of fingers in the region using geometric techniques. Later, as we progressed further, we identified some limitations with this approach. This resulted in us shifting to a new approach which made use of neural networks.

The first approach requires the webcam to be working which captures frames and gives the input to the algorithm. In the second approach, a data set of images of various hand gestures is used to train different models and later predict the gestures from the unseen test data set.

4 Approach

To achieve hand gesture recognition, we initially used the method of finding hand contours and counting the fingers. However, there are some limitations associated with it. This led to shifting to another approach using deep learning.

4.1 Using hand contours

A contour can be defined as a curved line indicating the boundary for same intensity regions in an image. The contour detection can be handled by using `cv2.findContours()` from `opencv` library. Figure 3 shows how the contours look on an image.

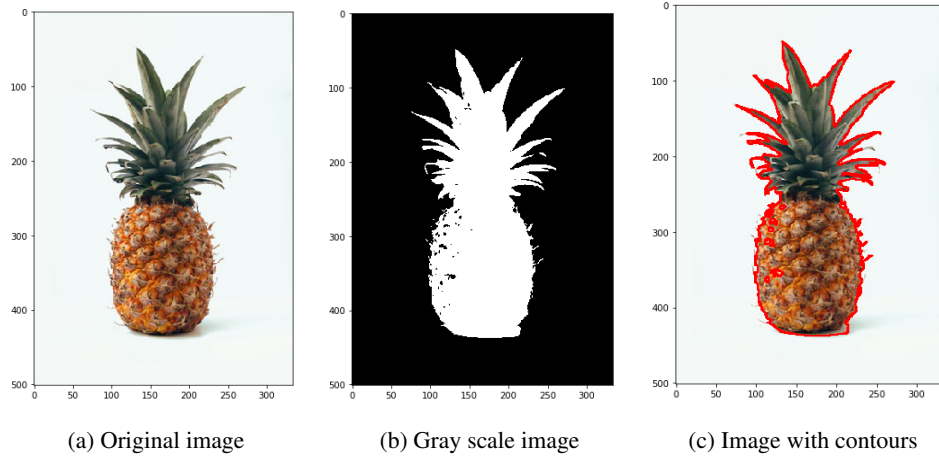


Figure 3: Images showing how the contours are identified and drawn[Jeong, 2019]

This approach is followed in two steps namely i) Background subtraction and thresholding ii) Contour extraction. The first step is the most crucial part of this approach because if the hand is segmented incorrectly, the entire algorithm would fail. As a part of the first step, to identify the hand region from the image / video sequence we use background subtraction techniques (frame differencing) and thresholding. The background of the image / video sequence is stored over a certain frame count. Now when a hand comes into the picture, the background image stored previously is subtracted from the current frame to get the foreground image. This way the presence of hand is identified and contours are drawn assuming that the contour with the largest area is our hand.

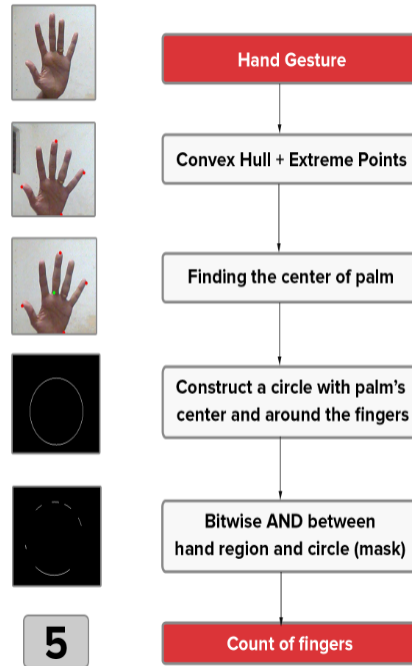


Figure 4: Flow chart for determining finger counts using hand contours [Ilango, 2017]

Figure 4 shows the flow chart for determining finger counts using hand contours. Using contour methods in openCV library, we found out the left and right, top and bottom extreme points. The centroid of the contour is calculated. With this point as the center and the radius as maximum of euclidean distance to extreme points times some threshold, a circle is constructed. Bitwise AND operation was performed between the thresholded hand image and the circular mask obtained above.

This gives the number of fingers in the frame by counting the number of transitions from black to white. This way, the finger counts of the hand gesture can be recognized.

However, this approach has a number of limitations. Firstly, this approach can detect up to five finger counts only and thus its scope cannot be expanded. It is too geometric and hand designed. One major drawback is with the assumption that the largest contour (maximum area) in the frame is our hand. If an object larger than the hand is brought inside the frame, then this algorithm fails. However, these drawbacks can be improved by changing the approach to the problem. This is the reason why we choose another approach for hand gesture recognition i.e., by using deep learning techniques.

4.2 Using deep learning models

Deep Learning is a sub-field of machine learning concerned with algorithms inspired by the structure and function of the brain and build using neural networks. A neural network takes in inputs images, processes them in the hidden layers by extracting features and using weights that are adjusted during training. Then the model makes prediction for the unseen data as seen in Figure 5. The weights are adjusted to find patterns in order to make predictions better.

In recent years deep learning has got an upper hand when it comes to computer vision problems. This improvement in the performance of the deep learning models is due to the recent advances in GPU design and architectures. They are well adjusted for training these types of models and are parallel in nature. We have trained models using several neural networks using a hand gestures data set and compared the performances.

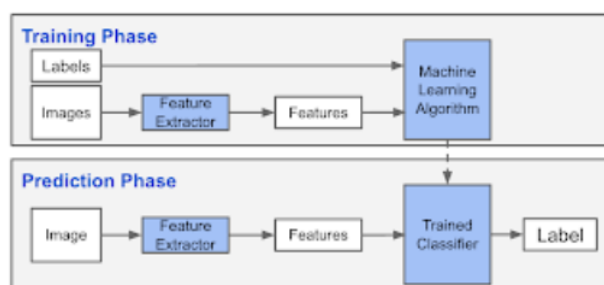


Figure 5: Different phases of using deep learning techniques [Moujahid, 2016]

Features are extracted from the input image data set and trained on various models both custom and pre-trained models. Predictions are made for unseen test images and the gesture is identified. This approach overcomes the limitations mentioned in the approach using hand contours as the scope of the problem can be increased by spiking the number of gestures. Also, its more sophisticated compared to the hand designed approach previously used.

5 Experimental setup

5.1 Using hand contours

The experimental setup used in this approach is as follows

opencv version - 3.4.2 imutils version - 0.4.6 numpy version - 1.18.1 sklearn version - 0.22.1

The experiment was carried out on a PC with no GPU and it did not give us any problems during the experiment as it was not a GPU intensive application and could be handled fairly by a CPU.

Table 1: Hyperparameters & Parameters details

Hyperparameters	Values
Batch Size	4
Epochs	20
Loss function	Categorical Cross-entropy
Optimizer	SGD
Learning Rate	0.01 [ResNet50, ResNet152], 0.0001[InceptionV3, VGG16]
Momentum	0.9
Data shuffling	Yes
Image Size	224*224*3

5.2 Using deep learning models

5.2.1 Dataset

Since the previous approach using contours had many limitations as discussed in 4.1, we decided to use deep learning models for the problem of gesture recognition. We retrieved a dataset available online for training and testing containing 10 different hand-gestures. The dataset originally was split up only in training and testing sets with 1000 and 50 images for each gesture in the training and testing sets, respectively. But for choosing better hyperparameters, we partitioned the training set into two datasets - training and validation sets. A validation data set is created by manually selecting some images from the training data. Finally, the data was shuffled before use.

5.2.2 Architecture

We used some famous convolutional neural networks such as ResNet50, ResNet152, InceptionV3 and VGG16. Additional layers namely Flatten, Relu activation and Softmax were added to the output from these architectures. The softmax classifier classified the input gesture into one of the 10 classes. We used Keras v2.3.1 for the implementation.

5.2.3 Hyperparameter tuning for training models

The hyperparameter tuning was done by running the model against different learning rates. We figured out the learning rates 0.01 was well-suited for the models ResNet50 & ResNet152. However, the other models VGG16 & InceptionV3 were giving poor accuracy with this learning rate. Finally, a much lower learning rate of 0.0001 gave very high performances even with VGG16 and InceptionV3. The training of the models was done using the GPU NVIDIA GeForce RTX 2080. Training time varied for different architectures. It was approximately 360 sec for each epoch for ResNet50, 890 sec for ResNet152, 110 sec for VGG16, 480 sec for InceptionV3. Refer table 1 for more details about the parameters and hyperparameters.

6 Results

6.1 Using hand contours

Once the python code is run, it opens up a window showing the live video capture on the webcam. After a particular frame count, the motion within the box is detected and foreground image is separated. Once this is done, the contours are obtained for the hand region as shown in the image below. Now the algorithm tries to get the finger count by counting the change in the transitions from black to white with the help of bit-wise AND operation. This finger count is displayed on the top left corner of the window. As the input gesture keeps changing, the algorithm keeps running and updates the finger count at each moment as seen in Figure 6.

However as previously mentioned when an object bigger than the size of the hand comes into the frame, then the algorithm fails. It can be seen in Figure 7 that when a bigger object comes, its count is predicted to be 1 instead of detecting the hand and giving a count of 4.

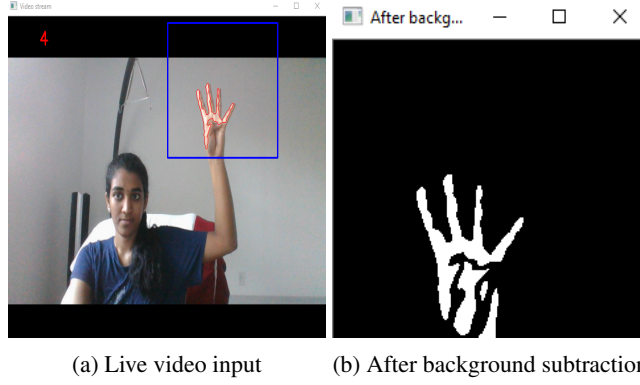


Figure 6: Counting the fingers and displaying the count on top left corner of the window

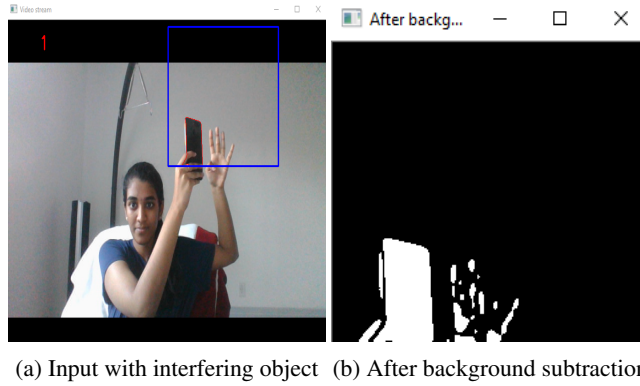


Figure 7: Wrong finger count due to an object interfering in the input

6.2 Using deep learning models

We evaluated the four different models ResNet50, ResNet152, VGG16, InceptionV3 by first training them for 20 epochs and then recording the final training accuracy and the accuracy on the testing set. The table 2 summarizes the different results. Also, we recorded the training and validation accuracy and the plots are depicted in figure 8. Some remarkable observations are as follows.

1. Plots in figure 8 depict that the losses keep on decreasing and the training & validation accuracies go on increasing with the number of epochs. This suggests that models are getting trained properly.
2. We can observe that the training and validation accuracies are almost 100% for all the models whereas the testing accuracy is in the increasing order of the following list: ResNet50 < ResNet152 < InceptionV3 < VGG16.
3. Since the training accuracy is 100%, it is quite possible that model is overfitting. We think that cross-validation could have helped us better tune our hyperparameters to get a better estimate of the training accuracy. Right now, we have a fixed validation set. Another way to get a good estimate of training accuracy could have been to use dropout layer during training time.
4. Another insight worth highlighting is the testing accuracy is also high. It is 97% and 95% for VGG16 and InceptionV3. So, it is quite possible that we are working with an easy dataset. If we manually inspect the images of the dataset, we can see that images are very similar to each other and therefore we are performing good on both training and testing datasets.

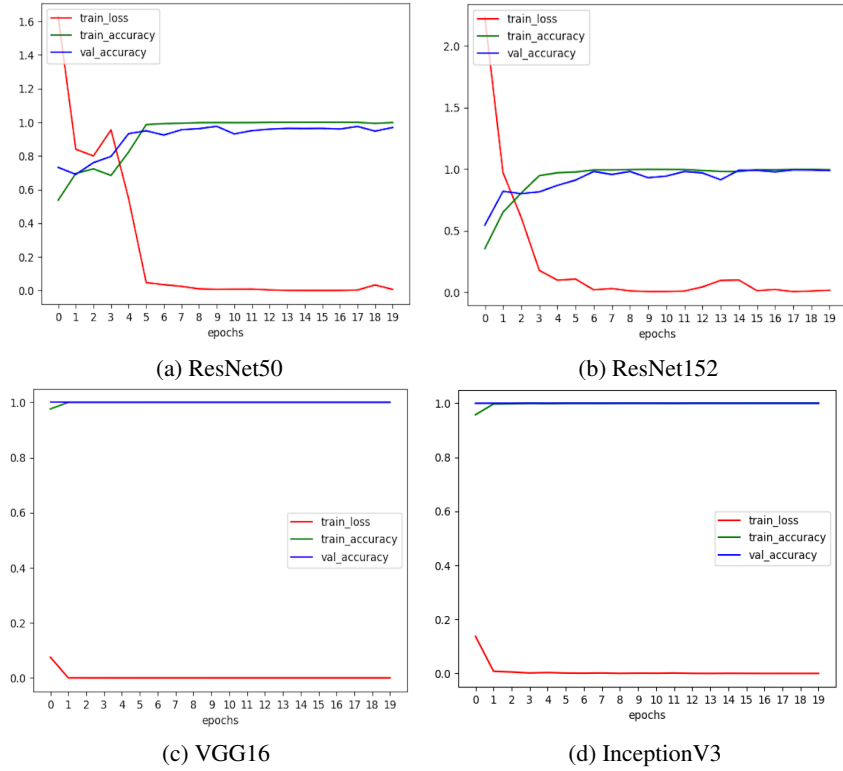


Figure 8: Plots for losses and accuracy vs epochs for 4 different CNN architectures (a) ResNet50 (b) ResNet152 (c) VGG16 (d) InceptionV3. Red line plots loss values recorded during training, green and blue lines plot training and validation accuracy, respectively, after every epoch.

Table 2: Results comparing different models

Models	Training Accuracy	Testing Accuracy
ResNet50	99.88%	85.04%
ResNet152	99.58%	87.21%
VGG16	100%	96.84%
InceptionV3	100%	94.85%

7 Conclusion

Hand gesture recognition has been an important challenge in Computer Vision. In this project, we tried to come up with a system that can recognize hand gestures from images / live video sequences. There are several solution to approach this problem of which we tried implementing two. From the first approach, hand gestures are identified using geometric techniques and it can be seen that the algorithm performs well. However, it is not sophisticated enough and has many limitations. In order to over come these limitations, we used deep learning models. We gathered a data set of hand gestures and tried training several models and later tested using unseen gestures. Out of the models we experimented on we observed that VGG16 performed the best and gives a very good accuracy on test data followed by InceptionV3.

References

- [Alam, 2020] Alam. Gesture Based Audio Video Control, 2020.
- [Dias, 2019] Rawini Dias. Sign Language Hand Gesture Recognition, 2019.
- [Freeman and Roth, 1995] William T. Freeman and Michal Roth. Orientation histograms for hand gesture recognition. 1995.

- [Gamal *et al.*, 2013a] Heba M. Gamal, Hatem M. Abdul-Kader, and Elsayed A. Sallam. Hand gesture recognition using fourier descriptors. *2013 8th International Conference on Computer Engineering & Systems (ICCES)*, pages 274–279, 2013.
- [Gamal *et al.*, 2013b] Heba M Gamal, HM Abdul-Kader, and Elsayed A Sallam. Hand gesture recognition using fourier descriptors. In *2013 8th International Conference on Computer Engineering & Systems (ICCES)*, pages 274–279. IEEE, 2013.
- [Heintz, 2018] Brenner Heintz. Action Control Hand Gesture Recognition, 2018.
- [Ilango, 2017] Gogul Ilango. Flow chart for determining finger counts using hand contours, 2017.
- [Ishiyama and Kurabayashi, 2016] Hidetoshi Ishiyama and Shuichi Kurabayashi. Monochrome glove: A robust real-time hand gesture recognition method by using a fabric glove with design of structured markers. In *2016 IEEE Virtual Reality (VR)*, pages 187–188. IEEE, 2016.
- [Jeong, 2019] Jiwon Jeong. Contour Detection, 2019.
- [Lee and Chung, 1999] Hyung-Ji Lee and Jae-Ho Chung. Hand gesture recognition using orientation histogram. In *Proceedings of IEEE. IEEE Region 10 Conference. TENCN 99.'Multimedia Technology for Asia-Pacific Information Infrastructure'(Cat. No. 99CH37030)*, volume 2, pages 1355–1358. IEEE, 1999.
- [Moujahid, 2016] Adil Moujahid. Different phases of using deep learning techniques, 2016.
- [Song *et al.*, 2016] Hongyong Song, Weijiang Feng, Naiyang Guan, Xuhui Huang, and Zhigang Luo. Towards robust ego-centric hand gesture analysis for robot control. pages 661–666, 08 2016.

A Contributions

Krishna Garg - Worked on multiple experiments using CNN approach. Tried out various models with different learning rates and batch sizes.

Nitya Mula - Worked on the approach using contours. Identified and put forth the data set to be used for neural networks. Contributed to some experiments using CNN approach.

Ramana Rao Akula- Worked on the approach using contours. Worked on skeleton code for experiments using Keras. Contributed to some experiments using CNN approach.