

Project Report

4th December 2023



INTRODUCTION

GestoVision is an innovative application designed for empowering sign language communication with the power of Machine Learning. The app accurately recognizes ASL (American Sign Language) hand gestures in real-time from a live video feed, and with a smart and neat overlay displays it over the hand in the video itself.

ASL allows communication with the hearing-impaired. Even though a generalized sign language helped with inclusion and communication, not everybody is proficient in this gesture based communication method or feels the need to understand it. The motivation behind this project and aim is to bridge this communication gap by creating a machine learning software that recognizes ASL, and in the process of doing so, ASL alphabet recognition was recognized as the crucial first step towards breaking down this barrier. GestoVision serves as a bridge, aimed towards facilitating communication in various fields such as education, sign-to-text conversion, healthcare, and more.

In the initial stages of the project, various solutions were explored to address the challenges of ASL gesture recognition. Multiple methodologies were considered, each with its distinct advantages and drawbacks. Among the solutions investigated were traditional computer vision approaches, deep learning models, and existing frameworks tailored for hand detection. In the quest for an optimal solution, the capabilities of cutting-edge technologies such as MediaPipe (a versatile library for hand detection) emerged as the backbone of the project.

The implemented software takes live video feed as input, employs hand detection using MediaPipe, and identifies the specific ASL alphabet. We used scikit-learn to train a machine learning model on a Kaggle dataset which is saved and used to detect the alphabets at high speed.

TOPIC DESCRIPTION

The project focuses on solving the challenges associated with ASL gesture recognition, as an elementary step towards bridging the communication gap with the hearing-impaired. Key elements in the topic description include:

Image Processing: The video feed is usually going to be from a webcam, which more often than not are of questionable quality and framerates. This raises a lot of concerns about how detailed the image is going to be, and gives rise to the possibility of blurs and the hands not being in focus. Again the feed might also have noise and in need of color correction. Some frames might also be overshadowed or have too high exposure. While finding a solution, all of these have to be kept in mind.

Hand Detection: Determining an effective method for detecting hands in a dynamic video environment can be challenging. Detecting a hand is a pretty big concern to begin with, not to mention there might be more than one hand in the video. The project would have to either accommodate multiple hands, or be specified to a single one.

Continuous Tracking: Even if detecting a hand on a still frame is possible, developing a mechanism for continuous tracking of the detected hand throughout the video feed after the first initial detection is a big concern. If the app is always detecting the hand from the entire frame in each frame it is going to be very computationally expensive. For efficiency and seamless hand detection in the final stage of the application, a method to track the hand after the initial detection needs to be implemented.

Data Acquisition: Obtaining relevant data for training, which takes into account the different possibilities of the background, lighting, skin tones etc. is necessary. Training a machine learning model has its advantages but it also has a big disadvantage, for it to accurately understand any new scenario, it has to be trained. So, when detecting the hand gestures, the model needs to be trained on various scenarios, like different lighting, gestures with people of different skin tones, different hand types, or even different hand gesture accents.

Data Representation: One of the most important parts of creating a machine learning model, if not the most important part is choosing the right data type, or preparing the data correctly. Choosing an appropriate representation for the acquired data(photos), proved to be challenging at first. If you train the model on an entire frame, the machine would have to read every pixel, every frame and it would be insanely heavy computationally and would possibly cause a lot of damage to the learning algorithm as well. Therefore, coming up with a better way to represent the data was essential.

SOLUTION DESCRIPTION

The solution involves the utilization of MediaPipe for hand recognition, Kaggle datasets supplemented by custom data for training, randomforest algorithm from scikit-learn library for the machine learning algorithm, and a combination of image processing techniques for optimal performance.

MediaPipe for Hand Recognition: Leveraging the capabilities of MediaPipe, a robust library for hand detection, the process of identifying hands in real-time was streamlined. After going through various methods of object detection and hand detection, MediaPipe was a game changer. MediaPipe accurately detects hand in multiple scenarios, even in different lighting and background. This library is also very computationally efficient which leads to our code being more streamlined and faster. MediaPipe also addresses the concern of continuous tracking in the video feed. It uses a very smart algorithm to follow the hand after an initial hand detection method, and if the hand is no longer detected in the video, the hand detection algorithm is initiated again followed by the tracking algorithm. The switching between a hand detection algorithm and a tracking algorithm significantly improves performance and makes it so that the app is not unnecessarily doing the search algorithm every frame.

Kaggle Dataset and Custom Data: Incorporated a Kaggle dataset along with custom images to ensure a diverse and comprehensive training set. A kaggle dataset was found with 3000 images for every letter with different light settings, backgrounds and skin tones to make the machine used to all the different situations. A few hundred images of the GestoVision creators for each alphabet were also used to ensure diversity.

MediaPipe Hand Landmarks to represent the data: MediaPipe library also has a method to create landmarks (points at each joint/junction of the hand) for the hands which can be used as a proper data representation of the hand positions. This ensures we are not using thousands of pixels as data to train the model and only have a few x,y coordinates, 21 points to be exact. The coordinates found on the training images are then strimmed to the lower left part of a graph plot, then when the model is trained, they are compared to the coordinates of the hand in the live video to accurately make a guess of the alphabet being shown.

Scikit-learn RandomForest algorithm for training the model: In the training phase, the power of machine learning was harnessed through the implementation of the RandomForest algorithm from the scikit-learn library. This algorithm, known for its robustness and versatility, played a pivotal role in training the model to recognize ASL alphabets. Leveraging the Kaggle dataset, we employed RandomForest to build a predictive model capable of efficiently categorizing hand gestures in real-time. The decision-making capabilities of the RandomForest algorithm, coupled with its ability to handle complex data patterns, contributed to the accuracy and speed essential for our ASL recognition framework. The model would not only provide the ASL alphabet that was recognized, but also a percentage which says how close the gesture is to the gesture used to train the data.

Image Processing: Implementing a sharpening kernel for enhancing hand features and blurring for noise reduction. These two image processing methods were what could fit into the app while balancing between an optimized performance without compromising the real-time nature of the application. A lot of other methods were experimented with but cut off due to performance issues. Besides, MediaPipe already deals with most pre-processing concerns.

EXPERIMENTS

Various machine learning methods were in consideration for the project. TensorFlow was initially planned to be used, but due to the complexity of its implementation, a much simpler library scikit-learn was employed. And for hand detection, various methods for object recognition were researched, such as– R-CNN Model to detect objects, however, towards the end, the mediaPipe library seemed to have an efficiency way better than any other methods.

During the initial stages of the project, various image processing methods were experimented with, such as– Sobel-Edge Detection for sharpening, Histogram Stretching and Gamma Correction. A lot of the photos generated during data gathering ended up being blurry. We decided to use Sobel for increasing the detail. Again, the webcam used during data gathering produced very poor results. Histogram Stretching was just one of the steps in our quest for image processing. Finally Gamma Correction technique was particularly valuable when dealing with diverse lighting conditions, ultimately contributing to a more robust and accurate machine learning model.

While various image processing methods were explored, a trade-off was established to maintain real-time performance. The use of a sharpening kernel and blurring for noise reduction was found to be efficient, considering that MediaPipe already addresses several image processing concerns.

REFERENCES

Nam liber tempor cum soluta nobis eleifend option congue nihil imperdiet doming id quod mazim placerat facer possim assum. Typi non habent claritatem insitam; est usus legentis in iis qui facit eorum claritatem. Investigationes demonstraverunt lectores legere me lius quod ii legunt saepius.

CONTRIBUTIONS

Trishir & Farhan: Equal Contribution

1. ASL Alphabet Dataset: <https://www.kaggle.com/dsv/29550>