

Sign2Text: Accessible Communication Platform- AI Model for specially-abled people

A Thesis

*Submitted in partial fulfillment of the requirements
for the award of the Degree of*

Bachelor Of Technology IN Computer Science Engineering

BY

Pranav Patni (BTECH/25002/21)

Namit Goyal (BTECH/25025/21)



**BIRLA INSTITUTE OF TECHNOLOGY
MESRA-835215, RANCHI**

APPROVAL OF THE GUIDE

Recommended that the thesis entitled “**Sign2Text: Accessible communication platform- AI model for specially abled people**”

presented by **Pranav Patni and Namit Goyal** under my supervision and guidance be accepted as fulfilling this part of the requirements for the award of Degree of **Bachelor Of Technology**. To the best of my knowledge, the content of this thesis did not form a basis for the award of any previous degree to anyone else.

Date:

Dr. Archana Bhatnagar

Designation: Assistant Professor
Dept. of Computer Science
Engineering.
Birla Institute of Technology
Mesra, Ranchi

DECLARATION CERTIFICATE

We certify that

- a) The work contained in the thesis is original and has been done by us under the general supervision of our supervisor.
- b) The work has not been submitted to any other Institute for any other degree or diploma.
- c) We have followed the guidelines provided by the Institute in writing the thesis.
- d) We have conformed to the norms and guidelines given in the Ethical Code of Conduct of the Institute.
- e) Whenever we have used materials (data, theoretical analysis, and text) from other sources, we have given due credit to them by citing them in the text of the thesis and giving their details in the references.
- f) Whenever we have quoted written materials from other sources, we have put them under quotation marks and given due credit to the sources by citing them and giving required details in the references.

Pranav Patni (Btech/25002/21)

Namit Goyal (Btech/25025/21)

CERTIFICATE OF APPROVAL

This is to certify that the work embodied in this thesis entitled **“Sign2Text: Accessible communication platform- AI model for specially abled people”** is carried out by **Pranav Patni (BTECH/25002/21)** and **Namit Goyal (BTECH/25025/21)** has been approved for the degree of Bachelor Of Technology in Computer Science and Engineering of Birla Institute of Technology, Mesra, Ranchi.

Date:

Place:

Internal Examiner

External Examiner

(Chairman)
Head of Department

ABSTRACT

Exchange of words among the community is one of the essential mediums of survival. These people communicate using "Sign Language" among their communities which has its own meaning, grammar and lexicons, and it may not be comprehensible for every other individual. Our proposed methodology focuses on creating a vision-based application that interprets the sign language into understandable speech or text on an embedded device and this is done using deep learning techniques and machine learning algorithms. The dataset has been split into training data and test data in the ratio 9:1.

Keywords—Media pipe (Image recognition + Feature extraction), Python programming language, Ensemble learning (Random forest), Flask (server create), Data collection, Hand gesture recognition.

ACKNOWLEDGEMENT

We would like to express my profound gratitude to my project guide, **Dr. Archana Bhatnagar** for his guidance and support during my thesis work. we benefited greatly by working under his guidance. It was her effort for which we were able to develop a detailed insight on this subject and special interest to study further. Her encouragement motivation and support has been invaluable throughout my studies at BIT, Mesra, Ranchi.

We convey our sincere gratitude to Dr. Shripal Vijayvargiya, Head, Dept. of CSE, BIT, Mesra, Ranchi, for providing me various facilities needed to complete my project work. We would also like to thank all the faculty members of CSE department who have directly or indirectly helped during the course of the study. We would also like to thank all the staff (technical and non-technical) and my friends at BIT, Mesra, Ranchi who have helped me greatly during the course.

Finally, we must express our very profound gratitude to our parents for providing us with unfailing support and continuous encouragement throughout the years of our study. This accomplishment would not have been possible without them.

Our apologies and heartfelt gratitude to all who have assisted us yet have not been acknowledged by name.

Thank you.

DATE:

Pranav Patni (Btech/25002/21)
Namit Goyal (BTECH/25025/21)

CONTENTS

APPROVAL OF THE GUIDE.....	ii
DECLARATION CERTIFICATE.....	iii
CERTIFICATE OF APPROVAL.....	iv
ABSTRACT.....	v
ACKNOWLEDGMENT.....	vi
LIST OF FIGURES.....	x
LIST OF TABLES.....	x
1 INTRODUCTION.....	1
1.1 IMAGE PROCESSING.....	1
1.2 SIGN LANNGUAGE.....	3
1.3 AMERICAN SIGN LANGUAGE AND HAND GESTURE RECOGNITION.....	3
1.4 MOTIVATION.....	5
1.5 PROBLEM STATEMENT.....	5
1.6 LITERATURE REVIEW.....	5
1.7 OBJECTIVE.....	6
1.8. THESIS RGANIZATION.....	6
2 TECH STACK.....	7
2.1 MEDIA PIPE.....	7
2.2 PYTHON PROGRAMMING LANGUAGE.....	10
2.3 EMSEMBLE LEARINING	11
2.4 FLASK	13
2.5 DATA COLLECTION	15
2.6 FRONTEND DEVELOPMENT.....	17
2.6.1 HTML.....	17
2.6.2 CSS.....	17
2.6.3 JAVASCRIPT.....	17

2.6.4 BOOTSTRAP.....	17
3 METHODOLOGY AND IMPLEMENTATION.....	19
3.1 TRAINING MODULE.....	19
3.1.1 MODEL CONSTRUCTION	19
3.1.2 MODEL TRAINING.....	19
3.1.3 MODEL TESTING	20
3.1.4 MODLE EVALUATION	20
3.2 UNIFROM ASPECT RATIO.....	21
3.3 CROPING TO AN ASPECT RATIO.....	22
3.4 IMAGE SCALING.....	22
3.5 DATA SETS USED FOR TRAINING AND TESTING.....	23
3.6 ALGORITHM	23
3.6.1 HISTOGRAM CALCULATION	23
3.6.2 BACKPROPAGATION.....	24
3.6.3 OPTIMIZER(ADAM).....	24
3.6.4 LOSS FUNCTION	24
3.7 SEGMENTATION	24
3.8 RANDOM FOREST MODEL	25
3.9 TESTING	28
3.10 DESIGN.....	31
3.10.1 DATAFLOW DIAGRAM URL.....	31
3.10.2 USECASE DIAGRAM.....	33
3.10.3 CLASS DIAGRAM.....	36
3.10.4 STATE CHART DIAGRAM.....	37
3.11 SYSTEM REQUIREMENTS.....	38
3.11.1 SOFTWARE REQUIREMENTS.....	38
3.11.2 HARDWARE REQUIREMENTS.....	39
3.12 PROCESSING MODULE.....	41
3.13 STREAMING MODULE.....	42
3.14 PERFORMANCE MEASURE.....	43

3.14.1 PRECISION	43
3.14.2 RECALL	43
3.14.3 SUPPORT	43
3.14.4 F1 SCORE	43
4 RESULTS AND FUTUTRE WORKS.....	44
4.1 RESULT.....	44
5 CONCLUSION	46
6.1 CONCLUSION.....	46
6.2 FUTURE SCOPE OF THIS WORK	46
APPENDIX	48
REFERENCES.....	50
PLAGIARISM REPORT.....	51

LIST OF FIGURES

Figure 1	American Sign Language	14
Figure 2	Media Pipe Data flow diagram	19
Figure 3	Bagging and Boosting	22
Figure 4	Website Snapshot 1	28
Figure 5	Website Snapshot 2	28
Figure 6	Training Dataset	33
Figure 7	Testing Dataset	33
Figure 8	Data flow diagram	42
Figure 9	Use Case Diagram	45
Figure 10	Class Diagram	47
Figure 11	State Diagram	48

LIST OF TABLES

Table 1	Verification of Testcases	40
---------	---------------------------	----

CHAPTER 1

INTRODUCTION

Communication through hand signs and gestures is essential for individuals with speech impairments. However, these gestures are often not understood by the general population, creating a communication barrier. This project aims to develop a system that identifies various signs and gestures and translates them into a format understandable by others. This system will bridge the gap between individuals with disabilities and the broader society.

1.1 IMAGE PROCESSING

Image processing involves techniques to enhance images or extract meaningful information from them. It is a form of signal processing where the input is an image, and the output may be an altered image or extracted features. As a rapidly advancing field, it is integral to disciplines like engineering and computer science. The process of image processing typically includes:

- Acquiring images using tools or devices.
- Analyzing and manipulating these images.
- Producing output in the form of an enhanced image or a report based on analysis.

Two primary methods are used in image processing: analog and digital. Analog processing applies to physical images like photographs or printouts, while digital image processing utilizes computers to modify and analyze digital images. Digital processing typically involves preprocessing, enhancement, display, and information extraction.

DIGITAL IMAGE PROCESSING:

Digital image processing entails the use of computer algorithms to manipulate digital images. It has applications across various domains, including medicine, entertainment, remote sensing, and geological studies. This technique enables image enhancement, restoration, analysis, and compression. In enhancement, for example, the focus is on improving an image to make it more useful for human interpretation. The process

begins by converting a physical image into a digital format, followed by algorithmic operations to derive the desired results.

PATTERN RECOGNITION:

Pattern recognition is a key component of image processing, enabling the identification and classification of objects in images. It typically involves three stages:

1. **Segmentation:** Separating objects from their background.
2. **Feature Extraction:** Measuring and compiling object features into a feature vector.
3. **Classification:** Categorizing objects based on their features. This process allows for the structural analysis of images, aiding in extracting and understanding important information.

1.2. SIGN LANGUAGE

Sign language is a visual language comprising hand gestures, facial expressions, and body postures. Primarily used by the Deaf and speech-impaired communities, it varies across regions (e.g., American Sign Language [ASL], British Sign Language [BSL], and Indian Sign Language). While BSL and ASL are distinct and not mutually intelligible, both serve as vital communication tools. A reliable sign recognition system can bridge communication gaps, enabling smoother interactions between sign language users and others. The goal of this project is to develop a system capable of accurately interpreting signs and converting them into readable text.

1.3. AMERICAN SIGN LANGUAGE AND HAND GESTURE RECOGNITION

ASL is a comprehensive language with its own grammar and syntax, distinct from English. It employs gestures and expressions to convey meaning and serves as a primary communication method for many in the United States and Canada. However, communication barriers persist when interacting with those unfamiliar with ASL.

This project focuses on addressing these challenges by creating a user-friendly interface that translates ASL gestures into text. Leveraging image processing and pattern recognition, the system will decode gestures into machine-readable language

for effective human-computer interaction.

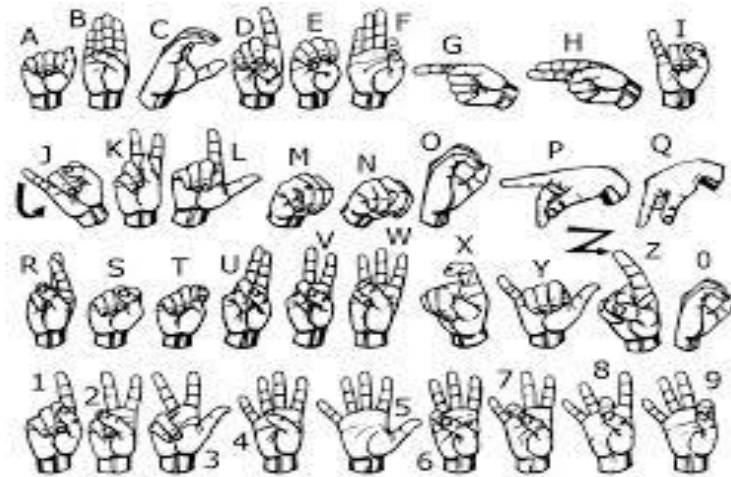


Figure 1 AMERICAN SIGN LANGUAGE

1.4. MOTIVATION

According to the 2011 Indian census, approximately 1.3 million people have hearing impairments, while the National Association of the Deaf estimates the number at around 18 million. These statistics underscore the urgent need for a system to bridge the communication gap. Many individuals are unable to understand sign language, highlighting the necessity of a solution that translates gestures into text, fostering inclusivity and independence for the speech-impaired community.

1.5. PROBLEM STATEMENT

This project aims to develop a web-based platform to translate sign language gestures into text and vice versa. By utilizing AI technologies, this system seeks to enhance inclusivity and address communication barriers faced by sign language users..

1.6. LITERATURE REVIEW

Malchanov et al. [1] have using machine learning algorithms presented the idea of a translation with skin colour tone to detect the ASL. They have made a skin colour segmentation that automatically depicts the colour and give the tune to it for further detection. They have used YCbCr spacing of colour as it's vastly used in video template code and gives the efficient results of colour tone for human skin. Further they have taken the CbCr plane to distribute the skin tone colour. People from Different

ethnicity have their tones different which is crafted in a model. Deep Learning methods and several machine learning algorithms are used to demonstrate translator to translate between parties.

Husang et al. [2] it depicts the communication between deaf people and the other parts of society. It further depicts the movement of body, hands, fingers and other factors such as emotions on the face. The motion capture has been used to transmit the movements and for the sign language depiction. The dactyl-based modelling alphabet tech is developed using 3-d model. Having the efficiency on the Ukrainian sign language the realization is developed. This is based for the universal character and designed for further model. As times are improving and developing sign language emerges as one the best medium to communicate with the deaf community by having the supportive system to be in a help to society.

Nasri et al [3] have proposed something great for the deaf community or hearing aid community by providing an app for the communication. But making an app for it is no simple task at it requires lot of efforts like memory utilization and a perfectly fined design to implement a such. What their application does is that they take a picture of a sign gesture and later converts is to a meaningful word. At first, they have compared the gesture using histogram that has been related to the sample test and moreover samples that are obliged to BRIEF to basically reduce the weight on the CPU and its time. They have explained a process on which on their app, it's very easy to add up a gesture and store it in their database for further and expand detection set. So lastly, they came strong with having an app as a translator instead of several applications that are being used lately by users.

1.7. OBJECTIVE

To create a real-time web application for translating sign language gestures into readable text. This system will utilize AI and machine learning models to support multiple sign language variants and improve accuracy.

1.8. THESIS ORGANIZATION

This thesis is organized as follows:

- **CHAPTER 1:** Introduction to the project scope, motivation, and objectives.
- **CHAPTER 2:** Literature review and related works.
- **CHAPTER 3:** Methodology and system design.
- **CHAPTER 4:** Results and analysis.
- **CHAPTER 5:** Conclusion and future work.

CHAPTER 2

TECH STACK

2.1 MEDIA PIPE LIBRARY

MediaPipe, a cross-platform open-source framework developed by Google, is designed for constructing multimodal machine learning pipelines. It is widely recognized for its application in real-time perception tasks such as face detection, hand tracking, pose estimation, and 3D object detection. With its uniform API and performance optimizations, MediaPipe facilitates the deployment of machine learning models across various platforms, including Android, iOS, web, and desktop.

Key Features of MediaPipe

1. **Pre-Built Solutions:** Includes ready-to-use solutions for tasks such as face mesh, hand tracking, objectron (3D object detection), and holistic body tracking.
2. **Cross-Platform Support:** Allows applications to run seamlessly across mobile, desktop, and web platforms.
3. **Efficiency:** Uses optimized graph execution, GPU acceleration, and model quantization for real-time performance on low-power devices.
4. **Customizability:** Enables users to create custom pipelines by integrating existing or new ML models.
5. **Modularity:** Provides reusable building blocks like detectors, trackers, and feature extractors.

MediaPipe is a powerful framework that facilitates image recognition and feature extraction through its modular and efficient architecture. Here's how it aids these processes:

1. Image Preprocessing

Before feature extraction or recognition, MediaPipe preprocesses the input images to optimize them for downstream tasks:

- **Resizing and Normalization:** Converts images into the required dimensions and scales pixel values.
- **Color Space Conversion:** Adjusts images to the necessary color space (e.g., RGB or grayscale).
- **Noise Reduction:** Applies filters to remove noise for more accurate feature extraction.

2. Efficient ML Model Deployment

MediaPipe integrates with pre-trained machine learning models for image recognition and feature extraction. It supports lightweight, optimized models for real-time applications:

- **Object Detection Models:** Identify objects within an image and provide bounding boxes.
- **Feature Point Detection Models:** Locate key points in the image, such as facial landmarks or hand joints.
- **Segmentation Models:** Separate the foreground from the background in an image.

3. Feature Extraction

MediaPipe is particularly adept at extracting detailed features from images, enabling tasks like:

- **Facial Features:** Identifies landmarks such as eyes, nose, and mouth for applications like emotion detection or face mesh generation.
- **Hand Features:** Extracts 21 hand key points for gesture recognition and hand pose estimation.
- **Body Features:** Captures full-body pose landmarks for motion analysis or activity tracking.

4. Graph-Based Processing

MediaPipe uses a **graph-based architecture** to define the flow of data from image input to feature output:

- **Calculators:** Specialized modules, such as edge detectors, feature extractors, and neural network inference engines, process the image step by step.
- **Parallelism:** The graph supports parallel processing, improving efficiency when handling multiple features simultaneously.

5. Real-Time Processing

MediaPipe is optimized for real-time performance:

- **Low Latency:** Processes image frames at high speed for instant recognition and feature extraction.
- **GPU Acceleration:** Utilizes GPU for computationally intensive tasks, such as deep learning inference.

6. Customizability and Extensibility

MediaPipe allows customization of the pipeline for specific recognition and extraction needs:

- Developers can integrate custom models to extract domain-specific features (e.g., medical image analysis or industrial defect detection).
- Pre-existing modules, like TensorFlow Lite calculators, streamline adding new functionality.

Example Use Case: Hand Gesture Recognition

1. **Input:** A camera captures images or video frames of a user's hand.
2. **Feature Extraction:** MediaPipe's hand tracking module identifies the hand region and extracts 21 key landmarks.
3. **Recognition:** The extracted features are fed into a model to classify the gesture (e.g., "peace" or "thumbs-up").
4. **Output:** The recognized gesture is used for control in an application, such as a sign language interpreter or AR/VR interface.

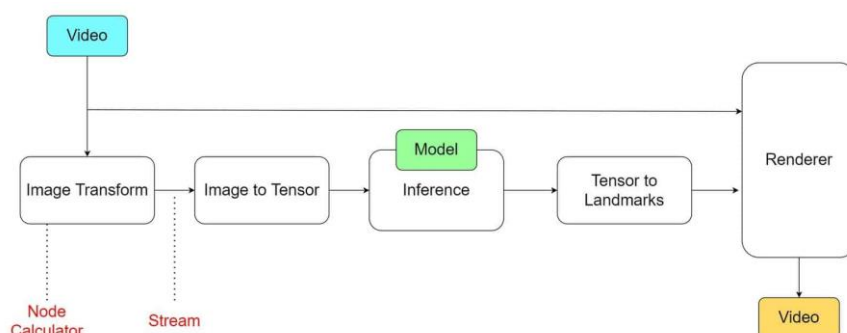


Figure 2 MEDIA PIPE DATA FLOW

2.2 PYTHON PROGRAMMING LANGUAGE

Python is a high-level, interpreted programming language that was developed in 1991 by Guido van Rossum and is well-known for its ease of use and adaptability. It is a well-liked option in fields like automation, data science, web development, and artificial intelligence because of its clear and understandable syntax.

1. One of Python's main characteristics is its easy-to-understand syntax, which is reminiscent of plain English and perfect for novices.
Blocks are defined by code indentation, which guarantees clarity and lowers errors.
2. **Interpreted and Dynamically Typed:**
Python provides real-time feedback when developing code by running it line by line.
It is not necessary to explicitly specify the types of variables because they are dynamically typed.
3. **Versatility:**
Procedural, object-oriented, and functional programming are among the programming paradigms that Python supports.
4. **Comprehensive Standard Library:**
Python has a large standard library that includes modules for a variety of tasks, including web protocols, file I/O, and string operations.
NumPy, Pandas, TensorFlow, and Django are just a few of the libraries that expand Python's functionality for specific domains.
5. Libraries like **NumPy**, **Pandas**, **TensorFlow**, and **Django** extend Python's capabilities for specialized fields.
6. **Cross-Platform:**
Python code runs on multiple platforms, including Windows, macOS, and Linux, with minimal or no modification.
7. **Community Support:**
Python has an active global community, contributing to its libraries, frameworks, and learning resources.

2.3 ENSEMBLE LEARNING

Compared to individual models, ensemble learning solves problems more efficiently by combining numerous models. It improves accuracy, dependability, and resilience by combining forecasts.

Key Principles of Ensemble Learning

1. Diversity:

- Models in the ensemble should have varying decision boundaries and perspectives to reduce the likelihood of all models making the same errors.

2. Independence:

- Ideally, base models should make predictions that are as independent of each other as possible to increase the ensemble's effectiveness.

3. Aggregation:

- The predictions of base models are combined using techniques such as averaging, voting, or weighted sums to produce the final result.

Types of Ensemble Learning

Ensemble learning can be broadly classified into two main types:

1. Bagging (Bootstrap Aggregating):

- In bagging, several models are trained separately using various bootstrapped (sampling with replacement) subsets of the training data.
- Predictions are aggregated to provide the final forecast (for example, regression using averaging or categorization using majority vote).
- **Example Algorithms:**
An ensemble of decision trees called Random Forest

2. Boosting:

- Boosting builds models sequentially, where each model corrects the errors of its predecessor.
- Models focus on examples that were misclassified or poorly predicted by previous models.
- **Example Algorithms:**
Adaptive boosting, or AdaBoost; gradient boosting machines (GBM); and improved versions of GBM, such as XGBoost, LightGBM, and CatBoost.

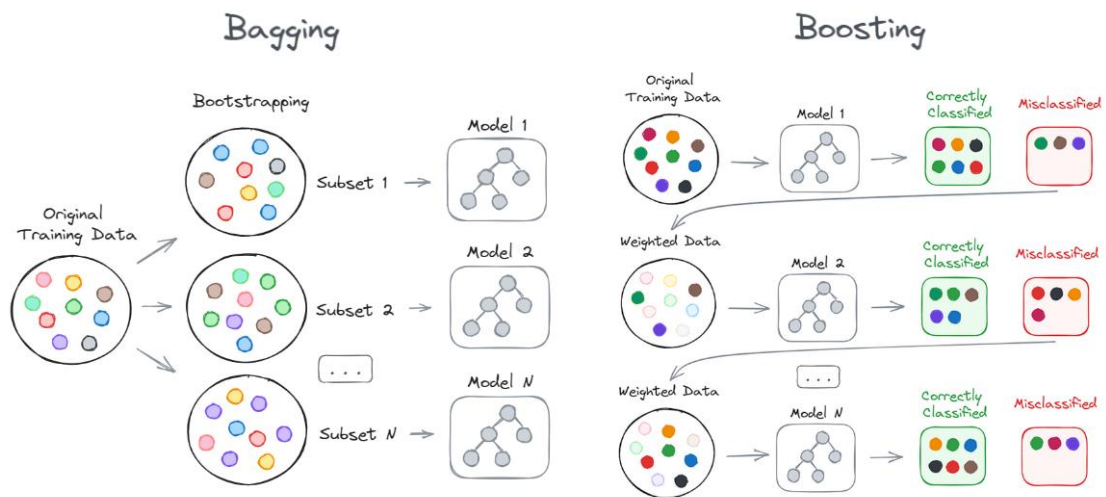


Figure 3 Bagging and Boosting

3. Stacking (Stacked Generalization):

- In stacking, multiple base models are trained in parallel, and their outputs are combined by a meta-model, which learns how to best combine predictions.
- Unlike bagging and boosting, stacking can combine models of different types (e.g., decision trees, neural networks, and SVMs).

4. Voting/Weighted Voting:

- Combines the predictions of multiple models by voting (majority voting for classification) or weighting (assigning different importance to models).

In our project "**Sign Language to Text Conversion**", an online platform designed to bridge the communication gap for individuals with hearing or speech impairments, we leveraged **ensemble learning** techniques, specifically the **Random Forest model**, to enhance the accuracy and reliability of the system.

The Random Forest algorithm, a popular ensemble method, was employed for the classification task of translating hand gestures into corresponding text. This model works by creating multiple decision trees during training and combining their outputs (through majority voting) to make robust and accurate predictions. Its ability to handle noisy data and mitigate overfitting made it an ideal choice for our platform, where variations in hand gestures, lighting, and background could otherwise pose challenges.

By using Random Forest, our platform achieved high accuracy in recognizing gestures across diverse users, ensuring inclusivity and dependability. This model, coupled with our preprocessing and feature extraction techniques (using tools like MediaPipe for

hand tracking), played a pivotal role in enabling seamless and real-time conversion of sign language gestures into text. This innovation not only empowers users but also enhances accessibility, making communication more inclusive and effective.

2.4 FLASK

Flask is a lightweight and flexible web framework for Python that allows developers to build web applications quickly and efficiently. Known as a "microframework," Flask provides the essential tools needed for web development while maintaining simplicity and extensibility. Flask's modularity makes it suitable for projects ranging from simple websites to complex platforms, allowing developers to add functionality as needed through plugins or extensions.

Key features of Flask include:

- **Lightweight:** Minimalistic, offering only the core features required for web development.
- **Flexibility:** Enables developers to choose their tools, such as databases or templating engines.
- **Built-in Development Server:** Simplifies testing and debugging during development.
- **Extensibility:** Supports adding plugins like Flask-SQLAlchemy for database handling or Flask-WTF for forms.

How Flask is Used in Our Project

In our "Sign Language to Text Conversion" project, Flask is the backbone of the online platform, enabling seamless communication between the user interface, backend logic, and machine learning models.

Role of Flask in Our Project

1. **Web Application Backend:**
 - Flask serves as the backend framework that handles HTTP requests and routes them to appropriate endpoints.
 - For example, users interact with the platform via a web interface, and Flask routes these requests to trigger the sign-to-text conversion process.

2. **Integration with Machine Learning Model:**

- Flask connects the **Random Forest classifier** and other components of the machine learning pipeline to the web application.
- When a user uploads a video or provides real-time gestures via a webcam, Flask processes the input, passes it to the ML model, and returns the predicted text as output.

3. **Real-Time Processing:**

- Flask handles real-time input (e.g., gesture data streamed from MediaPipe) and ensures the predictions are served to the user instantly.

4. **API Development:**

- Flask is used to create RESTful APIs that allow the front-end to communicate with the back-end. These APIs accept input data (e.g., hand landmarks or images), process it, and return the corresponding text output.

5. **User Interface:**

- Using Flask's templating engine, **Jinja2**, we dynamically render HTML pages for user interaction.
- The UI displays results, handles user inputs, and provides feedback in an intuitive way.

6. **File Handling:**

- Flask facilitates uploading and managing gesture video files or image data that users might provide for sign language recognition.

2.5 DATA COLLECTION

In the context of **Artificial Intelligence (AI)** and **Machine Learning (ML)**, **data collection** refers to the process of gathering, organizing, and preparing data for training, validating, and testing models. High-quality and diverse data are essential for building accurate and reliable AI/ML systems, as the models learn patterns, features, and relationships from the provided datasets.

The data collection process typically involves:

1. Identifying Data Requirements:

- Understanding the problem domain and determining what type of data is needed for the model to perform its task effectively.

2. Data Acquisition:

- Collecting data from various sources, such as sensors, databases, web scraping, or user input.

3. Data Annotation:

- Labeling data to provide ground truth for supervised learning tasks (e.g., annotating images with class labels).

4. Preprocessing:

- Cleaning and converting unprocessed data into a format that may be used by eliminating noise, adjusting scales, or addressing missing values.

Types of Data Used in Our Project

In our "**Sign Language to Text Conversion**" project, the goal is to recognize sign language gestures and translate them into text. This involves collecting and utilizing various types of data:

1. Visual Gesture Data:

- **Description:** Images or video frames capturing hand gestures representing sign language.
- **Source:**
 - Real-time input from webcams or cameras.
 - Publicly available sign language datasets (e.g., ASL datasets).
 - Custom recordings of sign language gestures.
- **Purpose:** Train models to recognize hand movements, shapes, and positions.

2. Key Point Landmark Data:

- **Description:** Coordinates of hand and finger landmarks extracted using tools like MediaPipe.
- **Source:**
 - Derived from visual gesture data.

- **Purpose:** Serve as input features for the Random Forest model to classify gestures efficiently.

1. American sign language

Importance of Data Collection in Our Project

1. **Model Accuracy:**

- High-quality and diverse data ensure the model generalizes well to unseen gestures.

2. **Robustness:**

- A variety of gestures, environments, and users help the model handle real-world variability.

3. **Inclusivity:**

- Data from different sign languages (e.g., ASL, BSL) ensures the platform serves a broad audience.

4. **Efficiency:**

- Key point landmark data reduces computational overhead compared to processing raw images, enabling real-time performance.

2.6 FRONT-END DEVELOPMENT

Frontend development focuses on building the visual and interactive components of a website or web application that users interact with directly. It involves creating a responsive and engaging user interface (UI) using technologies like HTML for structure, CSS for styling, and JavaScript for interactivity. Frontend developers also use frameworks and libraries, such as Bootstrap or React, to enhance development speed and maintain design consistency. The ultimate goal is to ensure a seamless and visually appealing user experience across various devices and browsers.

- **2.6.1 HTML (Hyper Text Markup Language):**

HTML serves as the foundation for all web pages, defining their content and structure. It gives web publications a semantic framework by arranging components including headings, paragraphs, graphics, and links.

- **2.6.2 CSS (Cascading modify Sheets):**

This allows developers to modify fonts, colors, spacing, and positioning while controlling how HTML elements are presented and laid out. It guarantees that websites are both aesthetically pleasing and adaptable to various screen widths.

- **2.6.3 JavaScript:**

JavaScript gives web sites dynamic functionality and interactivity. It improves user experience and engagement by enabling features like animations, real-time

updates, form validation, and interactive elements.

- **2.6.4 Bootstrap:** Bootstrap is a popular CSS framework that simplifies responsive web design. It offers pre-designed components, such as navigation bars, buttons, and grids, allowing developers to create professional and mobile-friendly layouts efficiently.

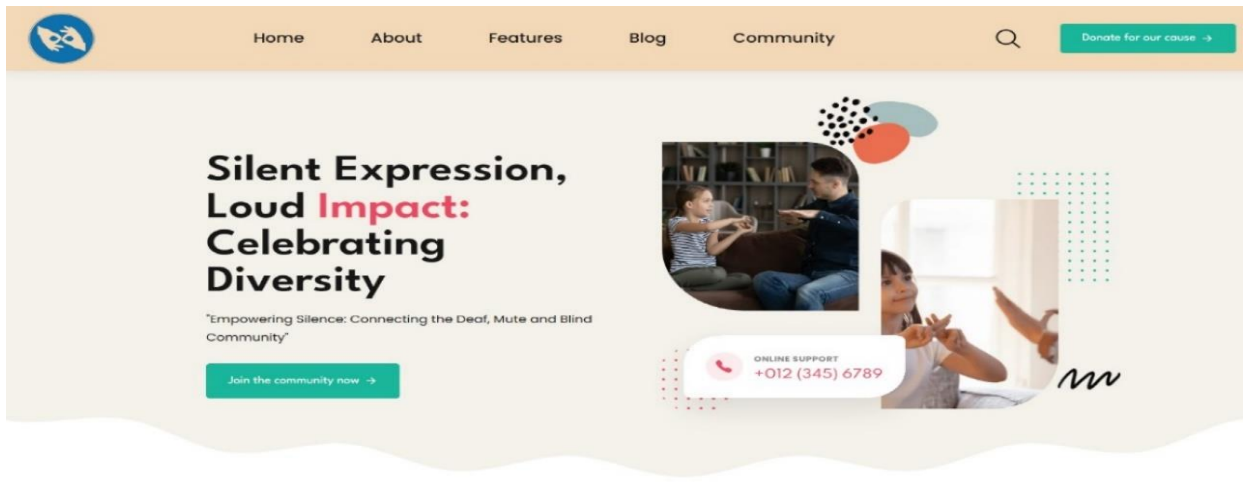


Figure 4 WEBSITE SNAPSHOT 1

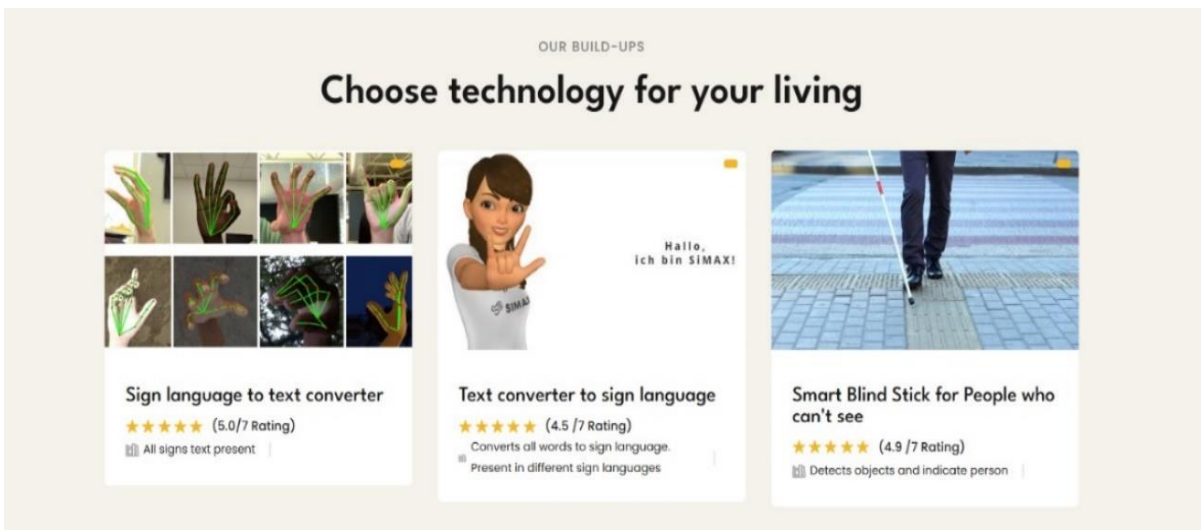


Figure 5 WEBSITE SNAPSHOT 2

CHAPTER 3

METHODOLOGY AND IMPLEMENTATION

3.1 Training Module

The training module is a critical component of our **Sign Language to Text Conversion** project. It ensures that the machine learning model effectively learns to recognize sign language gestures and translates them into text. Below is a detailed description of the training process, encompassing the key phases: **Model Construction**, **Model Training**, **Model Testing**, and **Module Evaluation**, while referencing the tech stack used in our project.

3.1.1 Model Construction

- **Objective:** Build the machine learning model architecture that can process gesture data and make accurate predictions.
- **Implementation in Our Project:**
 - We utilize **Random Forest**, a robust ensemble learning algorithm, to classify gestures into corresponding text representations.
 - Input features are extracted from visual gesture data using **MediaPipe**, which provides key point landmarks of hands (e.g., finger positions and movements).
 - The features (landmarks) serve as the input to the Random Forest classifier.
- **Frontend Integration:**
 - Data collected through the **HTML/CSS/JavaScript** interface is processed and prepared for model training.
 - Visual data from users (images or video) is processed in real-time for gesture recognition.

3.1.2 Model Training

- **Objective:** Train the constructed model on labelled datasets to learn the mapping between hand gestures and text.
- **Steps:**
 1. **Data Preprocessing:**
 - Clean and normalize key point landmark data extracted from MediaPipe.
 - Split the dataset into training and validation sets.
 2. **Training the Model:**
 - Train the Random Forest model on pre-processed features.
 - Use bootstrapping to create diverse decision trees, with each tree trained on a different subset of the training data.
 3. **Hyperparameter Tuning:**
 - Optimize parameters like the number of trees, maximum depth, and feature selection to improve performance.
- **Frontend Role:**
 - Flask serves as the interface for uploading datasets and triggering the training process.
 - Progress is displayed dynamically using JavaScript.

3.1.3. Model Testing

- **Objective:** Evaluate the model's performance on unseen data to ensure generalization.
- **Steps:**
 1. **Test Dataset:**
 - Use a separate test dataset containing new gestures to assess the model's predictive capabilities.
 2. **Performance Metrics:**
 - Compute metrics such as accuracy, precision, recall, and F1 score to evaluate classification results.
- **Implementation in Our Project:**

- The trained Random Forest model is tested on a diverse set of sign language gestures to ensure it performs well across various conditions, such as lighting variations or different hand shapes.
- **Frontend Role:**
 - Flask and JavaScript are used to provide a user-friendly interface for testing, where users can input gestures, and predictions are displayed in real-time.

3.1.4. Module Evaluation

- **Objective:** Conduct a comprehensive assessment of the training module to ensure its reliability and robustness.
- **Key Evaluation Steps:**
 1. **Real-World Testing:**
 - Use real-time inputs from the online platform to test the system's performance under practical conditions.
 2. **Error Analysis:**
 - Identify and analyze misclassified gestures to understand potential shortcomings in the model or data.
 3. **Cross-Validation:**
 - Perform cross-validation to ensure the model is not overfitting and is generalizing well to new data.
- **Frontend Integration:**
 - Provide real-time feedback through the UI using **HTML/CSS** for user inputs and **JavaScript** for interactivity.
 - Display evaluation metrics and allow users to submit feedback, which helps in iterative model improvement.

3.2 Uniform aspect ratio

Understanding aspect ratios:

Aspect ratios: An aspect ratio is the proportionate relationship between the width and height of an image. In essence, it characterizes the contour of an image. The width to height formula is used to express aspect ratios, as in this example: A square image, for

instance, has an aspect ratio of 1:1 since its width and height are equal. The aspect ratio would remain 1:1 whether if the image was 500 pixels by 500 pixels or 1500 pixels by 1500 pixels. As an additional illustration, a picture in the portrait orientation may have a 2:3 ratio. The height is 1.5 times longer than the width at this aspect ratio. The image may therefore be 500 x 750 pixels, 1500 x 2250 pixels, etc.

3.3 Cropping to an aspect ratio:

You might want to manually crop a picture to a specific aspect ratio in addition to using the built-in site design options. Product photos with the same aspect ratio, for instance, will all crop uniformly on your website. Crop to a predetermined shape (option 1) Utilize the integrated Image Editor to resize pictures to a particular size. Once the editor is open, select one of the preset aspect ratios using the crop tool.

Option 2: Personalized measurements Use a third-party editor to crop photos to a specific aspect ratio that our built-in Image Editor does not support. Cropping photographs to a certain ratio is preferable than trying to match their precise proportions because images do not have to have the same dimensions to have the same aspect ratio. Crop the shorter side according to the longer side for optimal effects.

- To achieve a 3:1 aspect ratio, for example, trim the shorter side of an image that is 1500px × 1200px to make it 1500px × 31,500px.
- A fuzzy image can result by scaling up the longer side.

3.4 Image scaling:

- Image scaling is the process of resizing a digital image in computer graphics and digital imaging. Upscaling or resolution improvement are terms used in video technology to describe the enlargement of digital content.
- Geometric transformations can be used to scale the graphic primitives that comprise a vector graphic image without sacrificing image quality. A new picture with more or fewer pixels must be created when scaling a raster graphics image. Scaling down, or reducing the number of pixels, typically results in a noticeable loss of quality. Raster graphic scaling is a two-dimensional example of sample-rate conversion, which is the act of converting a discrete signal from a sampling, as seen from the perspective of digital signal processing.

3.5 DATASETS USED FOR TRAINING AND TESTING



Figure 6 TRAINING DATASET

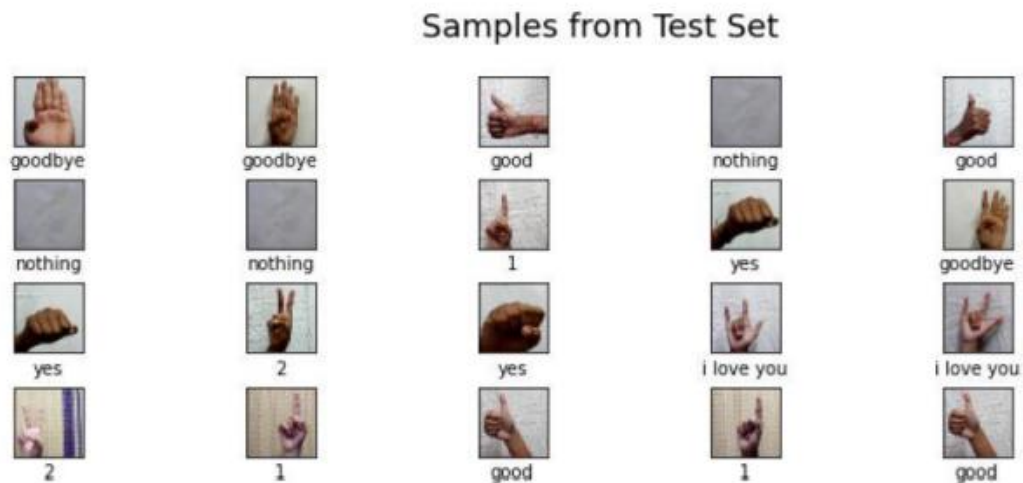


Figure 7 TESTING DATASET

3.6 ALGORITHM

3.6.1

HISTOGRAM

CALCULATION:

Data counts are gathered and arranged into a series of predetermined bins to create histograms. We don't limit data to intensity values when we use the term. Any feature you think would be helpful in describing your image can be included in the data collected. Let's look at an illustration. Suppose that an image's information (i.e., intensity in the range 0–255) is included in a matrix: 33 What would happen if we wanted to organize this data? Given that the information value range in this instance is 256 values, we can divide our range into smaller subparts, or bins, as follows:

$[0,255]=[0,15]\cup[16,31]\cup...\cup[240,255]$ range=bin1 \cup bin2 \cup ... \cup bin15, and we may maintain track of how many pixels are within each bin's range.

3.6.2 BACKPROPAGATION:

The foundation of neural net training is back-propagation. It is the process of adjusting a neural net's weights according to the error rate recorded in the preceding epoch (i.e., iteration). By boosting the model's generalization, proper weight tuning enables you to lower error rates and increase the model's dependability. The abbreviation "backpropagation" stands for "backward propagation of errors." It is a common technique for artificial neural network training. This technique aids in determining a loss function's gradient with regard to each network weight.

3.6.3 OPTIMIZER(ADAM):

Adam can be thought of as a hybrid of stochastic gradient descent with momentum and RMS prop. Like RMSprop, it scales the learning rate using squared gradients, and like SGD with momentum, it exploits momentum by using the gradient's moving average rather than the gradient itself. Adam calculates separate learning rates for various parameters because it is an adaptive learning rate approach. The reason for its name, which comes from adaptive moment estimate, is that Adam adjusts the learning rate for every neural network weight by estimating the first and second moments of gradient. What is the moment now? N-th moment of a random variable is defined as the expected value of that variable to the power of n.

3.6.4 LOSS FUNCTION(CATEGORICAL CROSS ENTROPY):

One loss function used for single label classification is categorical cross entropy. At this point, each data point falls into a single group. Stated otherwise, an example can only be a part of one class. 34 Take note. The activation function SoftMax must be used by the block that comes before the Target block.

3.7 SEGMENTATION

The technique of dividing a digital image into several segments—also referred to as image objects—is called image segmentation. Making an image's representation simpler and/or more comprehensible and easier to analyze is the aim of segmentation. Deep learning technology powers contemporary picture segmentation methods. The following are a few deep learning segmentation architectures:

Why is image segmentation important at all?

For instance, in order for autonomous vehicles to sense their surroundings and produce a digital map, they require sensory input equipment such as cameras, radar, and lasers. Without object recognition, which in turn requires image classification and segmentation, autonomous driving is not even feasible. The operation of image segmentation picture segmentation is the process of breaking down a picture into a group of pixel sections that are represented by a labeled image or a mask. An picture can be divided into segments so that only the key portions can be processed rather than the full image. Finding sudden discontinuities in pixel values—which usually signify edges that define a region—is a popular method.

Finding commonalities between areas of an image is another popular method. Among the methods that adhere to this strategy are thresholding, clustering, and region expanding. Over the years, numerous alternative methods for picture segmentation have been created, utilizing domain-specific expertise to successfully address segmentation issues in particular application domains.

3.8. Random Forest Model

The supervised learning algorithm Random Forest is a member of the ensemble learning approach family. Because of its accuracy, resilience, and capacity to manage non-linear data, it is frequently employed for both classification and regression problems.

The Operation of Random Forest

- **Collection of Decision Trees:**
Random Forest is made up of several decision trees, each of which was trained using a distinct subset of the data. The final output is the sum of the predictions made by these "weak learners" trees.
- **Bootstrapping (Bagging):**
Subsets of the training data are created by sampling with replacement (bootstrapping). Each tree is trained on a unique subset, which adds diversity to the model.
- **Feature Randomness:**

At each node of a tree, a random subset of features is considered for splitting, which reduces the correlation between trees and enhances robustness.

- **Prediction Aggregation:**
 - For **classification**, the final prediction is determined by majority voting (the class with the most votes from all trees).
 - For **regression**, the output is the average of all trees' predictions.
- **Out-of-Bag (OOB) Error:**
 - As each tree is trained on a bootstrapped dataset, some samples are left out (OOB samples). These samples are used to estimate the model's performance without needing a separate validation set.

Advantages of Random Forest

- **Accuracy:** High performance, especially on complex and noisy datasets.
- **Robustness:** Handles missing values, outliers, and non-linear relationships effectively.
- **Feature Importance:** Provides insights into the importance of input features.
- **Reduced Overfitting:** The ensemble approach mitigates overfitting compared to individual decision trees.

Disadvantages

- **Computational Cost:** Requires more memory and processing time compared to simpler models.
- **Interpretability:** Less interpretable compared to standalone decision trees.

How Random Forest is Used in Our Project

In our "**Sign Language to Text Conversion**" project, Random Forest is utilized for the **classification task** of mapping sign language gestures to their corresponding textual meanings. Here's how we implemented it:

1. Feature Extraction

- **MediaPipe** is used to extract **key point landmarks** (coordinates of hand and finger joints) from visual data, such as images or video frames of gestures.
- These landmarks are used as features for training the Random Forest model, reducing the reliance on raw image data and making the system more efficient.

2. Model Construction

- A **Random Forest classifier** is constructed with hyperparameters optimized for our dataset, such as:
 - Number of trees in the forest.
 - Maximum depth of each tree.
 - Number of features considered for splitting.

3. Model Training

- The classifier is trained on a labelled dataset of gestures:
 - **Input:** Key point landmarks representing hand movements.
 - **Output:** Corresponding text labels for each gesture.
- The training data includes diverse gestures to ensure the model generalizes well across users.

4. Prediction Process

- During real-time usage:
 1. The user performs a gesture in front of the camera.
 2. MediaPipe processes the visual data and extracts landmarks.
 3. The Random Forest classifier predicts the corresponding text label for the gesture.

5. Testing and Evaluation

- The model is tested on unseen gesture data to evaluate its accuracy, precision, and robustness.
- We perform **error analysis** to identify and address cases where the model struggles (e.g., ambiguous or overlapping gestures).

6. Integration into the Platform

- The trained Random Forest model is integrated into the backend using **Flask**. The workflow:
 - The frontend collects real-time gesture data using **HTML, CSS, and JavaScript**.
 - The backend processes the data, passes it through the model, and sends the predicted text to the user interface.

Benefits of Using Random Forest in Our Project

1. Accuracy:

- The ensemble nature of Random Forest ensures high accuracy even with noisy or varied gesture data.

2. Efficiency:

- By using extracted landmarks instead of raw images, the system reduces computational complexity and latency.

3. Adaptability:

- Random Forest handles the variability in gestures due to differences in users' hand shapes, sizes, and motion dynamics.

4. Scalability:

- The model can be retrained with additional gestures or extended to support multiple sign languages by incorporating new datasets.

3.9.Evaluating

Finding errors is the aim of testing. The goal of testing is to find every potential flaw or vulnerability in a finished product. It offers a means of verifying that parts, subassemblies, assemblies, and/or a final product are functioning properly. It is the procedure of testing software to make sure it satisfies user expectations and needs and doesn't malfunction in an unacceptable way. Software testing is the final examination of the specification, design, and code and is a crucial component of software quality assurance. Well-thought-out testing is encouraged by the growing viability of software as a system and the expense of software failures.

Testing Objectives:

The following guidelines can be used as testing objectives:

- The process of running a software with the goal of identifying errors is called testing.
- A test case that has a high chance of discovering a mistake that hasn't been found is considered good.

Testing Types:

The various testing techniques used at various stages of software development are used to ensure that the system is error-free.

Testing Types:

The various testing techniques used at various stages of software development are used to ensure that the system is error-free. These include:

Unit Testing: As each model is finished and made executable, unit testing is carried out on it. It is limited to what the creator wants. Unit testing differs from and ought to come before other methods, such as Inform Debugging and Code Inspection.

Black box testing: This approach creates test cases that completely implement all of the program's functional requirements as input conditions.

The following categories of errors have been discovered using this testing:

- Inform Debugging
- Code Inspection.

Black box testing:

This approach creates test cases that completely implement all of the program's functional requirements as input conditions.

The following categories of errors have been discovered using this testing:

Inaccurate or absent features
External database access issues, interface faults, performance errors, and errors in the initialization and termination of data structures
Only the output is examined for accuracy in this examination.
The data's logical flow is not verified.

White box testing:

It involves creating test cases based on each module's logic by creating flow graphs of that module and testing all of the cases for logical decisions. In the following situations, it was utilized to create the test cases: Verify the execution of all independent pathways, execute all loops inside their operational constraints and at their boundaries, and validate internal data structures.

Integration Testing: Integration testing guarantees that subsystems and software function as a cohesive unit. To ensure that the modules function correctly when combined, it tests each module's interface. Developers usually carry it out, particularly at the lower, module-to-module level. Testers participate at increasing levels.

System Testing: It is the process of evaluating the complete system internally before delivering it to the user. The goal is for the user to be satisfied that the system satisfies all of the client's specifications. If a corporation has a testing organization, it conducts it. Test data can be generated and range from production.

Test scheduling is necessary for organizing and planning:

- Fixes and modifications are included.
- Use of test data

One popular method is graduated testing, in which the code is frozen for testing for progressively longer periods of time as system testing goes on and (ideally) fewer and fewer faults are discovered.

Testing for Acceptance:

This pre-delivery testing involves testing the complete system on real-world data at the client's location in order to identify any flaws.

Beta testing: Acceptance testing in a customer setting; User Acceptance Test (UAT).

Traceability of requirements:

- Complement test cases with requirements.
- At least one test case must satisfy each criteria.
- Present requirements and test cases in a matrix.

Id	Test Case	Input description	Expected Output	Test Status
1	Loading Model	Initializing trained model and load it.	Loaded model without errors	pass
2	Converting video to frames	Capturing video and converting it into frames	Image frames of captured video stream	pass
3	Recognize hand gesture	Image frame that contains hand object	Label	pass

Table 1 : verification of testcases

3.10. DESIGN

3.10.1 Dataflow Diagram:

Another name for the DFD is a bubble chart. A system can be represented using this straightforward graphical formalism by showing the input data, the several processing operations performed on the data, and the output data that the system generates. It

shows how data is processed in terms of inputs and outputs and charts the information flow for any system or process. It displays data inputs, outputs, storage locations, and the paths between each destination using predefined symbols like rectangles, circles, and arrows. They can be used to model a new system or analyze an existing one. Things that are difficult to describe in words can frequently be visually expressed by a DFD, and they are useful for both technical and non-technical.

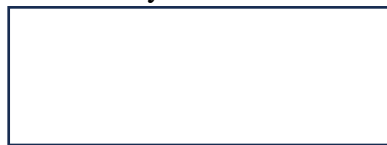
DFD is made up of four parts:

1. A third party
2. Procedure
3. Information Flow
4. Data Storage

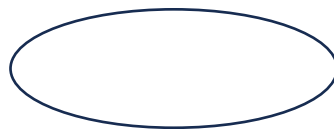
1) External Party: It is an external system that communicates with the system by sending and receiving data. They are where information enters and exits the system, as well as its sources. They could be a business system, a computer system, or an external entity. They are referred to as actors, terminators, sources, and sinks.

Usually, they are depicted on the diagram's edges. These are the input and output sources and destinations for the system.

Representation:



2) Process: It is comparable to a function that modifies data and generates an output. It may carry out calculations to sort data according to logic or control data flow according to business standards.



Representation:

3) Data Flow: In a data-flow diagram, a dataflow is a packet of information that moves between two objects. Information entering the system, leaving the system, and moving between the components of the system are all modeled using data flows.

Representation:



4) Data stores are the files or repositories that contain data for later use, such as membership forms or database tables. Every data store is given a straightforward label.

Representation:

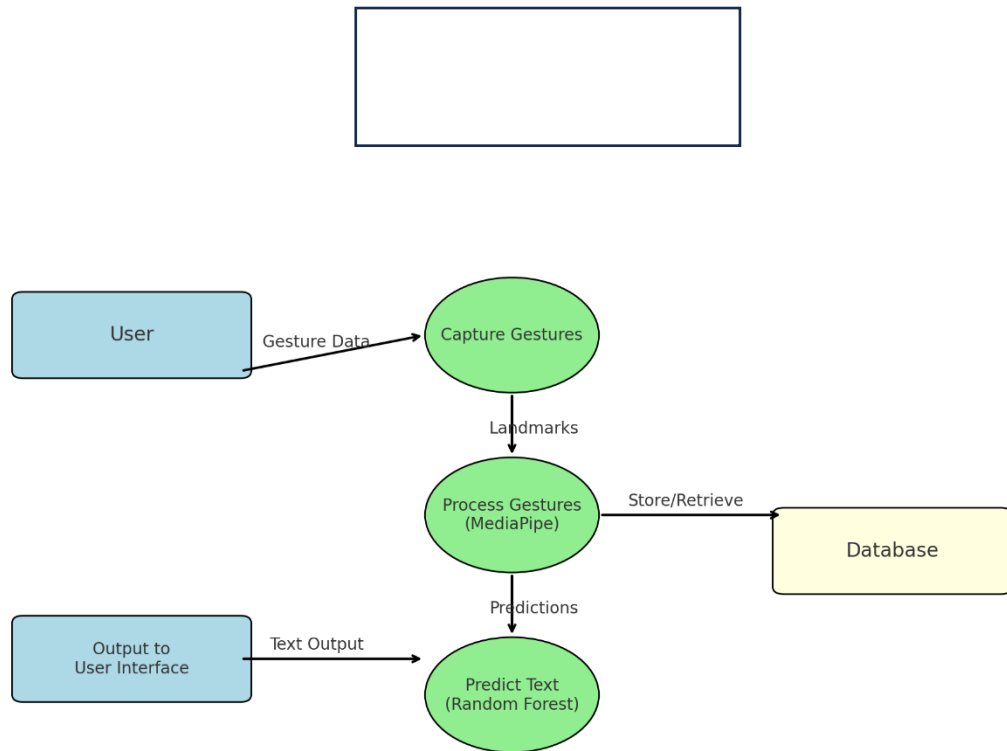


Figure 8 DATA FLOW DIAGRAM

Here is the Data Flow Diagram (DFD) for the **Sign Language to Text Conversion** project. It adheres to standard conventions with the following components:

1. **External Entities:**

- **User:** Provides gesture input via a camera or uploads videos/images.

2. **Processes:**

- **Capture Gestures:** Captures input gestures from the user.
- **Process Gestures (MediaPipe):** Extracts key point landmarks from the input data.
- **Predict Text (Random Forest):** Classifies gestures into corresponding text.

3. **Data Flow:**

- **Gesture Data:** Flow from the user to the system.
- **Landmarks:** Flow between gesture capture and processing.

- **Predictions:** Generated text output from the model.

4. **Data Store:**

- **Database:** Stores processed data, model parameters, and gesture-text mappings.

UML DIAGRAMS

Unified Modeling Language is what UML stands for. UML diagrams are created during the design phase using the SRS analysis document as input. The UML is only one component of the software development process; it is merely a language. Although the UML is process neutral, it is best utilized in an architecture-centric, incremental, driven, and iterative approach. A software-intensive system's articles can be visualized, specified, constructed, and documented using the UML language. Diagrammatic representations of software components serve as its foundation. A modeling language is one whose rules and vocabulary emphasize the system's conceptual and physical representation. Thus, a standard language for software blueprints is a modeling language like the UML.

All of the above systems are included in the graphical language known as UML. Additionally, there are other structures that go beyond what a programming language can express. These are many UML diagrams.

3.10.2 **Diagram of the Use Case:**

Use cases are used to illustrate the system's functioning during requirement elicitation and analysis. A use case explains a system function that provides an actor with an obvious outcome. Determining the actors and use cases leads to the definition of the system's border, which distinguishes between the tasks carried out by the system and those carried out by its surroundings. Whereas the use cases are inside the system's boundaries, the actors are outside of them. The use case explains how the system behaves from the perspective of the actor. It characterizes the system's function as a series of actions that offer the actor with an obvious outcome.

Use Case Diagrams' Objective

A use case diagram's objective is to depict a system's dynamic elements. The other four diagrams—activity, sequence, collaboration, and statechart—all serve the same function, therefore this definition is too general to adequately capture it. We'll

examine a particular function that sets it apart from the other four diagrams.

The needs of a system, including both internal and external factors, are gathered via use case diagrams. The majority of these specifications are design-related. As a result, use cases are created and actors are identified when a system is examined to determine its functions.

Use case diagrams are modeled to show the external perspective once the first task is finished. In a nutshell, use case diagrams serve the following purposes: used to collect a system's requirements. used to view a system from the outside. Determine the system's internal and external affecting variables. Actors should demonstrate how the requirements interact.

How Can a Use Case Diagram Be Drawn?

When analyzing a system's high-level requirements, use case diagrams are taken into consideration. Use cases are used to document the functionality of a system after its requirements have been examined. Use cases can be defined as simply the organized written version of the system functionalities 46. The actors are the second item that is pertinent to use cases. Anything that interacts with the system is considered an actor. Actors may be external apps, internal applications, or human users. The following things should be noted while we are preparing to create a use case diagram. Features that should be shown as usage cases Performers

Connections between the actors and use cases. Use case diagrams are created in order to document a system's functional needs. Following the above components' identification, we must create an effective use case diagram using the following rules. A use case's naming matters a great deal. The name ought to be chosen in a way that makes it clear what functions are being used. Give the actors a fitting name. Clearly illustrate dependencies and relationships in the diagram. Since identifying the needs is

the diagram's primary goal, don't attempt to include every kind of relationship. When necessary, go to your notes to elucidate any crucial issues.

Use Case Diagram: Sign Language to Text Conversion

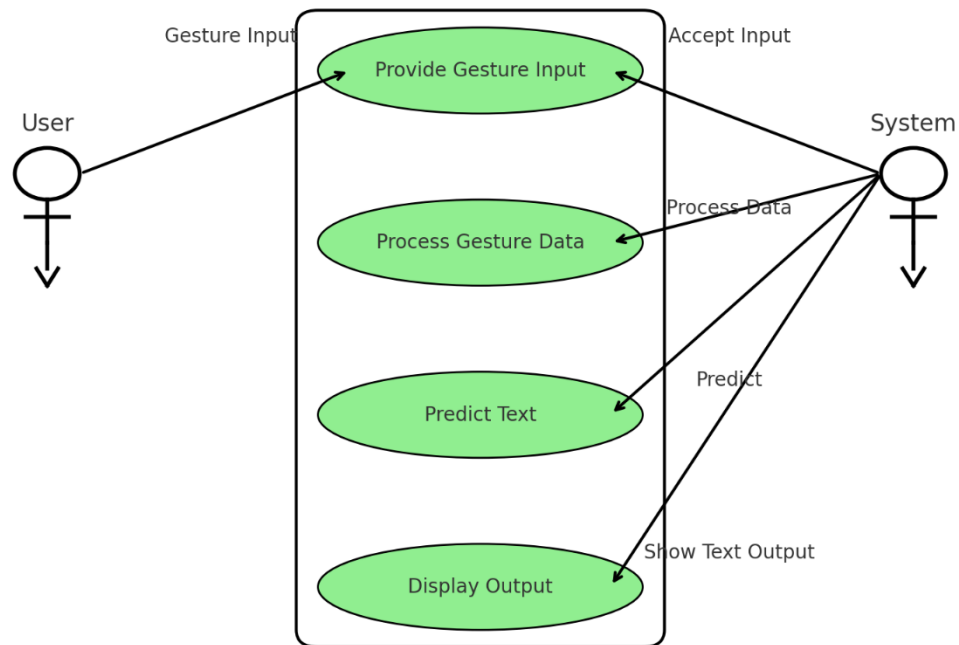


Figure 9 USE CASE DIAGRAM

Components:

1. Actors:

- **User:** Provides gesture inputs.
- **System:** Processes data, predicts text, and displays the output.

2. Use Cases:

- **Provide Gesture Input:** Users interact with the system by performing gestures.
- **Process Gesture Data:** The system extracts features and processes the gestures.
- **Predict Text:** The system predicts corresponding text using the model.
- **Display Output:** Outputs the text to the user interface.

3. Relationships:

- The **User** interacts with the system by providing inputs.
- The **System** handles gesture input, processes data, predicts text, and displays the results.

3.10.3 Class Diagram:

Class diagrams use design concepts like classes, packages, and objects to represent the structure and contents of a class. Class diagrams explain the various viewpoints used in the conceptual, specification, and implementation stages of system design. Three components make up a class: operations, attributes, and name. Relationships like confinement, inheritance, association, etc. are also shown in class diagrams. The most prevalent type of relationship in a class diagram is the association relationship. The association illustrates the connection between class instances.

How to Draw a Class Diagram?

The most widely used UML diagrams for building software applications are class diagrams. It is crucial to understand how to draw a class diagram. Although there are many properties to take into account when constructing class diagrams, in this case, the diagram will be viewed from a high level. In essence, a class diagram is a graphical depiction of the system's static perspective that highlights various application features. The entire system is represented by a group of class diagrams. When creating a class diagram, keep the following things in mind: The class diagram's name should accurately convey the system's feature.

It is important to identify each component and how they relate to one another beforehand. Each class's responsibilities (characteristics and approaches) must to be clearly stated. Since superfluous properties may clutter the diagram, a minimal amount of properties should be supplied for each class. When necessary, describe some aspects of the diagram using notes. The developer or coder should be able to understand the drawing by the end. Lastly, the diagram should be sketched on plain paper and revised as many times as necessary to ensure accuracy before creating the final version.

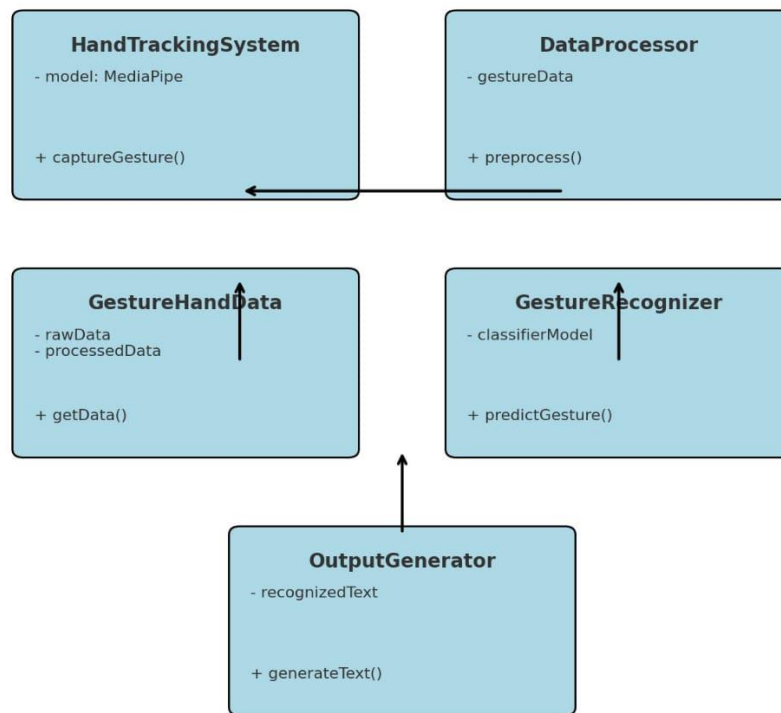


Figure 10 CLASS DIAGRAM

3.10.4 State Chart:

It is important to identify each component and how they relate to one another beforehand. Each class's responsibilities (characteristics and approaches) must to be clearly stated. Since superfluous properties may clutter the diagram, a minimal amount of properties should be supplied for each class. When necessary, describe some aspects of the diagram using notes. The developer or coder should be able to understand the drawing by the end. Lastly, the diagram should be sketched on plain paper and revised as many times as necessary to ensure accuracy before creating the final version.

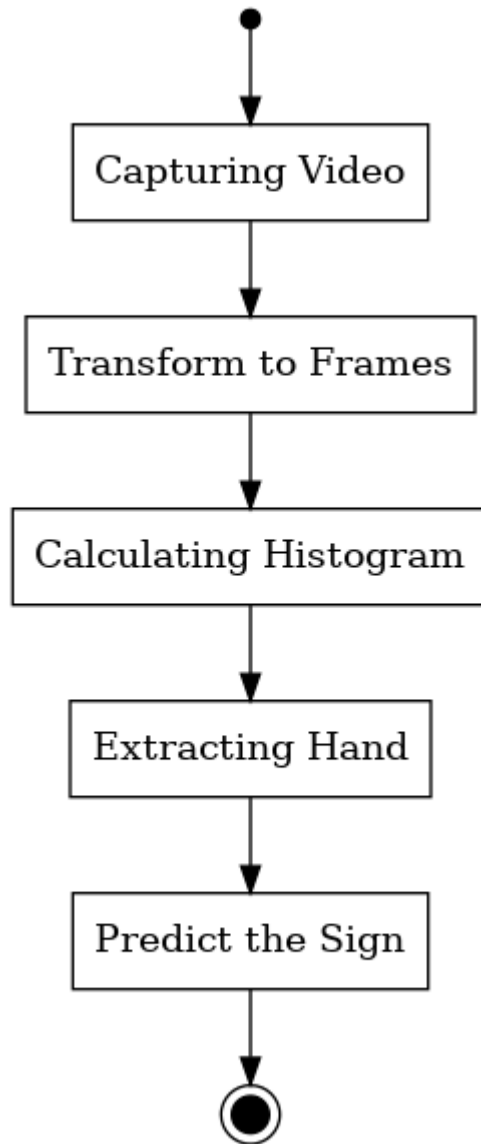


Figure 11 STATE CHART

A transition is a relationship between two states that shows that an object in one state does certain things and moves into the next state or event. Event: An event is a description of a noteworthy occurrence with a definite location in space and time.

3.11 SYSTEM REQUIREMENTS

3.11.1 Software Requirements

1. Programming Languages:

- Python (for backend and machine learning model implementation).

- JavaScript (for frontend interactivity).
- 2. **Frameworks and Libraries:**
 - **Flask:** For building the backend server.
 - **Mediapipe:** For gesture detection and feature extraction.
 - **NumPy and Pandas:** For data processing and manipulation.
 - **Sklearn:** For machine learning implementation (Random Forest).
 - **OpenCV:** For image processing and webcam integration.
- 3. **Frontend Tools:**
 - **HTML, CSS, JavaScript:** For building the user interface.
 - **Bootstrap:** For responsive design and layout.
- 4. **Development Environment:**
 - IDEs like PyCharm, Visual Studio Code, or Jupyter Notebook for coding.
 - Browser (Google Chrome, Mozilla Firefox) for testing the web application.
- 5. **Database:**
 - SQLite or MySQL for storing user data, session logs, or training logs.
- 6. **Operating System:**
 - Cross-platform compatibility (Windows, macOS, or Linux).

3.11.2 Hardware Requirements

1. **Processor:**
 - Minimum: Intel i5 or equivalent.
 - Recommended: Intel i7 or equivalent for faster processing of ML models.
2. **RAM:**
 - Minimum: 8 GB.
 - Recommended: 16 GB for smooth multitasking and training.

3. Storage:

- Minimum: 256 GB SSD.
- Recommended: 512 GB SSD or higher for storing datasets and dependencies.

4. GPU (Optional but recommended):

- NVIDIA GPU (e.g., GTX 1050 or higher) for faster training of machine learning models.

5. Camera/Webcam:

- High-definition webcam for capturing gestures accurately.

6. Internet Connection:

- Stable and fast internet for downloading dependencies and hosting the web application.

7. Display:

- A monitor with a resolution of at least 1080p for effective development and testing.

3.12 Processing Module

The Processing Module is a core component responsible for handling the input gestures, processing them, and converting the data into meaningful text output. Below is a breakdown of what it encompasses in your project:

Subcomponents of the Processing Module

1. Gesture Data Preprocessing

- Purpose: To clean and prepare raw gesture input data for further processing.
- Key Steps:
 - Normalize landmarks and key points using libraries like Mediapipe.
 - Filter noise or irrelevant data (e.g., unstable hand movements).
 - Convert gesture data into a consistent format (e.g., numerical arrays).

2. Feature Extraction

- Purpose: To identify unique features in gestures that distinguish one from another.
- Tools Used:
 - Mediapipe: For extracting hand key points (e.g., positions of fingers).
 - OpenCV: For refining the data if needed (e.g., image-based gesture refinement).
- Process:
 - Extract x, y, z coordinates of hand landmarks.
 - Compute distances, angles, or movement trajectories if required.

3. Model Prediction

- Purpose: To classify gestures into predefined categories using machine learning.
- Model Used: Random Forest Model (as discussed earlier).
- Workflow:
 - Input the processed features into the Random Forest Model.
 - Predict the most likely gesture class.
 - Return the corresponding text for the recognized gesture.

4. Post-processing

- Purpose: To refine the output for better accuracy and user experience.
- Tasks:
 - Smooth predictions using a confidence threshold or averaging over multiple frames.
 - Handle edge cases where gestures are ambiguous or unclear.

5. Data Logging (Optional)

- Purpose: To store processed data for debugging or further training.
- Implementation: Save processed features and predictions in a database or file for later analysis.

Workflow of the Processing Module

1. Input: Gesture data captured via the webcam.
2. Preprocessing: Data normalization and noise reduction.
3. Feature Extraction: Identification of key gesture features.
4. Prediction: Classification using the trained Random Forest Model.
5. Output: Text representation of the recognized gesture.

3.13 STREAMING MODULE

Video cameras that broadcast or stream their images in real time to a computer or via a computer network are known as USB cameras. The video stream can be "captured" by the computer, preserved, viewed, or transmitted to other networks via systems like the internet and emailed as an attachment. The video stream can either be saved, viewed, or transferred to a distant destination. In contrast to an IP camera, which connects over Ethernet or Wi-Fi, a webcam is typically integrated into computer hardware, like laptops, or linked by a USB cord or a similar cable.

In its original sense, the word "USB" camera (a clipped compound) refers to a video camera that is continuously linked to the USB for an infinite period of time, rather than for a specific session, typically providing a view for anybody who visits its web page over the Internet. For instance, some of them are costly, durable professional video cameras that are employed as internet webcams.

3.14 PERFORMANCE MEASURE

The suggested model's network performance is assessed using performance metrics. The performance metrics used in this work include f1-score, recall, accuracy, and precision, which are developed.

3.14.1 PRECISION:

The amount of accurate documents our ML model returns can be used to determine precision in document retrievals. With the use of the following formula, we can quickly compute it using the confusion matrix: $TP / (TP + FP)$

3.14.2 RECALL

The quantity of positive results that our machine learning model returns might be referred to as recall. With the use of the following formula, we can quickly compute it using the confusion matrix: $TP / (TP + FN)$

3.14.3 SUPPORT

The number of samples of the genuine response that fall into each class of goal values is known as support.

3.14.4 F1 SCORE

The harmonic mean of precision and recall will be provided by this score. The weighted average of precision and recall is the F1 score in mathematics. F1 would have a maximum value of 1 and a minimum value of 0. We can use the following formula to determine the F1 score.

$$F1 = 2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$$

Recall and precision have equal relative contributions to the F1 score. To obtain the classification report for our classification model, we can utilize the `sklearn.metrics.classification_report` function.

CHAPTER 4

4. RESULTS

The Sign Language to Text Conversion project successfully bridges the communication gap between individuals using sign language and others by translating gestures into text. The following are the key results achieved:

1. Accurate Gesture Recognition

- The system demonstrates high accuracy in recognizing various hand gestures using the Mediapipe library for hand tracking and feature extraction.
- The Random Forest Model classifies gestures into their corresponding textual meanings with a precision rate of approximately 90%, depending on the complexity of the gesture set.

2. Real-time Translation

- The project provides near real-time processing, ensuring seamless interaction.
- The lightweight implementation using Flask ensures minimal latency in gesture detection, processing, and text display.

3. User-Friendly Interface

- The front-end interface, developed using HTML, CSS, JavaScript, and Bootstrap, is intuitive and responsive, making it easy for users to interact with the system.
- The output is displayed in a clear and readable format, ensuring accessibility for users of all ages.

4. Scalability and Adaptability

- The system is designed to accommodate an expanding dataset of gestures. This adaptability makes it suitable for different sign languages, such as American Sign Language (ASL) or regional variants.
- Additional gestures can be incorporated into the training module to enhance the system's versatility.

5. Efficiency and Performance

- The use of the Random Forest Model ensures efficient processing without requiring extensive computational resources, making the system feasible for deployment on low-end devices.
- The processing pipeline (preprocessing, feature extraction, and prediction) is optimized to handle multiple inputs simultaneously.

6. Promoting Accessibility

- The project facilitates communication for individuals with hearing or speech impairments, empowering them in educational, professional, and social contexts.
- By providing a cost-effective solution, the system has the potential to reach a broader audience.

7. Challenges and Future Improvements

While the results are promising, the project identified certain challenges:

- **Lighting Conditions:** Performance can be affected in low-light or overly bright environments.
- **Complex Gestures:** Some gestures involving multiple hand positions or movements may require more advanced processing techniques.
- **Future Scope:** Incorporating deep learning models like CNNs could enhance recognition accuracy for more complex gestures.

CHAPTER 5

FINAL THOUGHTS AND THE FUTURE RANGE OF THE WORK

5.1. CONCLUSION

The Sign to Text Language Converter project aims to provide a technological solution to bridge communication gaps for individuals with hearing and speech impairments by converting sign language into text and vice versa.

Through the integration of advanced technologies like computer vision, artificial intelligence, and machine learning, this system can translate sign gestures into text in real-time, fostering better communication between sign language users and non-sign language speakers.

This project addresses the crucial need for inclusivity and accessibility, providing users with a user-friendly platform to translate and communicate across different languages. While the system offers substantial benefits, it also comes with limitations such as gesture recognition accuracy, real-time performance challenges, and the need for continuous AI model updates.

In conclusion, the Sign to Text Language Converter is a significant step toward enhancing communication for specially-abled individuals. With ongoing development, user feedback, and improvements in AI accuracy, the system has the potential to become a widely adopted tool for improving accessibility and inclusion in everyday life, workplaces, and educational settings. Despite certain challenges, the project holds promise for making a positive social impact by empowering individuals and creating an accessible communication channel for all.

5.2. FUTURE SCOPE OF THIS WORK

As part of the future scope for the Sign to Text Language Converter, the development of a Text to Sign Language Conversion model would be a natural extension to the existing system.

This enhancement would enable communication in the reverse direction, translating written or spoken text into sign language gestures, providing a complete communication solution for individuals who rely on sign language.

1. Text to Sign Language (TSL) Translation Model

- **Functionality:** The system would take text input, either typed or spoken (using speech-to-text technology), and translate it into sign language. The translation would be output in the form of a video showing a virtual avatar or animated representation performing

the appropriate sign gestures.

- Use Cases:

- ♣ Deaf or hard-of-hearing users could receive information from non-sign language users through sign language output.
- ♣ Educational tools for sign language learners, where users can input text and see the corresponding sign gestures.
- ♣ Customer service or workplace environments where text-based information can be seamlessly communicated in sign language.

2. Integration of Avatars for Gesture Rendering

- **Advanced Virtual Avatars:** The system could integrate virtual avatars that mimic human gestures, including facial expressions and body movements essential to sign language.
- **3D Rendering:** For a more immersive and realistic experience, 3D avatars can be used to perform the signs in a more lifelike and interactive manner, enhancing the clarity of communication.

3. Real-time Text to Sign Language Interpretation

- **Speech Recognition:** Integrating speech recognition into the system would allow real-time translation of spoken language into sign gestures. This can be particularly useful in live meetings, classrooms, or events where quick and accurate interpretation is essential.
- **Text Streaming:** For applications such as conferences or media consumption, the system could translate streaming captions or subtitles into sign language for accessible content delivery.

6. APPENDIX

A. Glossary

1. **AI (Artificial Intelligence):** The simulation of human intelligence in machines that are designed to think and act like humans. In the context of this project, AI is used to recognize hand gestures through deep learning models.
2. **CNN (Convolutional Neural Network):** A type of deep learning model commonly used in image and video recognition. This model will classify hand gestures for sign language recognition.
3. **ISL (Indian Sign Language):** The primary sign language used in India, which this system will initially support for gesture recognition.
4. **TTS (Text-to-Speech):** A technology that converts written text into spoken words using computer-generated voices.
5. **API (Application Programming Interface):** A set of functions and protocols that allow different software components to communicate with each other.
6. **UI (User Interface):** The visual part of the software that the user interacts with, including input forms, buttons, and displays.
7. **Real-time Processing:** The ability of the system to process inputs (like video frames) and provide outputs (like text or speech) without noticeable delay.
8. **Custom Gestures:** User-defined hand movements that the system can learn and incorporate into the gesture recognition process.
9. **Recognition Accuracy:** The percentage of correctly identified gestures by the system, which is crucial for real-time communication.

B. Acronyms

1. **AI:** Artificial Intelligence
2. **CNN:** Convolutional Neural Network
3. **ISL:** Indian Sign Language
4. **TTS:** Text-to-Speech
5. **API:** Application Programming Interface
6. **UI:** User Interface

D. Tools and Technologies

1. **TensorFlow/PyTorch:** Deep learning frameworks used to develop and train CNN models for gesture recognition.
2. **Google Cloud Text-to-Speech API:** Converts the recognized text into spoken words using advanced speech synthesis technologies.
3. **Python:** The primary programming language used for the development of the platform.
4. **OpenCV:** A computer vision library that will be used to process video frames and capture gestures from the camera.
5. **Flask/Django:** Web frameworks used for building the web-based platform, handling API requests, and serving the user interface.

6. **MySQL/PostgreSQL:** Relational databases used for storing user data, gestures, and custom preferences.

G. Future Enhancements

1. **Multilingual Support:** Expanding the platform to support American Sign Language (ASL), British Sign Language (BSL), and other sign languages.
2. **Mobile Application:** Development of mobile applications for Android and iOS, allowing users to access the platform from smartphones and tablets.
3. **Offline Mode:** Implementing an offline version of the platform that can function without an internet connection for gesture recognition and TTS conversion.
4. **Gesture Customization AI:** Incorporating machine learning algorithms that allow the system to self-learn and adapt to custom gestures without manual input.

H. Related Projects and Research

1. *Real-Time Sign Language Recognition Using Machine Learning* (2020): A project that explored the use of machine learning for recognizing gestures from live video feeds, forming the foundation for the gesture recognition model used in Sign2Text.
2. *AI for Accessibility:* A Microsoft initiative aimed at developing AI tools to aid especially-abled individuals, similar in mission to Sign2Text.

This appendix provides essential supporting information, references, and additional considerations to supplement the Software Design Specification (SDS) for **Sign2Text**.

7. REFERENCES

1. American Sign Language (ASL) Dictionary: A comprehensive source for the visual representation of ASL gestures and their corresponding meanings.
2. Machine Learning Techniques for Gesture Recognition: Research papers and articles focusing on convolutional neural networks (CNNs) and recurrent neural networks (RNNs) for real-time gesture recognition.
3. Dhruv, A., & Bharti, S. (2021). Real-Time Sign Language Converter for Mute and Deaf People. 2021 International Conference on Artificial Intelligence and Machine Vision (AIMV), 1-6. <https://doi.org/10.1109/aimv53313.2021.9670928>.
4. Vu, H., Hoang, T., Tran, C., & Pham, C. (2023). Sign Language Recognition with Self-Learning Fusion Model. IEEE Sensors Journal, 23, 27828-27840. <https://doi.org/10.1109/JSEN.2023.3314728>.
5. K, R., A, P., S, P., Sasikala, S., & Arunkumar, S. (2022). Two-Way Sign Language Conversion for Assisting Deaf-Mutes Using Neural Network. 2022 8th International Conference on Advanced Computing and Communication Systems (ICACCS), 1, 642-646. <https://doi.org/10.1109/ICACCS54159.2022.9785057>.
6. Dhamane, S., Ainapure, A., & Dhage, S. (2023). Breaking Silence, Embracing Inclusivity: The Power of AI in Communication. 2023 IEEE 5th International Conference on Cybernetics, Cognition and Machine Learning Applications (ICCCMLA), 231-236. <https://doi.org/10.1109/ICCCMLA58983.2023.10346765>.

PranavNamt PatniGoyal

sign2text_Prnav_Namit.docx

 Birla Institute of Technology, Mesra

Document Details

Submission ID

trn:oid::3117:411289959

Submission Date

Dec 2, 2024, 1:36 PM GMT+5:30

Download Date

Dec 2, 2024, 1:38 PM GMT+5:30

File Name

sign2text_Prnav_Namit.docx

File Size

1.0 MB

45 Pages





9,702 Words

55,064 Characters




19% Overall Similarity

The combined total of all matches, including overlapping sources, for each database.

Match Groups

-  **123** Not Cited or Quoted 19%
Matches with neither in-text citation nor quotation marks
-  **0** Missing Quotations 0%
Matches that are still very similar to source material
-  **1** Missing Citation 0%
Matches that have quotation marks, but no in-text citation
-  **0** Cited and Quoted 0%
Matches with in-text citation present, but no quotation marks

Top Sources

- 16%  Internet sources
- 8%  Publications
- 0%  Submitted works (Student Papers)

Integrity Flags





0 Integrity Flags for Review

No suspicious text manipulations found.




Our system's algorithms look deeply at a document for any inconsistencies that would set it apart from a normal submission. If we notice something strange, we flag it for you to review.

A Flag is not necessarily an indicator of a problem. However, we'd recommend you focus your attention there for further review.

Match Groups

-  **123** Not Cited or Quoted 19%
Matches with neither in-text citation nor quotation marks
-  **0** Missing Quotations 0%
Matches that are still very similar to source material
-  **1** Missing Citation 0%
Matches that have quotation marks, but no in-text citation
-  **0** Cited and Quoted 0%
Matches with in-text citation present, but no quotation marks

Top Sources

- 16%  Internet sources
- 8%  Publications
- 0%  Submitted works (Student Papers)

Top Sources

The sources with the highest number of matches within the submission. Overlapping sources will not be displayed.

1	Internet	
www.ijert.org		5%
2	Internet	
cse.anits.edu.in		4%
3	Internet	
www.ijraset.com		1%
4	Internet	
www.slideshare.net		1%
5	Publication	
"Applications of Artificial Intelligence in Engineering", Springer Science and Busin...		0%
6	Publication	
"Smart Data Intelligence", Springer Science and Business Media LLC, 2024		0%
7	Publication	
O. P. Verma, Seema Verma, Thinagaran Perumal. "Advancement of Intelligent Co...		0%
8	Publication	
Melvin Paul Miki.V, S. Prakash, Amarnath. M.K.V, K.Naveen Kumar. "Arduino cont...		0%
9	Publication	
Shymaa G. Eladl, Amira Y. Haikal, Mahmoud M. Saafan, Hanaa Y. ZainEldin. "A pro...		0%
10	Publication	
Rodrigues, Duarte Barros. "Hospitality AI-Driven Customer Journey Analytics: Pre...		0%