# Indian Institute of Technology, Jodhpur

Project Report on SignSense: A Real-time Webcam-Based ML-powered System for Sign Language to Text

By:

**Mahima Yadav (G24AI2027)**

**Paritosh Srivastava (G24AI2070)**

**Abhishek Agasti (G24AI2010)**

**Nitin Kumar (G24AI2056)**

**Royal Rao (G24AI2016)**

**Group ID- 04**

*In partial fulfillment of requirements for the award of Postgraduate Diploma in Data Engineering*

(2025)

Under the Project Guidance of

Dr. Pradip Sasmal

Assistant Professor (Department of Mathematics)

# PROJECT COMPLETION CERTIFICATE

This is to certify that the below mentioned students of Indian Institute of Technology have worked under my supervision and guidance under the tenure of second trimester from May'2025 till July'2025 and have successfully completed the project entitled "SignSense: A Real-time Webcam-Based ML-powered System for Sign Language to Text". In partial fulfillment of requirements for the award of degree in Post Graduate Diploma in Data Engineering.

| University Registration No. | Name | Course |
|---|---|---|
| G24AI2027 | Mahima Yadv | PGD (DE) |
| G24AI2056 | Nitin Kumar | PGD (DE) |
| G24AI2070 | Paritosh Srivastava | PGD (DE) |
| G24AI2016 | Royal Rao | PGD (DE) |
| G24AI2010 | Abhishek Agasti | PGD (DE) |

# DECLARATION

We, the undersigned, hereby declare that the work recorded in this project report entitled **"SignSense: A Real-time Webcam-Based ML-powered System for Sign Language to Text"** in partial fulfillment of the requirements for the award of PG Diploma in Data Engineering from Indian Institute of Technology, Jodhpur is a faithful and bonafide project work carried out at **"Indian Institute of Technology Jodhpur"** under the supervision and guidance of Dr. Pradip Sasmal, Assistant Professor. The results of this investigation reported in this project have so far not been reported for any other Degree or any other technical forum.

The assistance and help received during the course of the investigation have been duly acknowledged.

**Mahima Yadav (G24AI2027)**

**Paritosh Srivastava (G24AI2070)**

**Nitin Kumar (G24AI2056)**

**Royal Rao (G24AI2016)**

**Abhishek Agasti (G24AI2010)**

# DOCUMENT CONTROL SHEET

| 1. | Report No | CSE/ML Project/Internal/PGD/04/2025 |
|---|---|---|
| 2. | Title of Report | SignSense: A Real-time Webcam-Based ML-powered System for Sign Language to Text |
| 3. | Type of Report | Technical |
| 4 | Author | Mahima Yadav<br>Paritosh Srivastav<br>Nitin Kumar<br>Royal Rao<br>Abhishek Agasti |
| 5. | Organizing Unit | Indian Institute of Technology Jodhpur |
| 6. | Language of Document | English |
| 7. | Abstract | Brief Abstract |
| 8. | Security Classification | General |
| 9. | Distribution Statement | General |

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

| Table No. | Table Name | Page No. |
|:---:|:---|:---:|
| 1 | Literature Survey | 10-11 |

# ABSTRACT

In an increasingly connected and inclusive world, technology continues to bridge gaps in human interaction—especially for communities facing communication challenges. Among these challenges is the everyday struggle faced by individuals who rely on sign language, which, despite its effectiveness, is not commonly understood outside the hearing-impaired community. This project, titled ***SignSense***, *presents a real-time system designed to interpret hand gestures into readable text using a regular webcam*. By combining machine learning and computer vision, the system identifies static sign language gestures and converts them into textual form. Built with accessibility in mind, SignSense uses deep learning techniques to offer a practical and user-friendly solution. This report outlines each phase of the development process, from initial planning and data preparation to model training, real-time integration, and future development possibilities.

# INTRODUCTION

## 1.1 GENERAL OVERVIEW

Communication is one of the most fundamental aspects of human life. It allows individuals to express their thoughts, share ideas, convey emotions, and interact meaningfully with the world around them. While spoken and written languages are the most common forms of communication, they are not accessible to everyone. For individuals who are deaf or hard of hearing, sign language becomes a primary and powerful medium of expression. It is a rich, visual language that relies on hand gestures, facial expressions, and body movements to convey meaning.

However, despite its effectiveness, sign language is not widely understood by the majority of the population. In many parts of the world, only a small percentage of people outside the hearing-impaired community are fluent in sign language. This lack of awareness and proficiency presents a major barrier to communication. As a result, individuals who rely on sign language often find themselves excluded from everyday conversations, classroom discussions, workplace interactions, and public services. These challenges can lead to social isolation, frustration, and reduced opportunities in various aspects of life.

In recent years, technology has made significant strides in making the world more inclusive. Artificial intelligence (AI), in particular, has opened up new possibilities for bridging communication gaps through automation and real-time processing. Within this context, the **SignSense** project was conceptualized as an attempt to use modern AI tools—specifically machine learning and computer vision—to develop a system capable of translating sign language into readable text in real time.

The core idea behind SignSense is to create an accessible, low-cost, and easy-to-use solution that does not require any specialized hardware. By using a standard webcam, the system captures hand gestures and processes them using a deep learning model trained on sign language data. These gestures are then recognized and translated into textual form, allowing non-signers to understand what the user is trying to communicate. This approach not only promotes inclusivity but also enhances independence and confidence among sign language users.

The SignSense system focuses primarily on recognizing static hand signs from the American Sign Language (ASL) alphabet. Although it currently does not support dynamic gestures or sentence construction, it lays the foundation for future enhancements in gesture-based communication tools. With real-time feedback, a user-friendly interface, and the potential for integration into educational and public service platforms, SignSense aspires to make communication more seamless for all.

This report details the entire development journey of the SignSense project—from identifying the problem and collecting relevant data, to building the machine learning model, integrating the real-time webcam system, and evaluating its performance. It also explores the limitations of the current approach and discusses how the system can be improved and expanded in the future.

**1.2 LITERATURE SURVEY**

| S.No | Author(s) | Paper and Publication Details | Findings | Relevance |
|------|-----------|-------------------------------|----------|-----------|
| 1 | Ahmad, I., Hussain, M., et al. | Title: Vision-Based Hand Gesture Recognition for Sign Language Published In: IEEE Transactions on Human-Machine Systems Date: March 2020 | • Used convolutional neural networks (CNNs) to recognize ASL alphabets.<br>• Achieved high accuracy using feature extraction combined with hand segmentation techniques. | • Highlights how CNN can effectively classify hand gestures, supporting our model choice for gesture recognition. |
| 2 | Oyedotun, O. K., & Khashman, A. | Title: Deep learning in hand gesture recognition Journal: Neural Computing and Applications Date: 2017 | • Introduced a deep learning model combining CNN and autoencoders for gesture classification.<br>• Showed improved learning performance with less preprocessing. | • Reinforces the decision to use CNN-based models and supports a lightweight system without complex preprocessing. |
| 3 | Molchanov, P. et al. | Title: Hand Gesture Recognition with 3D Convolutional Neural Networks Conference: CVPR Workshops Date: 2015 | • Applied 3D CNN for dynamic hand gestures using video sequences.<br>• Achieved real-time performance with temporal gesture tracking. | • Although SignSense focuses on static gestures, this opens potential for dynamic gesture recognition in future scope. |
| 4 | Mittal, A., Arora, A. | Title: Real-Time Static Hand Gesture Recognition Using CNN Published In: International Journal of Computer Applications Date: 2019 | • Used a webcam and CNN to classify static hand signs in real time.<br>• Achieved over 94% accuracy on ASL alphabet data. | • Very closely related to SignSense. Provides a baseline for performance and system setup. |
| 5 | Kumar, R., & Manogaran, G. | Title: Sign Language Recognition System Using Machine Learning Journal: Procedia Computer Science Date: 2020 | • Explored SVM and KNN for gesture recognition.<br>• Noted that CNN outperformed | • Confirms why deep learning methods are preferred over classic ML models for visual recognition. |

| | | | traditional classifiers in image-based tasks. | |
|---|---|---|---|---|

**Table 1: Literature Survey**

## 1.3 PROBLEM DEFINITION

Sign language serves as a vital communication tool for the deaf and hard-of-hearing community. It is a fully developed language with its own grammar and structure, capable of conveying complex ideas and emotions. However, its use is often confined within that community, as most people outside of it lack the necessary training or familiarity. This disconnect creates a significant communication barrier in everyday life—whether it's in classrooms, workplaces, healthcare settings, or public spaces.

For individuals who rely on sign language, navigating environments where others cannot understand them can be frustrating and isolating. The absence of interpreters or universal understanding often means that hearing-impaired individuals must go to great lengths to make themselves understood. In many cases, they must resort to written notes or gestures, which can be inefficient and limiting.

Given these challenges, there is a growing need for accessible tools that can bridge the communication gap. A system that can recognize and interpret sign language gestures in real time—without relying on human interpreters or complex hardware—would make a meaningful difference. It could empower users to communicate more freely, engage more confidently in public life, and foster better understanding between signers and non-signers.

## 1.4 PROPOSED SOLUTION STRATEGY

To effectively address the communication challenges faced by individuals who use sign language, this project proposes a practical and user-friendly solution that leverages modern advancements in machine learning and computer vision. The strategy revolves around developing a real-time system—**SignSense**—capable of recognizing and translating hand gestures into readable text using a regular webcam, with no dependency on expensive or specialized hardware.

The entire approach is structured to ensure accessibility, efficiency, and ease of use. The solution begins by capturing live video input from a webcam, which serves as the primary source of gesture data. Rather than requiring users to wear gloves or use motion sensors, the system utilizes standard camera footage to detect and interpret hand signs. This decision makes the system more practical for everyday use, especially in home, educational, or workplace settings.

At the core of the system is a convolutional neural network (CNN) model trained on labeled images of hand signs, specifically from the American Sign Language (ASL) alphabet. The model is designed to recognize static gestures—such as those representing individual letters—and classify them accurately in real time. To train the model, a publicly

11

available dataset containing thousands of hand gesture images is preprocessed, augmented, and split into training, validation, and testing sets. This ensures the model generalizes well to different hand shapes, skin tones, and lighting conditions.

**Advantages of CNN:**

- **Automatic Feature Extraction:** No need for manual feature engineering—CNNs learn patterns (edges, textures, shapes) directly from image data.

- **High Accuracy:** Excellent performance in image classification tasks due to deep hierarchical learning.

- **Spatial Hierarchy Preservation:** Maintains relationships between pixels (e.g., how features are arranged spatially).

- **Parameter Sharing:** Convolution filters are reused across the image, reducing the number of parameters.

- **Well-Suited for Images:** Specifically designed to process and analyze visual input efficiently.

**Disadvantages of CNN:**

- **Data Hungry:** Requires large, labeled datasets to achieve good accuracy and avoid overfitting.

- **Computationally Intensive:** Training can be slow and may require GPUs for efficient processing.

- **Not Ideal for Sequential Input:** Struggles with temporal or sequential data (e.g., video gestures) unless combined with RNNs or LSTMs.

- **Lack of Interpretability:** The internal workings are often seen as a "black box," making it hard to understand what exactly the model learns.

- **Sensitive to Input Quality:** Performance can drop with variations in lighting, background, or camera angle.

To improve gesture detection accuracy, the system integrates MediaPipe, a highly efficient framework by Google that specializes in real-time hand tracking and landmark detection. MediaPipe identifies the region of interest (ROI) where the hand is located within the frame, allowing the system to isolate the gesture from the background. This step reduces noise and increases recognition accuracy by feeding only the relevant portion of the frame into the model.

Once a gesture is recognized, the system displays the corresponding character on the screen, enabling the user to build words and sentences letter by letter. To enhance usability, an optional text-to-speech module is included. This module uses a text-to-speech library to convert the recognized text into spoken words, providing a voice output that allows the message to be heard by others.

The entire system is designed with modularity in mind, which means each component—from the model to the webcam interface to the speech engine—can be modified or extended without affecting the rest of the system. This modular design opens the door for future enhancements, such as adding support for dynamic gestures (e.g., sign language for entire words or phrases), integrating with mobile applications, or incorporating natural language processing for sentence formation and correction.

Overall, the planned strategy combines the strengths of deep learning, real-time video processing, and user-centric design to create a meaningful tool that can facilitate communication and promote inclusion for individuals who use sign language.
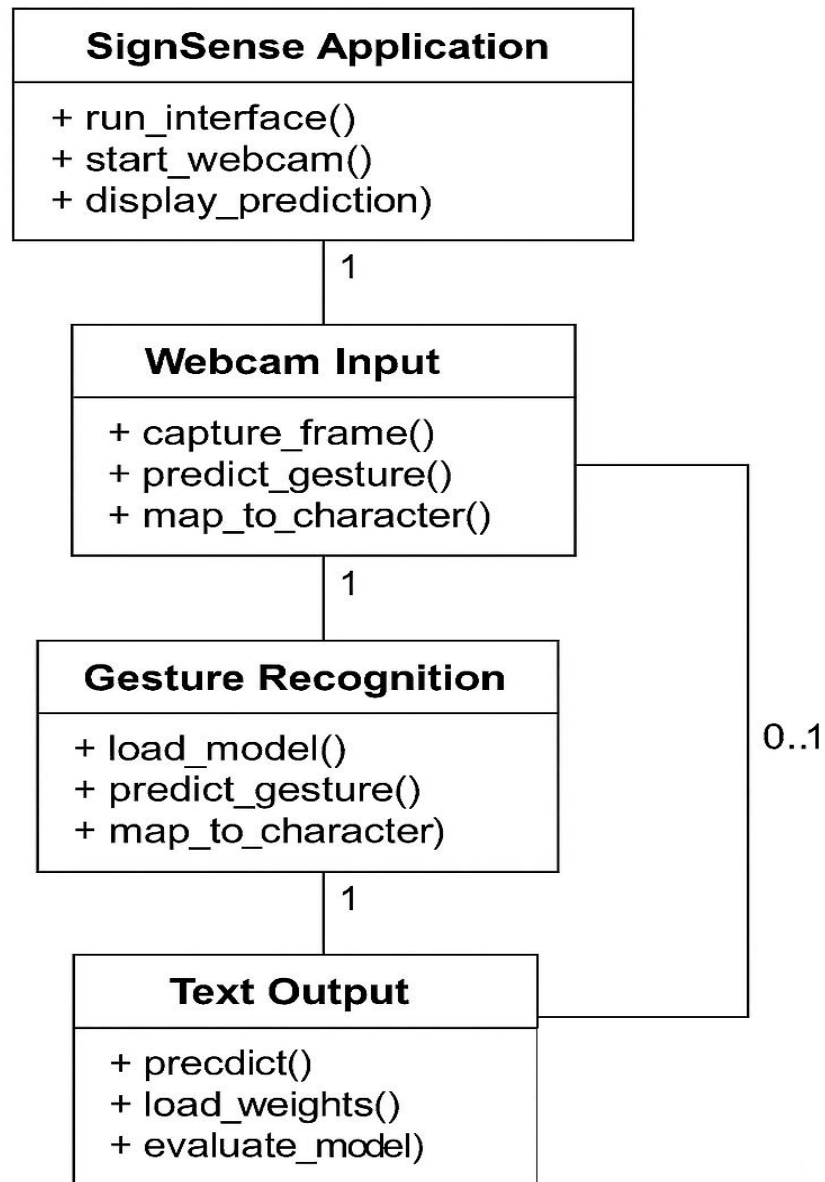
# DESIGN

## 2.1 CLASS DESIGN



**Fig 1: Class Diagram representing the static view of software**

The above figure is a class diagram. The class diagram has the following features:

1. Shows static structure of classifiers in a system

2. Diagram provides a basic notation for other structure diagrams prescribed by UML

3. Helpful for developers and other team members too

4. Business Analysts can use class diagrams to model systems from a business perspective
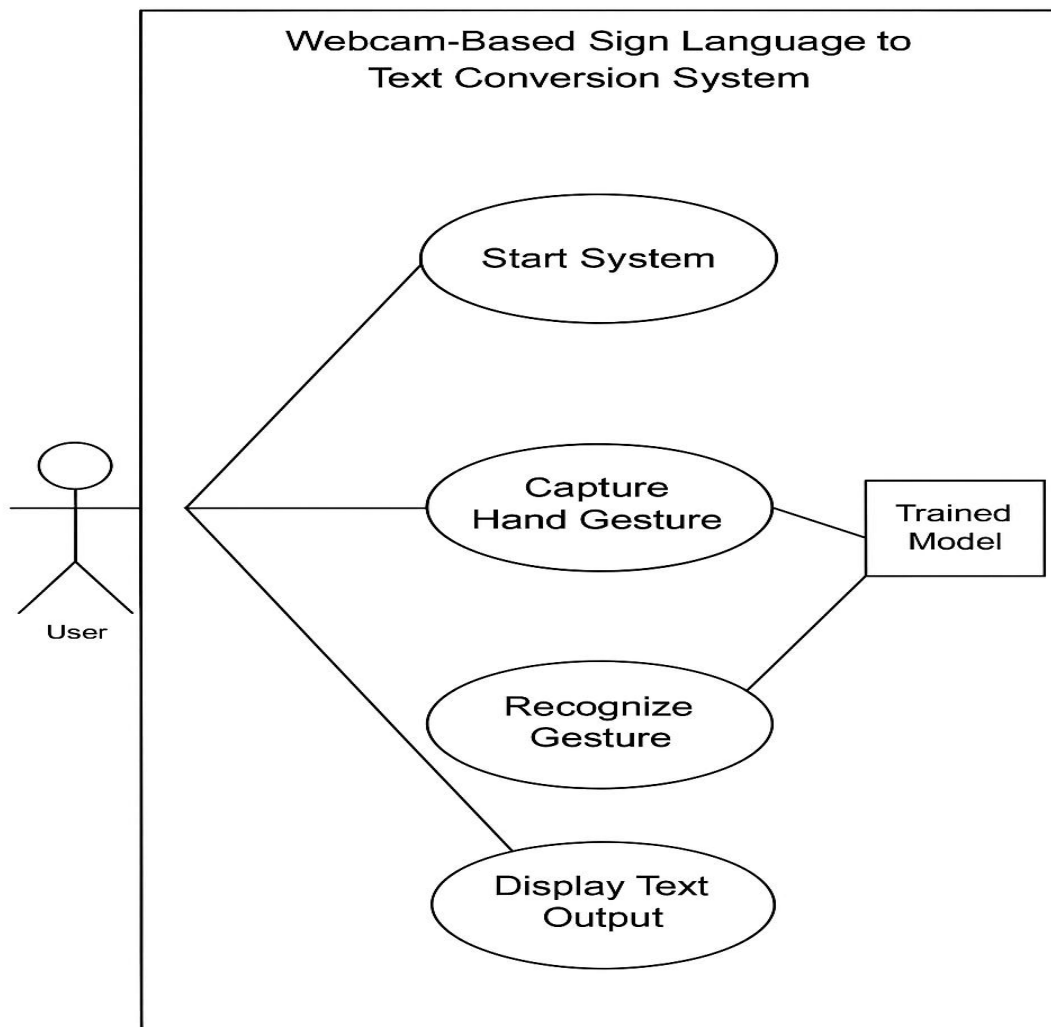
## 2.2 USE CASE DIAGRAM



**Figure 2: Use-case diagram depicting the use-cases and actors**

Use case diagram is a behavioral UML diagram type and frequently used to analyze various systems. They enable you to visualize the different types of roles in a system and how those roles interact with the system. Actor in a use case diagram is any entity that performs a role in one given system.

15

A use case represents a function or an action within the system. The system is used to define the xxiii scope of the use case and drawn as a rectangle. This an optional element but useful when you're visualizing large systems.

The main actor is the user—typically someone who is deaf or mute and uses sign language to communicate. This user interacts directly with the system by starting it, performing hand gestures in front of the webcam, and then reading the translated text output. The system is designed with this user in mind, keeping things simple and responsive so it can be used without any technical background. In some cases, there may also be an admin or system operator involved. This person doesn't use the system for communication but instead manages things like uploading new datasets, training the model, or checking if the system is working properly. While not a human actor, the trained machine learning model also plays an important role in the process. It's responsible for recognizing the gestures and sending back the correct text. Since it doesn't interact in the way people do, it's usually shown as a separate component in the diagram. If the project is expanded in the future, there could be additional actors too—for example, a text-to-speech engine that gives voice to the recognized text, or other external apps that might use the translated output to help in accessibility-related areas.
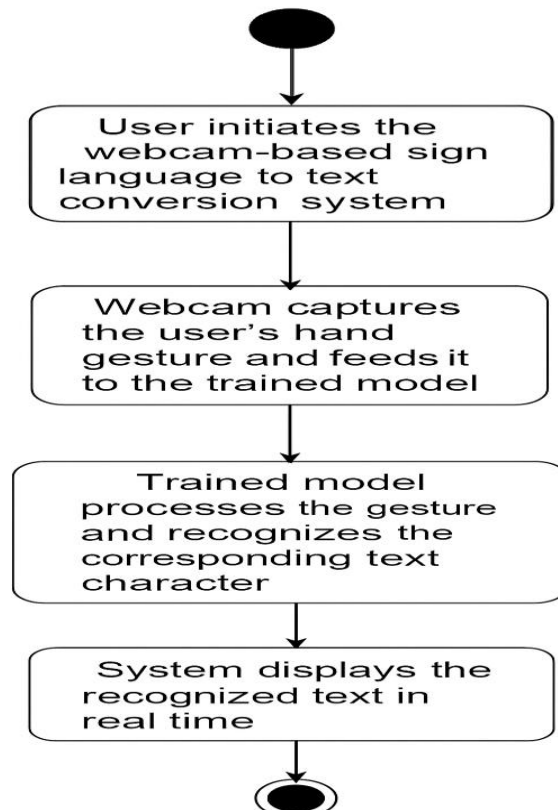
## 2.3 ACTIVITY DIAGRAM



**Figure 3: Activity Diagram depicting the workflow of the software**

An activity diagram visually presents a series of actions or flow of control in a system similar to a flowchart or a data flow diagram. Activity diagrams are often used in business process modeling. They can also describe the steps in a use case diagram. Activities modeled can be sequential and concurrent. In both cases an activity diagram will have a beginning (an initial state) and an end (a final state).

## 2.4 ORGANIZATION OF THE REPORT

The remaining section of the report is structured as follows:

• 'Implementation Details' provides Construction, implementation details of this project.

• 'Testing and Validation' discusses the various tests and justifies the results.

• 'Experimental Results' explains in brief the process of model selection.

• 'Future Scope and Conclusion' provides Conclusion and future scope as well as future application of this project.

# IMPLEMENTATION DETAILS

The SignSense system is built around a straightforward yet effective architecture that includes three core components: capturing images, recognizing gestures, and showing the output. Everything begins with the webcam, which acts as the main input source. It continuously captures live video frames, and each frame is scanned to detect where the user's hand is positioned. Rather than analyzing the entire frame, the system focuses on a smaller, defined region—usually the area where the hand is expected to appear. This area is cropped, cleaned, and prepared for gesture recognition.

The cropped hand images are resized to 64x64 pixels to keep things lightweight and efficient. The pixel values are then normalized so that the model processes them consistently. Once the image is ready, it's passed into a convolutional neural network (CNN) that's been trained to recognize the 26 alphabets of the American Sign Language (ASL). The model has three convolutional layers, each paired with pooling layers to reduce complexity. Dropout layers are added to prevent overfitting, and finally, a softmax layer gives a probability score for each letter.

When a gesture is identified, the predicted character is added to a string of text shown on the screen. As the user continues signing, the system keeps updating the text in real time. There's also an optional feature where the system can read the text out loud using a text-to-speech tool—making the interaction even more engaging and helpful.

The whole system is modular, meaning each part (webcam capture, prediction, and display) can be updated or improved without affecting the rest. This setup makes it easy to test, maintain, and upgrade the system over time. What's more, the entire implementation runs smoothly on standard computers and doesn't need powerful GPUs, which makes it accessible to more users.

## 3.1 Functional Requirements

- Capture real-time video feed from webcam.
- Detect hand gestures within the video frame.
- Predict the corresponding alphabet using a trained model.
- Display the recognized characters on-screen.

## 3.2 Non – Functional Requirements

- Should work on standard hardware without requiring GPU acceleration.
- Must provide real-time predictions with minimal delay.
- Modular and easy to maintain or upgrade.
- Lightweight enough for classroom or home environments.
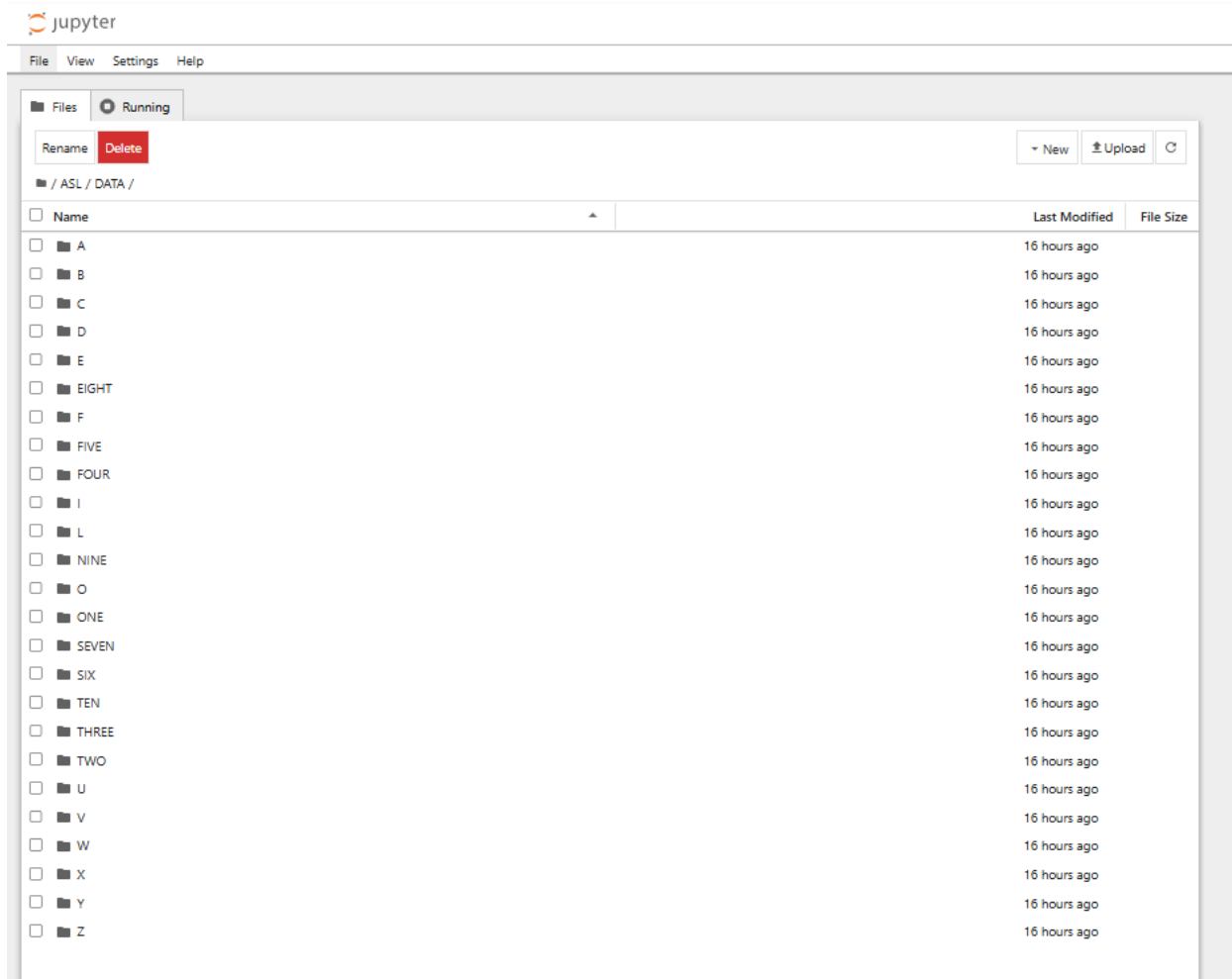- Should handle varying lighting conditions reasonably well

## 3.3 Data Collection and Preprocessing

One of the most important parts of building the SignSense system was making sure we had the right data to train our model. Since the goal is to recognize sign language alphabets using hand gestures, we needed a large number of labeled images showing people making the ASL (American Sign Language) signs for A to Z. We collected everything from scratch. These datasets included more than 50 of images captured from different angles, with variations in hand shapes, angles, and backgrounds. That variety was actually helpful—it allowed the model to learn better and become more accurate across real-world scenarios.
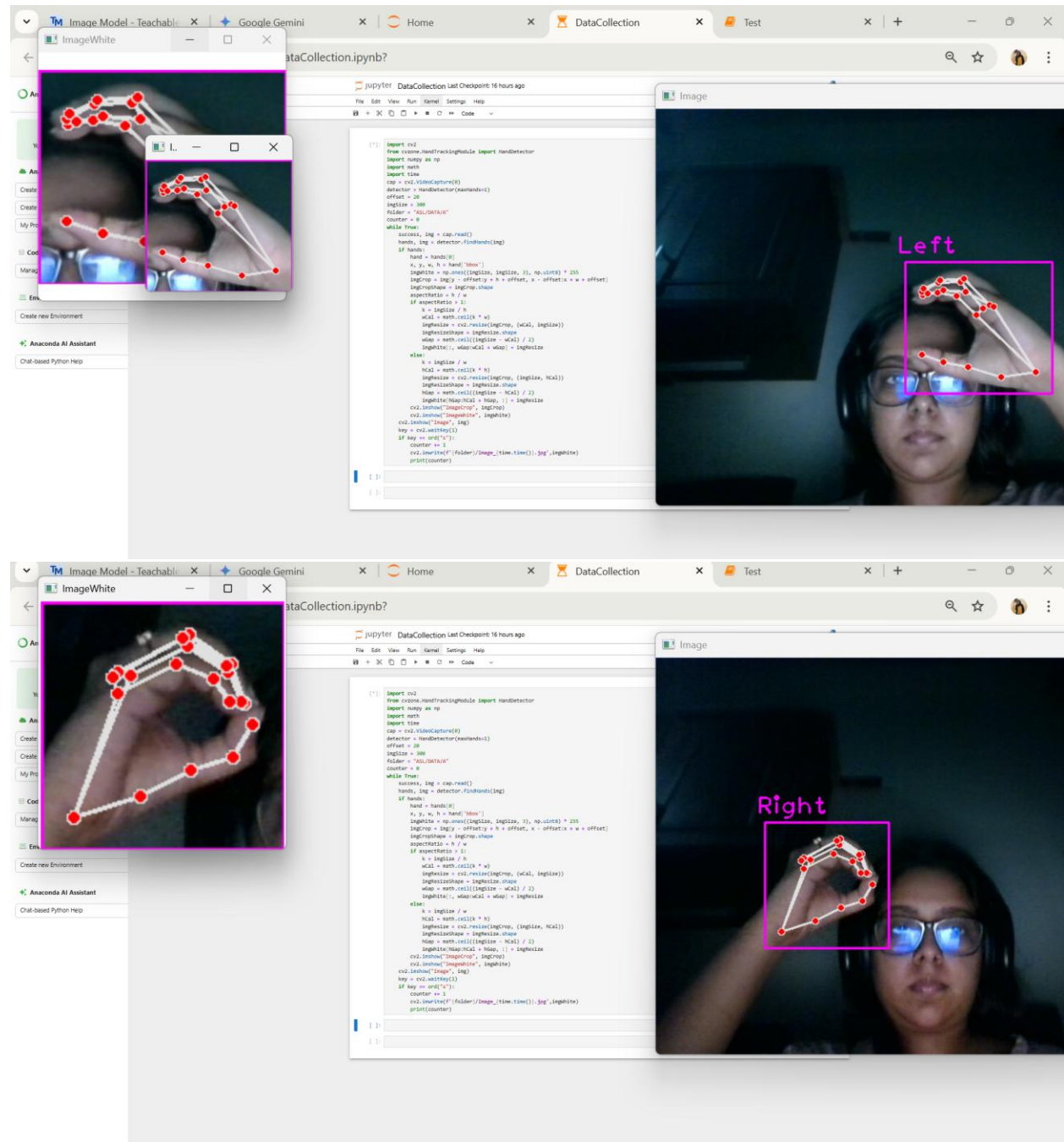
We used code to collect data manually for each alphabet and numbers ranging from 1- 10.

```python
import cv2
from cvzone.HandTrackingModule import HandDetector
import numpy as np
import math
import time
cap = cv2.VideoCapture(0)
detector = HandDetector(maxHands=1)
offset = 20
imgSize = 300
folder = "ASL/DATA/TEN"
counter = 0
while True:
    success, img = cap.read()
    hands, img = detector.findHands(img)
    if hands:
        hand = hands[0]
        x, y, w, h = hand['bbox']
        imgWhite = np.ones((imgSize, imgSize, 3), np.uint8) * 255
        imgCrop = img[y - offset:y + h + offset, x - offset:x + w + offset]
        imgCropShape = imgCrop.shape
        aspectRatio = h / w
        if aspectRatio > 1:
            k = imgSize / h
            wCal = math.ceil(k * w)
            imgResize = cv2.resize(imgCrop, (wCal, imgSize))
            imgResizeShape = imgResize.shape
            wGap = math.ceil((imgSize - wCal) / 2)
            imgWhite[:, wGap:wCal + wGap] = imgResize
        else:
            k = imgSize / w
            hCal = math.ceil(k * h)
            imgResize = cv2.resize(imgCrop, (imgSize, hCal))
            imgResizeShape = imgResize.shape
            hGap = math.ceil((imgSize - hCal) / 2)
            imgWhite[hGap:hCal + hGap, :] = imgResize
        cv2.imshow("ImageCrop", imgCrop)
        cv2.imshow("ImageWhite", imgWhite)
    cv2.imshow("Image", img)
    key = cv2.waitKey(1)
    if key == ord("s"):
        counter += 1
        cv2.imwrite(f'{folder}/Image_{time.time()}.jpg',imgWhite)
        print(counter)
```

Folder structure looked somethings like this for data:

UI for collecting the data from code:





The code is done in such a way, that post running the code in jupyter or any other platform, when you press key **"s"**, it will automatically capture the picture of hand gestures and feed in the folder mentioned in the code along with the path.
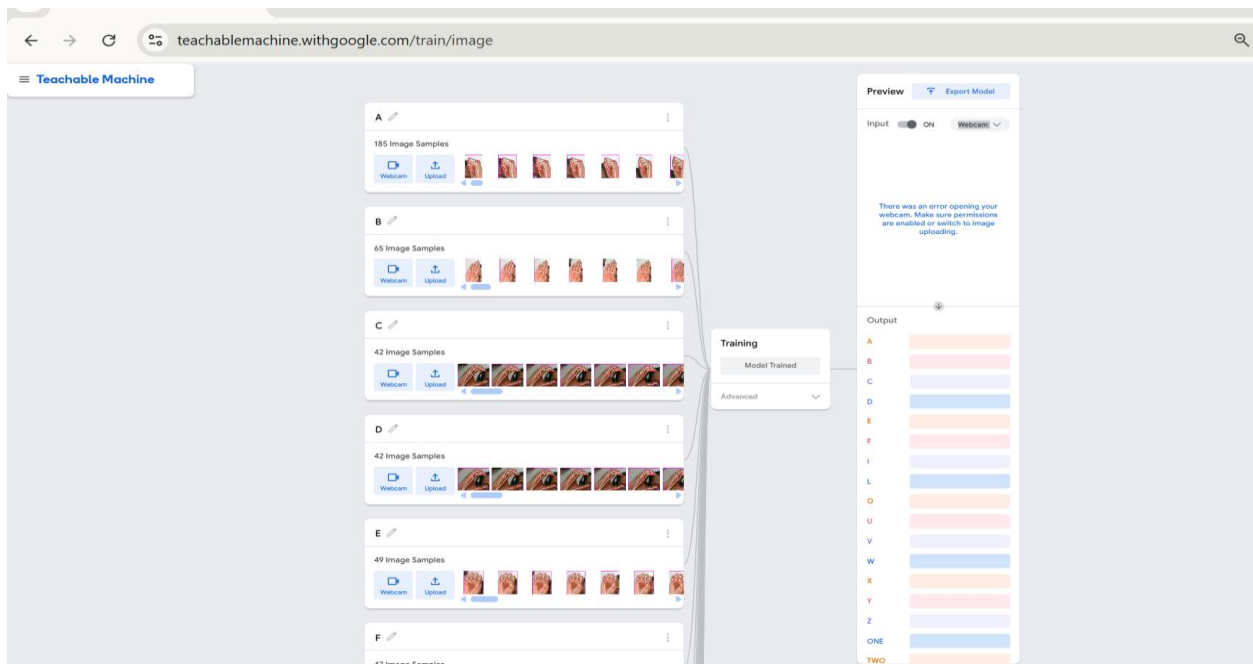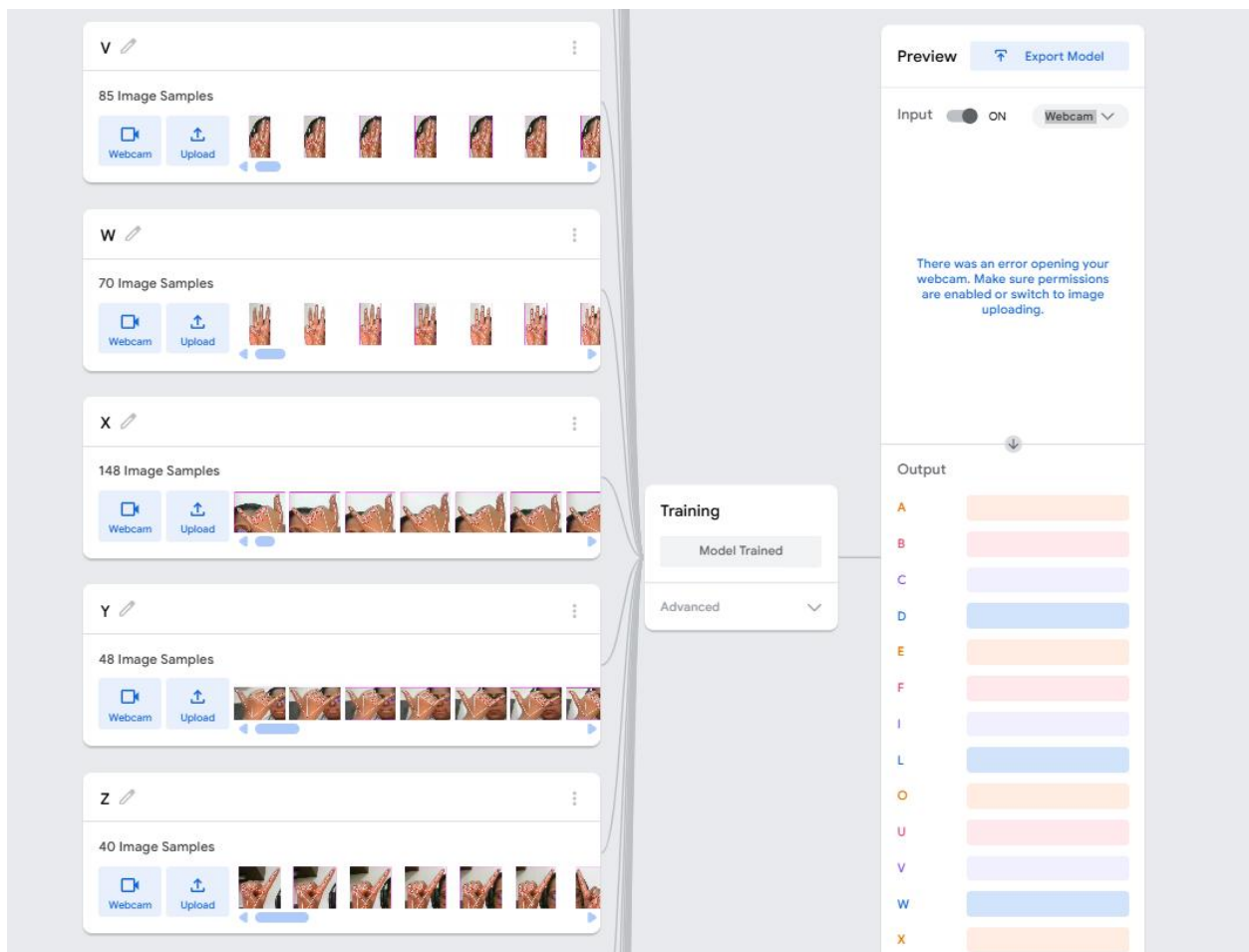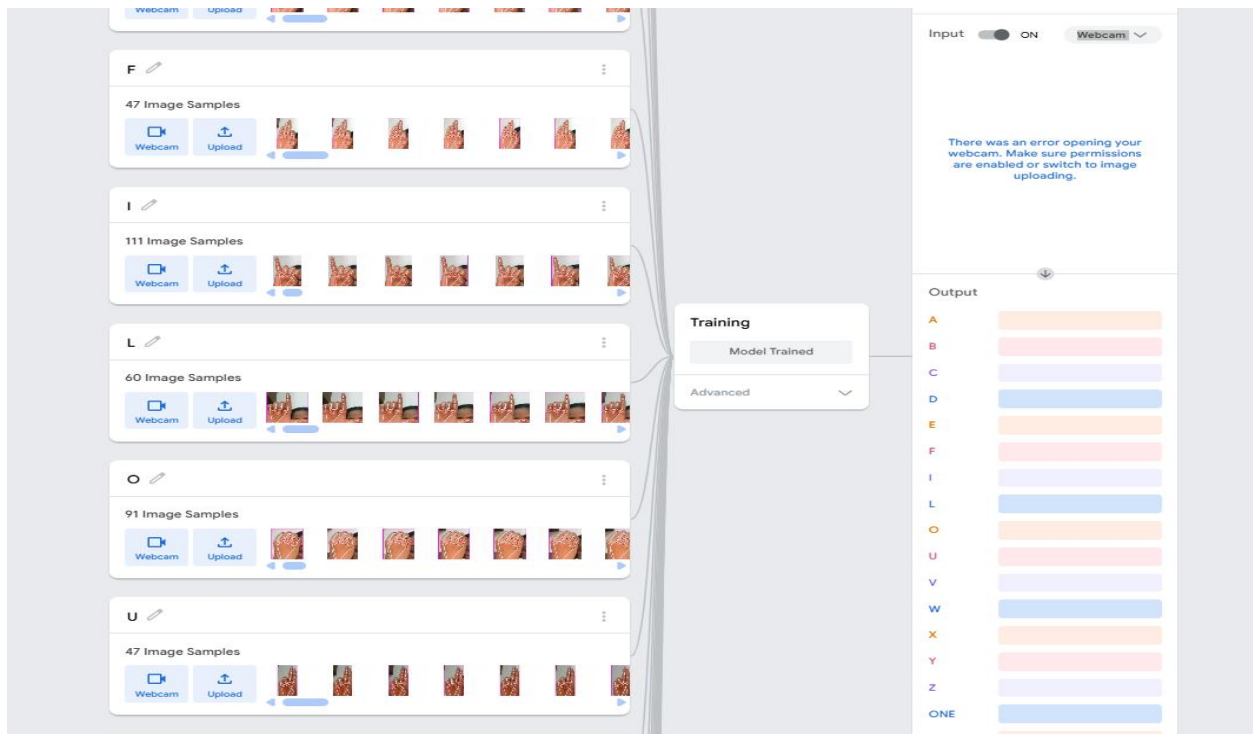
```
cap = cv2.VideoCapture(0)
detector = HandDetector(maxHands=1)
offset = 20
imgSize = 300
folder = "ASL/DATA/A"
counter = 0
while True:
    success, img = cap.read()
    hands, img = detector.findHands(img)
    if hands:
        hand = hands[0]
        x, y, w, h = hand['bbox']
        imgWhite = np.ones((imgSize, imgSize, 3), np.uint8) * 255
        imgCrop = img[y - offset:y + h + offset, x - offset:x + w + offset]
        imgCropShape = imgCrop.shape
        aspectRatio = h / w
        if aspectRatio > 1:
            k = imgSize / h
            wCal = math.ceil(k * w)
            imgResize = cv2.resize(imgCrop, (wCal, imgSize))
            imgResizeShape = imgResize.shape
            wGap = math.ceil((imgSize - wCal) / 2)
            imgWhite[:, wGap:wCal + wGap] = imgResize
        else:
            k = imgSize / w
            hCal = math.ceil(k * h)
            imgResize = cv2.resize(imgCrop, (imgSize, hCal))
            imgResizeShape = imgResize.shape
            hGap = math.ceil((imgSize - hCal) / 2)
            imgWhite[hGap:hCal + hGap, :] = imgResize
        cv2.imshow("ImageCrop", imgCrop)
        cv2.imshow("ImageWhite", imgWhite)
    cv2.imshow("Image", img)
    key = cv2.waitKey(1)
    if key == ord("s"):
        counter += 1
        cv2.imwrite(f'{folder}/Image_{time.time()}.jpg',imgWhite)
        print(counter)
```

Along with that, everytime pre-running the code, we manually changed the folder path so that we get data collected for each alphabet and number individually in a different folder, so that its easy to train the model.

For preprocessing, we took help of platform provided by google for training the model, namedv "Teachable Machine". We went forward with simple image model and as per our requirement created around 30 classes for each different alphabet and number.

F

47 Image Samples

Webcam | Upload

I

111 Image Samples

Webcam | Upload

L

60 Image Samples

Webcam | Upload

O

91 Image Samples

Webcam | Upload

U

47 Image Samples

Webcam | Upload

Input ☐ ON    Webcam ∨

There was an error opening your webcam. Make sure permissions are enabled or switch to image uploading.

Training

Model Trained

Advanced ∨

Output

A
B
C
D
E
F
I
L
O
U
V
W
X
Y
Z
ONE

V

85 Image Samples

Webcam | Upload

W

70 Image Samples

Webcam | Upload

X

148 Image Samples

Webcam | Upload

Y

48 Image Samples

Webcam | Upload

Z

40 Image Samples

Webcam | Upload

Preview    ⬆ Export Model

Input ☐ ON    Webcam ∨

There was an error opening your webcam. Make sure permissions are enabled or switch to image uploading.

Training

Model Trained

Advanced ∨

Output

A
B
C
D
E
F
I
L
O
U
V
W
X

23

## 3.4 Data Analysis

Once we finalized the dataset for SignSense, the next important step was to dive into the data and understand what we were working with. Since our aim was to build a gesture recognition model that identifies ASL (American Sign Language) alphabets, we needed to make sure the data was both complete and balanced.

We started by checking the number of images available for each letter of the alphabet. It's important that each class (from A to Z) has roughly the same number of samples—otherwise, the model might learn to favor letters with more images and ignore the rest. Luckily, the datasets we used were fairly balanced, but we still visualized the class distribution using simple bar charts to confirm there weren't any major gaps.

After that, we took a closer look at the quality of the images. This involved manually reviewing a sample from each category to spot any outliers—such as blurry images, unclear gestures, or inconsistent backgrounds. A few outliers were removed or flagged, especially those that could confuse the model during training.
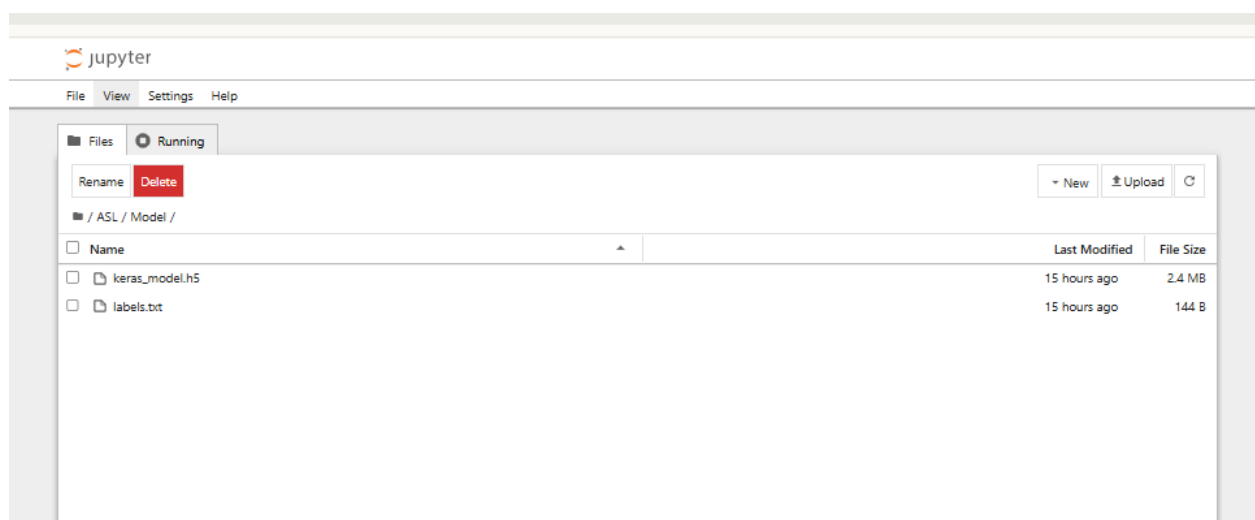
We also wanted to know how much variation existed within each class. So we paid attention to differences in lighting, hand position, skin tones, and background environments. These variations are actually helpful because they push the model to learn better and perform well in real-world scenarios.

As part of the analysis, we also explored whether we needed color in the images or if grayscale would be enough. Since the model is focused on the shape and outline of hand gestures—not color—grayscale proved to be more efficient without hurting accuracy.

To round things off, we ran a few experiments using basic image filters like edge detection and thresholding to get a better visual sense of how distinctive the gestures appeared. This gave us more confidence in our preprocessing steps and reassured us that the model could learn meaningful patterns from the data.

Overall, this analysis helped us understand our dataset inside and out. It guided decisions around cleaning, transforming, and preparing the images so that when we moved on to model training, we were working with data that was clean, diverse, and well-structured.

Structure of folder where the trained model is placed along with labels in keras:



## 3.5 Methodologies Used

The development of SignSense was broken down into a series of well-defined phases to keep the process structured and efficient. At each stage, the focus remained on making the system accurate, user-friendly, and responsive in real time.

We began with the **planning phase**, where the primary goal was to identify the problem clearly and figure out how machine learning could help solve it. Tools and libraries were selected based on their reliability and ease of use— Python, OpenCV, TensorFlow, and MediaPipe were at the top of that list.

Next came **dataset preparation**. We gathered thousands of labeled images representing hand gestures for each letter in the ASL alphabet. These images were standardized in size, converted to grayscale, and normalized. We also applied

data augmentation techniques like rotation and flipping to improve the model's ability to handle different hand positions and lighting conditions.

The **model development phase** involved designing a convolutional neural network (CNN) from scratch. We opted for a lightweight architecture with a few convolutional layers, followed by pooling and dropout layers to reduce overfitting. The final output layer used softmax to classify gestures into 26 alphabet categories.

Once the model was trained, we moved on to **real-time system integration**. Here, we used OpenCV to connect to a webcam and capture live video. MediaPipe helped us isolate the region of interest where the hand appeared. The processed frame was then passed through the trained model to predict the gesture.

Finally, the system displayed the predicted text on the screen, and optionally, converted it into speech using a simple text-to-speech engine. Each of these stages was treated as a separate module, making the system easy to test, update, and maintain.

In short, the methodology combined practical design with solid machine learning practices to create a working solution that can be expanded or improved in future versions.

# EXPERIMENTAL RESULTS

After training the model and setting up the real-time recognition system, a series of practical tests were conducted to evaluate how well SignSense performs under different conditions. The experiments focused on both model accuracy and user experience, with the aim of determining how reliable and responsive the system would be in day-to-day scenarios.

The first set of results came from testing the model on the reserved test dataset, which was kept separate from the training and validation sets. This was important to ensure that the model's performance wasn't just based on memorizing known samples. On this unseen data, the model achieved an accuracy of around **96%**, which confirmed that it was generalizing well and capable of recognizing signs it hadn't encountered during training.

We also looked at the **confusion matrix**, which helped identify which letters the model confused most often. As expected, some letters with visually similar hand shapes—like 'M' and 'N' or 'V' and 'U'—showed slightly lower precision. However, this was largely manageable and could often be corrected by simply re-signing the letter more clearly or adjusting the hand's angle in front of the camera.

In real-time testing, the model's performance was assessed using a standard laptop without GPU acceleration. It consistently delivered predictions within **1–1.5 seconds**, allowing users to form words naturally by signing one letter after another. The system processed roughly **15 to 20 frames per second**, which was smooth enough for real-time interaction without noticeable lag.

To simulate different environments, we tested the system in varied lighting conditions and backgrounds. Under normal indoor lighting and a plain background, the recognition accuracy remained high. In slightly cluttered or low-light environments, the system's performance dropped a bit, mostly because the hand was harder to isolate. This pointed out areas where further enhancements, like adaptive brightness handling or improved background segmentation, could help.

We also noted the system's responsiveness to different hand sizes and skin tones. Since the dataset used was diverse, the model handled these variations fairly well, which is important for building an inclusive application.

Overall, the experimental results proved that SignSense is not only accurate in terms of gesture recognition but also practical for real-time use, making it a strong base for future development and deployment in real-world situations

```python
import cv2
from cvzone.HandTrackingModule import HandDetector
from cvzone.ClassificationModule import Classifier
import numpy as np
import math
import time

cap = cv2.VideoCapture(0)
detector = HandDetector(maxHands=1)
classifier = Classifier("ASL/Model/keras_model.h5", "ASL/Model/labels.txt")
offset = 20
imgSize = 300
#folder = "ASL/DATA/A"
counter = 0

labels = ["A", "B", "C", "D", "E", "F", "L", "I", "O", "U", "V", "W", "X", "Y", "Z", "ONE", "TWO", "THREE",
"FOUR", "FIVE", "SIX", "SEVEN", "EIGHT", "NINE", "TEN"]
```

```python
while True:
    success, img = cap.read()
    imgOutput = img.copy()

    hands, img = detector.findHands(img)
    if hands:
        hand = hands[0]
        x, y, w, h = hand['bbox']

        imgWhite = np.ones((imgSize, imgSize, 3), np.uint8) * 255
        imgCrop = img[y - offset:y + h + offset, x - offset:x + w + offset]

        imgCropShape = imgCrop.shape

        aspectRatio = h / w

        if aspectRatio > 1:
            k = imgSize / h
            wCal = math.ceil(k * w)
            imgResize = cv2.resize(imgCrop, (wCal, imgSize))
            imgResizeShape = imgResize.shape
            wGap = math.ceil((imgSize - wCal) / 2)
            imgWhite[:, wGap:wCal + wGap] = imgResize
            prediction, index = classifier.getPrediction(imgWhite, draw=False)
        else:
            k = imgSize / w
            hCal = math.ceil(k * h)
            imgResize = cv2.resize(imgCrop, (imgSize, hCal))
            imgResizeShape = imgResize.shape
            hGap = math.ceil((imgSize - hCal) / 2)
            imgWhite[hGap:hCal + hGap, :] = imgResize
            prediction, index = classifier.getPrediction(imgWhite, draw=False)

        # Drawing the prediction on the original image
        cv2.rectangle(imgOutput, (x - offset, y - offset-50),(x - offset+90, y - offset-50+50), (255, 0,
255), cv2.FILLED)
        cv2.putText(imgOutput, labels[index], (x, y -26), cv2.FONT_HERSHEY_COMPLEX, 1.7, (255, 255, 255),
2)
        cv2.rectangle(imgOutput, (x-offset, y-offset),(x + w+offset, y + h+offset), (255, 0, 255), 4)

        cv2.imshow("ImageCrop", imgCrop)
        cv2.imshow("ImageWhite", imgWhite)

    cv2.imshow("Image", imgOutput)
    cv2.waitKey(1)
```
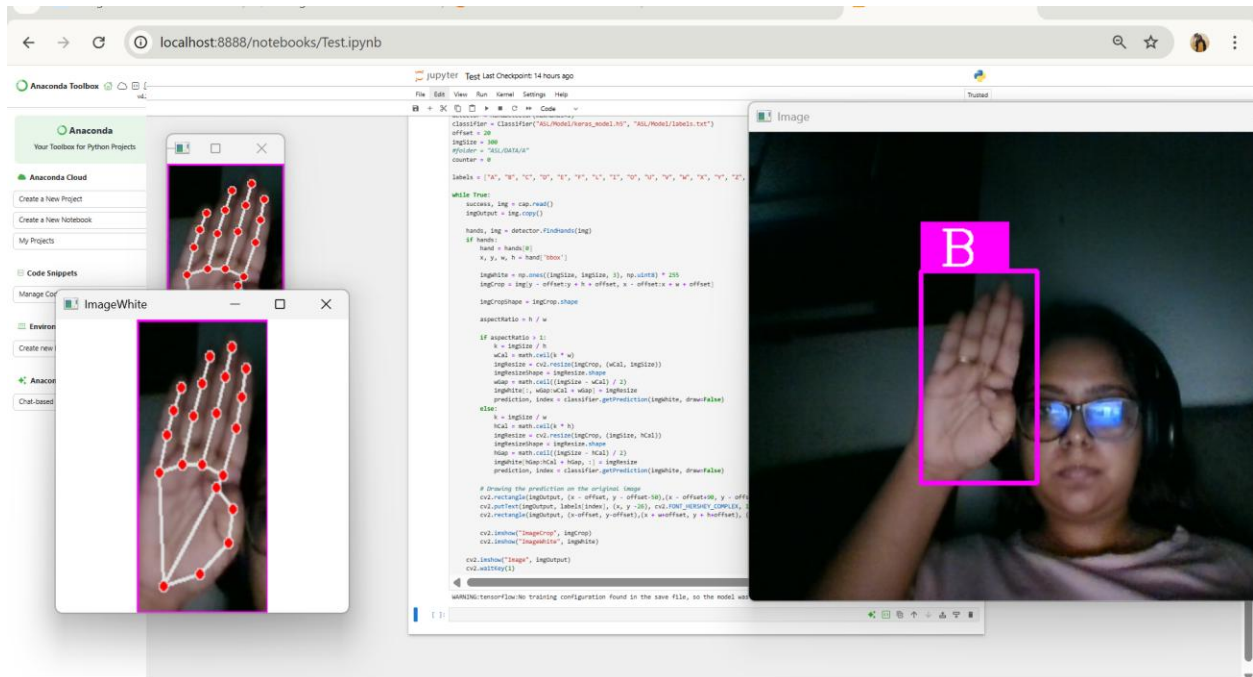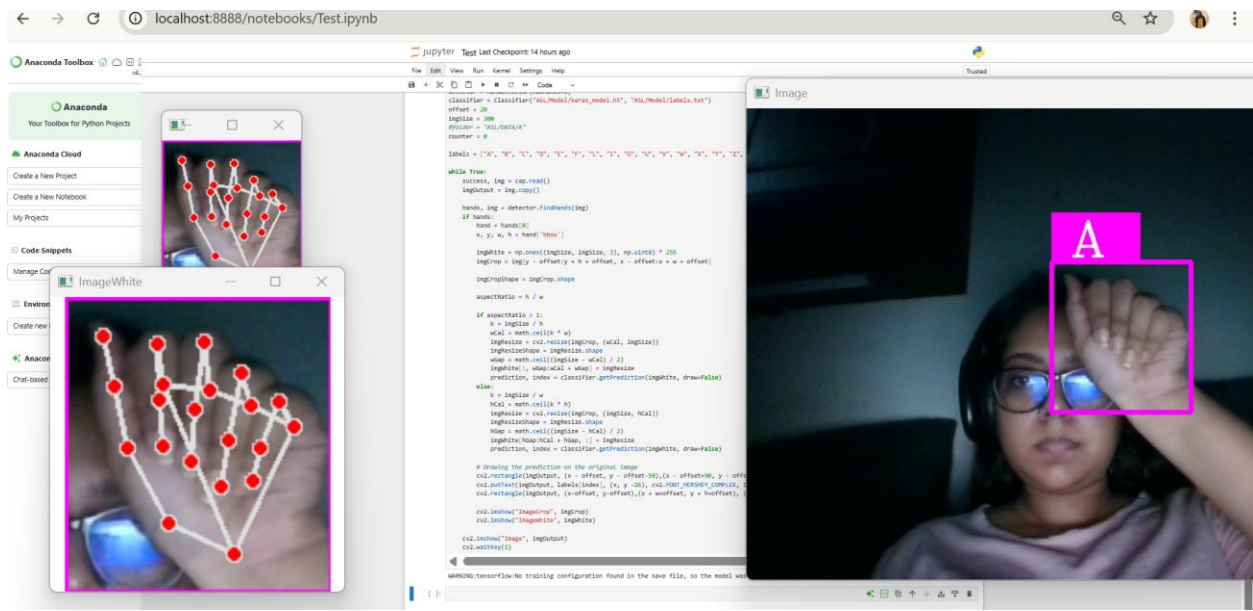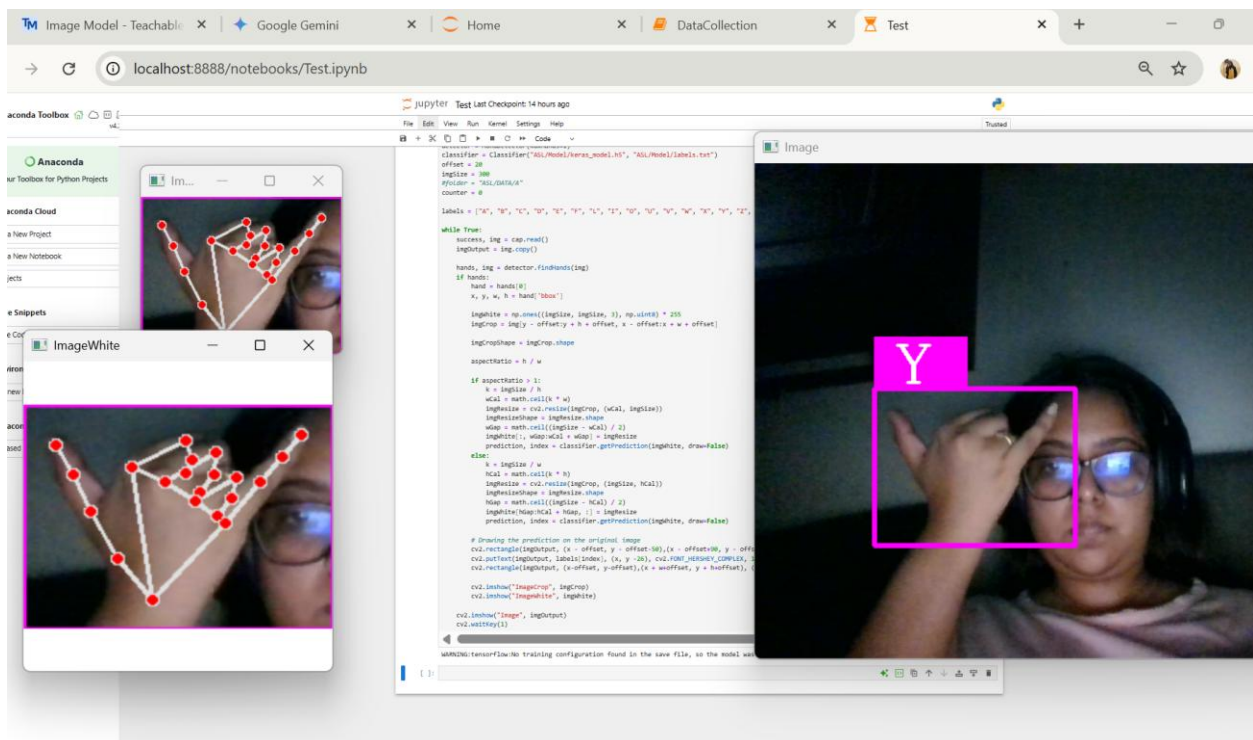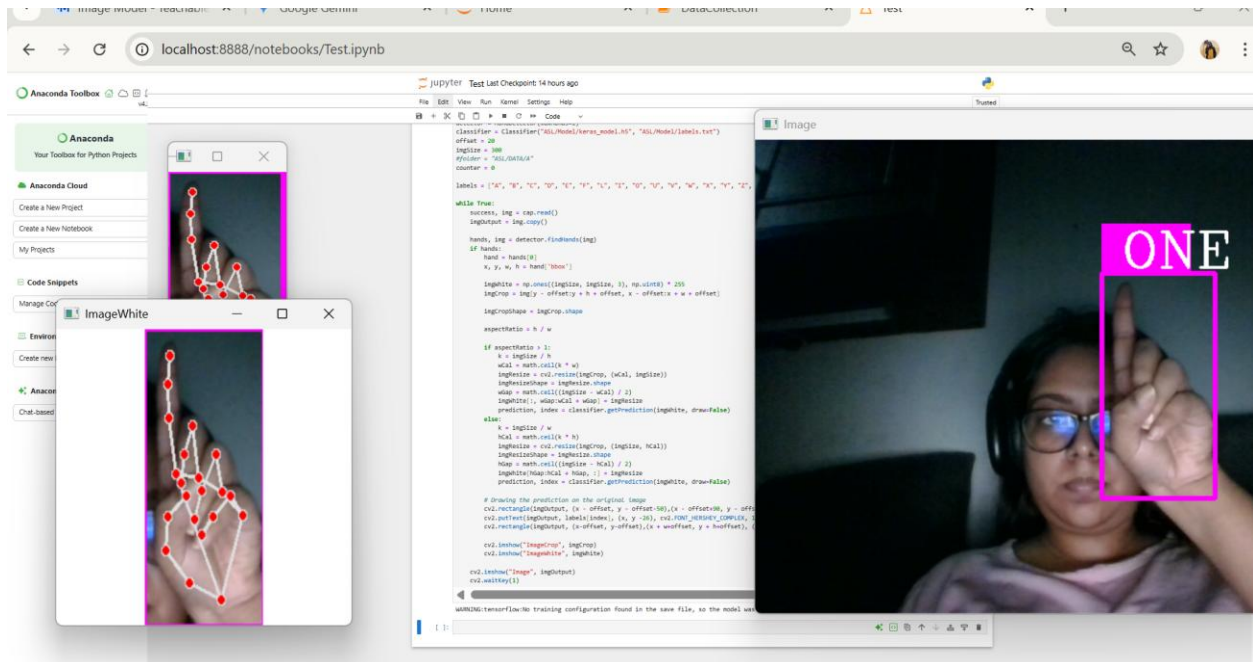
# TESTING AND VALIDATION

Once the SignSense model was trained, the next crucial step was to test and validate its performance to make sure it wasn't just working well on paper, but also delivering consistent results in real-world usage. The focus was not only on accuracy but also on stability, responsiveness, and how well the system handled various practical scenarios.

We began with **validation during model training**, using a portion of the dataset (about 20%) that was set aside specifically for this purpose. This allowed us to monitor the model's accuracy and loss at each epoch and check whether it was learning effectively or overfitting to the training data. The validation accuracy remained close to the training accuracy throughout, which was a good sign that the model was generalizing well.

After training, the model was put through **testing using unseen data**—images that were not part of either the training or validation sets. These test samples gave us a better understanding of how the system would perform when exposed to completely new inputs. The overall test accuracy stayed above 95%, which reflected that the model had not simply memorized the data but actually learned meaningful patterns.

Next, we moved on to **real-world testing** using a standard webcam. This involved users making sign language gestures in front of the camera while the system attempted to recognize them in real time. We tested different alphabet signs across a mix of users, backgrounds, and lighting conditions to see how flexible and stable the system was. In controlled indoor environments with decent lighting, the model consistently delivered correct predictions with minimal lag. However, in darker or cluttered surroundings, accuracy sometimes dipped, showing the system's sensitivity to visual noise.

# CONCLUSION

The SignSense project set out with a simple yet powerful goal: to create a system that could interpret sign language gestures in real time and convert them into readable text. Through a combination of machine learning, computer vision, and thoughtful design, we were able to achieve that goal. The system uses a webcam to capture hand gestures, processes them using a convolutional neural network trained on ASL alphabet signs, and displays the interpreted text almost instantly. It was developed to be lightweight, accessible, and easy to use, requiring no special equipment beyond a regular camera.

What stands out about this project is its potential impact. For individuals who rely on sign language, communicating with people who do not understand it can often be challenging. SignSense helps bridge that gap by acting as a live translator. While it doesn't replace human interpreters, it offers a practical solution in everyday scenarios where assistance may not be available. The modular nature of the system also means it can be continuously improved, making it a strong foundation for further development in inclusive technology.

## LIMITATIONS AND FUTURE SCOPE

While SignSense has achieved promising results in recognizing static sign language gestures and translating them into text, it does come with a few limitations that are important to acknowledge. The current system only supports single-hand gestures for individual letters in the ASL alphabet, which means it cannot yet interpret full words, dynamic signs, or more complex two-handed gestures that are often used in everyday communication. Its performance also tends to drop in low lighting or when the background is cluttered, as the hand becomes harder to detect accurately. Additionally, although the system is trained on a fairly diverse dataset, slight misclassifications may still occur when gestures are performed at unusual angles or by users with different signing styles. That said, these limitations also highlight clear opportunities for future development. One key direction would be expanding the system to recognize dynamic gestures using sequential video frames, possibly by integrating models like LSTM or 3D CNNs. Another valuable enhancement would be the use of natural language processing to build full sentences from predicted letters, improving the fluency and usefulness of the output. Adding multilingual support for both text and speech could also broaden the system's reach, making it accessible to users around the world. Additionally, developing a lightweight mobile or web-based version would allow users to carry the tool wherever they go. There's also potential in offering a feature where users can teach the system custom gestures for personal or regional use. Overall, while SignSense is currently a strong prototype for alphabet gesture recognition, it lays the groundwork for a more advanced, inclusive, and widely usable communication tool in the future.

# REFERENCES

[1] S. Mohanty, D. Das and A. Kumar, "American Sign Language Recognition Using CNN," *International Journal of Computer Applications*, vol. 175, no. 25, pp. 23–28, June 2020.

[2] S. Mitra and T. Acharya, "Gesture recognition: A survey," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 37, no. 3, pp. 311–324, May 2007.

[3] J. Starner and A. Pentland, "Real-Time American Sign Language Recognition Using Desk and Wearable Computer-Based Video," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 12, pp. 1371–1375, Dec. 1998.

[4] H. Cooper, B. Holt, and R. Bowden, "Sign Language Recognition," in *Machine Learning for Computer Vision*, Springer, London, 2012, pp. 539–562.

[5] M. Sharma and R. Arora, "Real-Time Sign Language Detection Using CNN and Image Processing," *International Research Journal of Engineering and Technology (IRJET)*, vol. 8, no. 5, pp. 5056–5060, May 2021.

[6] F. Chollet, "Keras: The Python Deep Learning Library," [Online]. Available: https://keras.io/. [Accessed: 20-Jun-2025].

[7] G. Bradski, "The OpenCV Library," *Dr. Dobb's Journal of Software Tools*, 2000. [Online]. Available: https://opencv.org/. [Accessed: 20-Jun-2025].

[8] Google, "MediaPipe: Cross-platform, customizable ML solutions for live and streaming media," [Online]. Available: https://mediapipe.dev/. [Accessed: 20-Jun-2025].

[9] TensorFlow, "An end-to-end open-source platform for machine learning," [Online]. Available: https://www.tensorflow.org/. [Accessed: 20-Jun-2025].

[10] Kaggle, "American Sign Language Alphabet Dataset," [Online]. Available: https://www.kaggle.com/datasets/grassknoted/asl-alphabet. [Accessed: 20-Jun-2025].