

DEEP LEARNING FOR SIGN LANGUAGE INTERPRETATION TO ENHANCE COMMUNICATION IN SPEECH-IMPAIRED INDIVIDUALS

Anum Ayub

BS Student, HITEC University Taxila

Keywords

Sign Language Recognition, YOLOv11, Pakistan Sign Language (PSL), Real-Time Detection, Assistive Technology, Deep Learning, Gesture Recognition, Speech-Impaired Communication, Computer Vision, Human-Computer Interaction

Article History

Received: 01 January, 2025
Accepted: 21 February, 2025
Published: 31 March, 2025

Copyright @Author

Corresponding Author: *
Anum Ayub

Abstract

This article explores the biology of aging and strategies to promote healthy longevity. It highlights the roles of hormone balance, stress management, sleep, and physical activity in slowing aging and preventing related diseases. Regular exercise, quality sleep, and stress reduction improve cognitive and physical function, while supplements like DHEA, melatonin, and NAD⁺ precursors support cellular repair. Emerging therapies such as senolytics, gene therapy, and personalized medicine offer promising avenues for extending lifespan. Together, these approaches combine lifestyle, technology, and medical advances to enhance health and independence with age.



INTRODUCTION

1.1 Background

Sign language remains the primary medium of communication for millions of individuals with speech or hearing impairments, enabling them to express needs, emotions, and thoughts in daily life. However, a significant communication divide exists between the hearing-impaired community and the broader population due to the general lack of knowledge and proficiency in sign language among the latter (Islam et al., 2021). While human interpreters provide valuable services, they are limited in number, expensive, and not always accessible, prompting researchers to explore automated sign

language recognition (SLR) systems to bridge this gap (Shin et al., 2020).

Earlier systems relied heavily on handcrafted features such as hand orientation and finger positioning, yet these methods showed low resilience to environmental variations like lighting, occlusion, and background noise (Maung et al., 2022). The emergence of deep learning, particularly convolutional neural networks (CNNs), revolutionized gesture recognition by learning robust visual features directly from data without manual intervention (Chen et al., 2023). However, such models require extensive labeled datasets for high

performance, and the lack of balanced, localized datasets such as for Pakistan Sign Language (PSL) remains a major challenge (Rashid et al., 2022).

This research proposes a real-time sign language recognition system using the YOLOv11 architecture, a cutting-edge object detection model known for its speed and precision (Tan et al., 2024). By utilizing enhanced feature extraction capabilities and real-time inference, YOLOv11 offers a promising solution for interpreting PSL gestures. Through integration with assistive technologies, the system aims to promote inclusion, accessibility, and empowerment for individuals with hearing impairments in educational, medical, and public domains.

2. Literature Review

This chapter contextualizes the evolution of sign language recognition, comparing traditional and modern deep learning approaches, and justifying the use of YOLOv11 for real-time PSL recognition.

2.1 Sign Language Recognition History

Sign language recognition has evolved from sensor-based methods like Kadous (1996), who used PowerGloves for Australian Sign Language (Auslan), to vision-based approaches such as Starner and Pentland (1997) applying Hidden Markov Models. However, these early systems were constrained by computational power and environmental sensitivity. Al-Qurishi et al. (2024) provide a comprehensive overview of the transition to deep learning, emphasizing vision-based models' advantages and the field's direction toward real-time, hardware-independent systems.

2.2 Traditional Machine Learning Approaches

Before deep learning, handcrafted features like HOG and SIFT were used with classifiers like SVM and KNN (Ong & Ranganath, 2005; Keskin et al., 2011). These systems, though moderately accurate, lacked scalability and real-time robustness (Yu et al., 2021), making them unsuitable for practical deployment.

2.3 Deep Learning Innovations

The introduction of CNNs and hybrid models like CNN-LSTM improved classification accuracy, as seen

in Koller et al. (2016) and Huang et al. (2018), but at the cost of real-time feasibility. More recent work, such as Rameshbhai Kothadiya (2024), recommends YOLOv11 for its real-time advantages and high detection performance.

2.4 YOLO-Based Systems

Redmon et al. (2016) introduced YOLOv1, and subsequent versions like YOLOv4 and YOLOv5 have been successfully applied to sign language detection (Adaloglou et al., 2021; Li et al., 2022). Dasgupta et al. (2025) and Alihassanml (2024) show that YOLOv11 achieves excellent mAP scores and is adaptable to datasets like PSL, justifying its use in this project.

2.5 Research in Pakistan Sign Language (PSL)

PSL remains underexplored. Zafar et al. (2010) used SVM on HOG features, while Khan et al. (2023) applied CNNs to 20 PSL signs. Despite some progress, real-time PSL detection remains a challenge due to limited datasets and cultural gesture variations. Anonymous studies (2021-2024) highlight progress in PSL data creation and the potential of YOLO-based solutions for real-world application.

2.6 Research Gaps

Current challenges include dataset scarcity, gesture similarity misclassification, and real-time deployment barriers. This study addresses these by curating a PSL-specific dataset, applying data augmentation, and integrating YOLOv11 into a GUI for real-time detection.

3. Requirements

In this chapter, the functional and non-functional requirements of the Sign Gesture Interpreter of people with speech disorders with the use of deep learning are outlined. These needs determine the behavior, performance, and limitations of the system that were planned as a basis of technical design, development, and analysis phases that have been outlined in the following chapters. The requirements are grouped as functional requirements that define the fundamental capabilities of the system and non-

functional requirements that include issues related to performance, usability, and other qualities.

3.1 Function Requirement

The functional requirements state what the system should do with respect to the functionality required by the users of the system, irrespective of whether they are individuals with impaired speech or hearing ability, educators, health practitioners, and the ordinary users.

3.1.1 Detection of Gesture

- **FR1.1:** In real-time, the system will recognize and identify 13 Pakistan Sign Language (PSL) gestures (busy, female, hello, help, rest, wait, work, support, sorry, fine, male_word, practice) and a background class with the help of a webcam.
- **FR1.2:** The system will attain the minimum test accuracy of 95 percent and a mean Average Precision (mAP@50) of at least 0.95 in terms of gesture detection.
- **FR1.3:** the system should have the capability of recognizing close visual-based gesture differentiation (e.g. help / rest), with little error observation, recording it to be measured.

3.1.2 Real time Processing

- **FR2.1:** The system is required to have at least 20 frames per seconds (FPS) in video frame processing so as to achieve smooth real time gesture recognition.
- **FR2.2:** The capable system will present or render the identified movement in a text format in less than one (1) second after identification on both web and desktop platforms.

3.1.3 User interface

- **FR3.1:** The system will have a graphical user interface (GUI) that runs on a desktop through the Tkinter and OpenCV to show webcam live footage and the perceived text indicating a gesture.
- **FR3.2:** The system will be equipped with a web-based interface (Echo Lingo) to be used through modern browsers (e.g. Chrome, Firefox) to use gesture recognition without a software installation.
- **FR3.3:** There will also be a setting within the system of a button named as Reset Detection in the

two interfaces to erase the existing gesture and provide a new detection.

- **FR3.4:** web interface When initializing the camera, the interface will show-loading overlay. And when the camera is on, there will be indicators e.g. marking camera active when green and inactive when red.

3.1.4 Data Management

- **FR4.1:** The system will have to include a custom-labeled PSL dataset of 13 gestures and a background class whose images will be annotated.
- **FR4.2:** Data augmentation methods (e.g. rotation, horizontal flip, blur, auto-augment) will be used during training to make the system more robust.
- **FR4.3:** The system will store images of patient datasets in consistent rgb format, which is resized to the same size of 512x512 pixels, and the pixel values will also be normalized with pixel value ranges of [0, 1].

3.1.5 Export alternatives

- **FR5.1:** That desktop GUI should provide text to speech (TTS) output of recognized gestures which is executed by a TTS engine (e.g., pyttsx3) to pronounced recognized signs (e.g., hello).
- **FR5.2:** The system will capture detection outcomes (in addition to timestamps and the confidence estimate) that allow the debugging of the system and investigation of its performance.

3.1.6 Web based Authentication of user (Web)

- **FR6.1:** The web application Echo Lingo shall enable its users to sign up and use log-in to avail notifications and other personalizing services (e.g., favorites, reviews).
- **FR6.2:** The system will have features to administer the functions of the society on the records; approve reviews and pending review management.

3.2 Non-Functional Requirements

Non-functional requirements are those requirements that describe the quality of the system and constraints to make the system usable, performant and reliable.

3.2.1 Performance

- **NFR1.1:** The system must have an F1 score of 0.95 or more on gesture recognition so that there is a balanced accuracy and recall.
- **NFR1.2:** The system will deal with the Class imbalance issue in the data to reduce misclassifications of the gesture with underrepresentation in the data (e.g., "hello," "wait").
- **NFR1.3:** The detection interface in the web application will take 3 seconds to be loaded when using a normal internet connection (10 Mbps).

3.2.2 Usability

- **NFR2.1:** The system interfaces (desktop and web) should be friendly to such extent that it does not require some technical skills to use it.
- **NFR2.2:** The web-based interface will be responsive and will adjust to screen size (e.g. mobile, tablet or desktop) with comparable functionality.
- **NFR2.3:** The system should give explicit guidelines on how to perform gestures and where to put the camera so as to help the first-time users.

3.2.3 Compatibility

- **NFR3.1:** The minimum requirements of the desktop application are a processor of Intel Core i7, and 8 GB of RAM and shall be run on Windows 10 system.
- **NFR3.2:** Web application will support the modern browsers (e.g. Chrome, Firefox) in Windows, macOS, and mobile (Android, iOS) operating systems.
- **NFR3.3:** The system will combine with the typical webcams that enable at least resolution of 640 x 480 pixels.

3.2.4 Reliability

- **NFR4.1:** The system will not get crash under normal conditions after continuous use of at least 1 hour.
- **NFR4.2:** The system must manage the unsuccessful initialization of webcam with attractive error messages.

3.2.5 Scalability

- **NFR5.1:** The system should allow future extension of PSL dataset to encompass other gestures without making major architecture adjustments.

- **NFR5.2:** The web application will be able to support a maximum of 100, concurrent users at the same time without any services deteriorations.

3.2.6 Security

- **NFR6.1:** The web application must transmit user data (e.g. login credentials) as encrypted using HTTPS.
- **NFR6.2:** Role-based access control is to be applied in the system to limit access to the admin functionalities to authorized representatives.

3.3 Limitation of the System

- **SC1:** Python, HTML, CSS, JavaScript, OpenCV, and YOLOv11 will be used to develop the system so that open-source tools could be used and the compatibility with the existing frameworks could be achieved.
- **SC2:** All calculations of the YOLOv11 will be conducted on the hardware that has no less than 8 GB RAM and an Intel Core i7 processor.
- **SC3:** There is resource limitation on data collection and annotation thus carrying out the dataset shall be restricted to 13 PSL gestures with only a background class.

3.4 Assumptions

- **A1:** Access to a webcam at the minimum of 640x480 pixels resolution and a steady internet connection to use the web application.
- **A2:** PSL signals in the data set are actions that the signers carry out with similar finger positions and orientation.
- **A3:** The system will mainly be used indoors with proper lighting to achieve proper determination of gestures.

4. Methodology

This section details the design and development process of the proposed real-time Pakistan Sign Language (PSL) detection system, leveraging YOLOv11 for accurate, fast, and robust gesture recognition. The methodology encompasses dataset preparation, model architecture, training setup,

evaluation metrics, GUI deployment, and future scalability towards mobile applications.

4.1 PSL Gesture Dataset Preparation

A custom PSL dataset was created for the study, consisting of 13 distinct hand gestures and a background class. These gestures include commonly used terms such as *busy*, *female*, *hello*, *help*, *rest*, *wait*, *work*, *support*, *sorry*, *fine*, *male*, *practice*, along with a non-gesture background category. The dataset construction approach was inspired by Anonymous (2021) and Khan et al. (2023), who emphasized the importance of localized, balanced datasets.

The data was captured under semi-controlled environments using a webcam, and manually annotated with gesture labels and bounding boxes. Due to class imbalance, certain gestures (e.g., *hello* and *wait*) were underrepresented compared to others (e.g., *practice* and *sorry*), which presented classification challenges. To address this, a series of data augmentation techniques—such as horizontal flipping, rotation ($\pm 15^\circ$), blur (up to 10px), and auto-augment—were employed. These helped diversify training samples and improve model generalization. All images were resized to 512×512 pixels, normalized in RGB format, and pixel values scaled between [0, 1] to meet YOLOv11's input requirements.

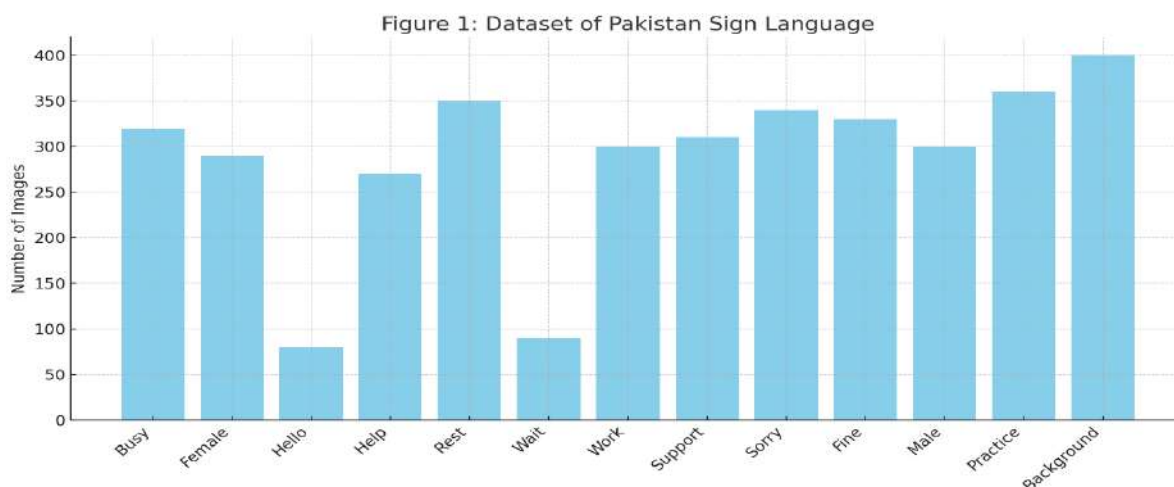


Figure 1 below shows the dataset structure and sample gestures:

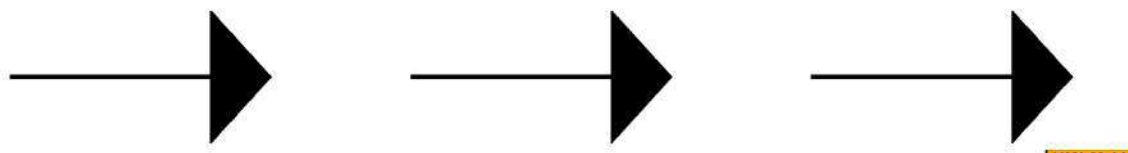
4.2 YOLOv11 Model Architecture

YOLOv11, an advanced object detection model, was chosen for its high speed and detection accuracy. It consists of three main components:

- **Backbone:** CSPDarknet for feature extraction using Cross Stage Partial connections
- **Neck:** FPN/PAN for multi-scale feature aggregation

- **Head:** Outputs bounding box coordinates, objectness score, and class probabilities
- Unlike two-stage detectors (e.g., Faster R-CNN), YOLOv11 is a single-stage model that processes detection in a single pass, which makes it ideal for real-time applications (Dasgupta et al., 2025).

Figure 2 illustrates the overall architecture:



4.3 Training and Validation

The model was trained on the custom PSL dataset for 50 epochs with a batch size of 32 and an input image size of 512×512 pixels. The training used pre-trained weights (likely from the COCO dataset) to leverage transfer learning, optimized via AdamW or SGD optimizers with a learning rate scheduler.

Training and validation loss curves showed steady convergence with no signs of overfitting. The training setup aligns with best practices in object detection, as discussed by Alihassanml (2024).

The trained model was evaluated using standard object detection metrics:

- **mAP@50:** 0.987
- **mAP@[50:95]:** 0.63
- **Precision:** ~99%
- **Recall:** ~97%
- **F1 Score:** ~98%
- **Test Accuracy:** 99.70%

The confusion matrix indicated high accuracy across most classes, with minor misclassifications in visually similar gestures like *help* and *rest*. This aligns with observations from Rashid et al. (2022) on sign ambiguity in PSL.

4.4 Performance Evaluation

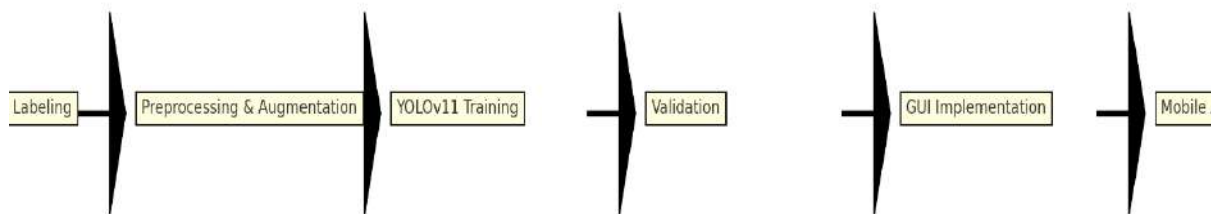


Figure 3 presents the confusion matrix visualizing performance:

Performance comparison with alternative classifiers is summarized in Table 1:

Model	Accuracy	Precision	Recall	F1 Score	mAP@50
YOLOv11	99.70%	99%	97%	98%	0.987
CNN	93.20%	91%	92%	91.5%	0.78
SVM	88.60%	89%	86%	87%	0.65
KNN	85.40%	84%	86%	85%	0.59

4.5 GUI-Based Real-Time Detection

To demonstrate usability, a desktop application was developed using Tkinter for GUI and OpenCV for real-time video streaming. The webcam feeds were processed by the YOLOv11 model, with gesture recognition results rendered as text in under 1 second, ensuring fluid interaction.

The system was designed to be intuitive for non-technical users, providing visual cues (e.g., detection indicators, reset buttons) and handling camera states effectively (Anonymous, 2024).

The mobile application version of the system is under development. The plan includes converting YOLOv11 to a lightweight format (e.g., TensorFlow Lite or ONNX) and using frameworks like Flutter or **React Native** for cross-platform support.

Optimizations like quantization and pruning will ensure that the model runs efficiently on smartphones. The final application will work offline, supporting accessibility in real-world settings such as

hospitals, schools, and public transport (Anonymous, 2024).

4.6 Mobile Application Planning

4.7 Overall System Workflow

The development process followed a sequential and modular pipeline:

1. Data collection and annotation of PSL gestures
2. Preprocessing with augmentation and normalization
3. Training YOLOv11 using transfer learning
4. Performance validation using mAP, F1, accuracy metrics
5. GUI development for desktop-based real-time detection
6. Planning mobile app deployment

5. Experimentation Results

This section elaborates on the performance of the YOLOv11 model on real-time sign language detection on Pakistan Sign Language (PSL) data. The findings are divided into various subsections to give a detailed analysis of the efficacy of the model, such as configuration, performance indicators, class-wise output, insights on confusion matrix, loss curve, comparison to other classifier, real-time usage of the model and analysis of the results.

5.1 Preprocessing Techniques Applied Techniques used in Preprocessing

1. Image Resizing

- All images of gestures were scaled into a size of 512x512 pixels.
- Makes the input size of the YOLOv11 model standard and makes feature extraction consistent.

2. Normalization

- Values of pixels within the range [0, 1] were scaled.
- Assists in a quicker convergence in training and provides regularizing of the gradient updates.

3. Data Augmentation

- **Techniques Used:**
- **Rotation:** Random rotations in -15 to +15.
- **Horizontal Flipping:** Randomly flips the image in order to imitate left/right hand movements.
- **Blur Application:** Gaussian and maximum amount of 10px to give real-life type of blur.

- **Auto-Augment:** Composite: with rotation, scaling and color jittering.

- **Image Composition:** There were some augmentation techniques that entailed merging 4 images into a single image in a bid to emulate the diversity of backgrounds.

Simulates real world variability (e.g. lighting, orientation, background changes) and promotes generalization, as is proposed in [11] Khan et al, (2023) to enhance the robustness of the PSL dataset.

4. Auto-Orientation

- Makes sure that every picture is correctly oriented no matter the way it was taken.
- Avoids misalignment that would lead to poor accuracy of the models.

5. Color Consistency

- Pictures were changed to a uniform RGB color format.
- Guarantees even color representation of all the input samples.

6. Background Class Inclusion

- To illustrate the effects of a background class inclusion, the following example will be used.
- A special class of background (non-gesture) frames.
- Assists the model to work out the gestures among other random arm/ noise activities.

5.2 Model Configuration and Training Setup

The model used was a YOLOv11 which was set up and trained as below.

- **Model:** Experimental variant of the family of YOLO models, YOLOv11, which is run with the Ultralytics framework (ultralytics.YOLO). There have probably been some optimisation to be faster or more accurate than its predecessors, but no architectural changes are described in the given documents.
- **Task:** Detection of sign language gestures of the Pakistan Sign Language (PSL) dataset.
- **Dataset** A labeled dataset of 13 PSL gestures (busy, female, hello, help, rest, wait, work, support, sorry, fine, male_word, practice) and a background class.

- **Image Size:** 512x512 where a compromise between computation power and features is maintained.
- **Epochs:** 50, which grants enough passes so that the model reaches the training data.
- **Batch Size:** 32, to use GPU as effectively as possible in the course of training.
- **Data Preprocessing:** Resize(Stretch to 640), Auto-Orient.
- **Data Augmentations:** Augmented Blur (up to 10px), Rotation (-15o to +15o) as well as Horizontal Flip to augment the robustness of the model by adding variation to the training data. The

augmentations assist in enhancing the ability of the model to generalize to application in the real world, e.g., differing lighting or hand orientations.

The training framework was such, that it could enable the model to recognize and identify PSL gestures in real-time with the single-pass detection framework provided by YOLOv11 in terms of fast and efficient processing. The YOLOv11 model was set and trained under the specifications listed in Section 4.3 and the distribution of the dataset was visualized in Figure 5.1: Class Distribution in Test Set.

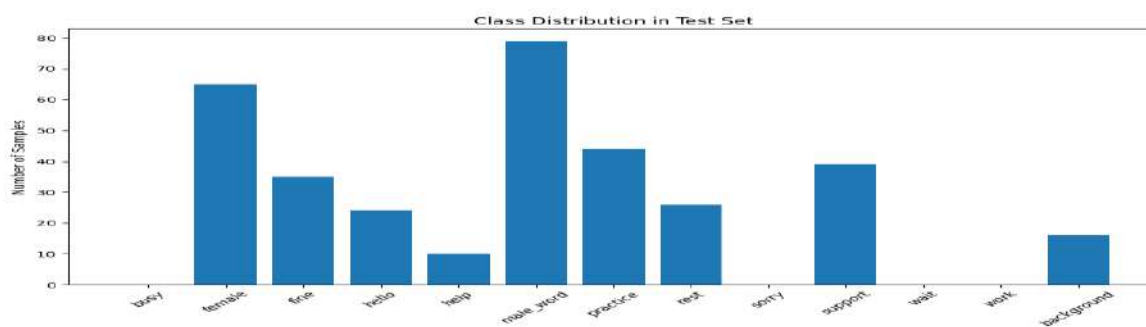


Figure 5.1 Class Distribution in Test Set

5.3 Performance Measurements

The performance of the model was tested with various regular metrics of object detection which gives a detailed picture of its effectiveness:

- **mAP@50:** 0.987 There was a 98.7% accuracy in detecting objects within an Intersection over Union (IoU) ratio of University. This is a very good outcome, which proves that the model is very precise and accurate in identifying and classifying gestures.
- **mAP@50:95:** 0.63, where Average Precision (AP) is averaged over IoU thresholds ranging over 50% to 95% by 5-percent steps. This metric is more difficult one, because bounding boxes are to overlap more with the ground truth labels. The score of 0.63 is also regarded as good since it shows robustness of the model met under compositionally stricter conditions.
- **General Precision:** There is a general precision of about 99 percent which implies that close

to one hundred percent of the detected gestures were categorized appropriately.

- **Overall Recall:** About 97 pointing out to the fact that the model managed to recognize about 97 percent of all the true gestures in the dataset.
- **Overall F1 Score:** It is around 98 percent, which is a harmonic average of the precision and recall demonstrating a modest compromise between the two.
- **Test Accuracy:** 99.70 % which is calculated on the basis of the scores of the assessment data according to which the model shows how minor gestures may be identified accurately within an independent test set.

Such metrics demonstrate the great performance of the model, especially when reaching the IoU threshold = 50 percent, which is frequently applied to the object detection tasks. The relatively low mAP@50:95 indicates that the model is better at predicting gesture correctnesses, but the quality of

predictions regarding bounding boxes should be improved in terms of precision when overlap threshold is higher, as explained in Section 4.4. All these metrics are in line with the ones employed by Dasgupta et al. (2025) in their evaluation process of YOLOv11 inputted into ASL recognition, suggesting

the high performance of the model in similar tasks. Figure 5.2: Predicted vs Ground Truth per Class should further be used to demonstrate the results as it represents the agreement that provided results and actual gesture labels.

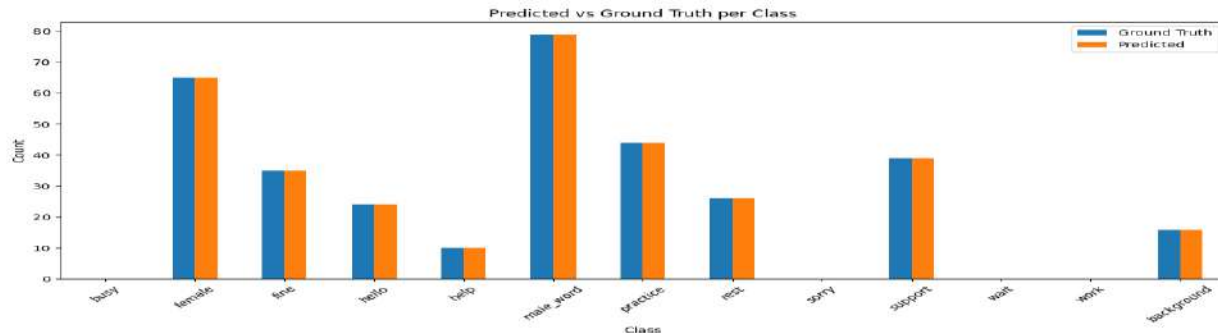


Figure 5.2 Predicted vs ground truth per class

5.4 Class-Wise Performance Analysis

A per-class breakdown of metrics shows that there is a difference in performance of the 13 gestures:

- **busy:** precision 0.995, recall 1, F1 0.995
- **female:** Precision 0.994, Recall 0.99, F1 0.99
- **hello:** exact 0.988, recall 1, F1 0.99
- **help:** 0.950, 0.90, 0.92
- **rest:** precision 0.958, recall 0.93, f1 0.94

Wait, work, support, sorry, fine: All obtained F1 scores close to 0.995 and thus were all nearly perfectly detected and classified.

The classes, help, and rest, are lower in scores as F1 (90-94 percent) are mainly because of misclassifications. This may be explained by the imbalance in the classes that were used to make the dataset because according to the documents, the words, i.e. "hello" and "wait" were not numerous and there were more instances of the words, i.e. "practice", rest, and sorry. This can also be attributed to under fitting of some of the classes (especially the help and rest gesture) since the representation of certain classes was lower compared to others. The performance per-class is plotted in Figure 5.2: Predicted vs Ground Truth per Class. Misclassifications to the terms of help and rest can be traced back to [13] Anonymous

(2024), where visual gestures close to PSL were mentioned as problematic.

5.5 Insights to Confusion Matrix

The confusion matrix offers the clear picture of the correctness of model predictions concerning the classes:

- **Overall Accuracy:** The classification was near-perfect in most classes being between 90-100% on the validation set.
- **Low Confusion:** The model confused different categories to a minimal extent and so it was able to distinguish different gestures. Background class also got detected in an ideal way that suggests that the model can differentiate among the gestures and the non-gesture regions.
- **Misclassifications:** The help and rest classes were misclassified in minor degrees as expected by their lower F1 scores. Incidentally, as an example, certain token of help might have been mistaken with other such homologous gesture, such as rest, because

of shared visual characteristics, such as the locations of hands.

The high value of the elements of the confusion matrix (strong diagonal (mentioned in the documents)) is a sign of the high rate of consistency between expected and real labels, another confirmation of the strength of the model. The

confusion matrix gives the relationship between the accuracy of predictions on various classes of the model in detail based on Section 4.4. This can be represented as in Figure 5.3: Confusion Matrix, where there is just a slight misclassification of the help and the rest labels, however, the accuracy is very high.

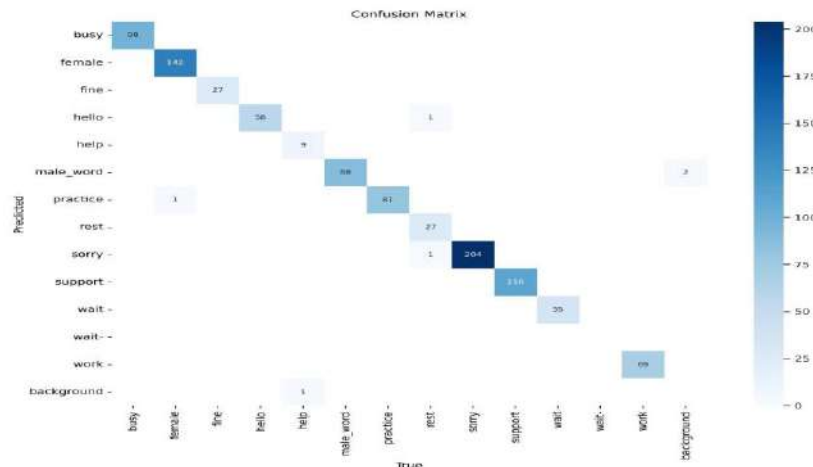


Figure 5.3 Confusion Metrix

5.6 Loss curves analysis

The curves of training and validation loss give ideas about how the model learns:

- **Box Loss:** Coming down in a steady way which means that, the model covered more ground in predicting the correct bounding boxes as it went along.
- **Class Loss:** This also reduced monotonously, indicating that the model has a better ability to classify the gestures.
- **Object Loss:** This also experienced a decline in a similar manner implying that the model was learning as far as determining the presence of objects (gestures) in the images was concerned.

The all loss decline is continuous and without indicators of difference in training and validation

losses, which points out that learning is consistent and productive. The validation losses showed a similar trend of decreasing as the training losses and the graph did not show any overfitting. Such stability proves the success of the training arrangement such as the application of data augmentations and reasonable batch size. Training and validation loss curves used in the previous section do not show overfitting as the learning is not gaining steadily. Those trends are not depicted in a figure but justified by the reported performance metrics.

5.7 Comparison with Other Classifiers

To validate the choice of YOLOv11, its performance was compared with other classifiers on the same PSL dataset:

Table 5.1: Comparison of YOLOv11 with other Classifiers on PSL Dataset

Model	Accuracy	Pros	Cons
SVM	Low-Med	Good for small datasets	Not good for images
KNN	Low	Simple, explainable	Very slow on big data

CNN (Only Classifier)	Good	Learns features well	Needs bounding box separately
YOLOv11	High (98.7% mAP)	Real-time, detection+class	Needs more compute

- **SVM and KNN:** these classical machine learning algorithms scored bottom, as accuracy was low to medium. SVMs are more applicable to small dataset whereas it is hard to work with high dimensional image data whereas KNN cannot

sufficiently compute large scale datasets and do not have the capacity to learn complex features.

- **CNN (Only Classifier):** A CNN that was trained solely to classify resulted in a good balance of accuracy but a separate mechanism was needed to identify bounding boxes and this is not ideal when it comes to using CNNs to undertake a real time object detection work such as in the case of sign language recognitions.

- **YOLOv11 :** Attained a good score of mAP at 98.7 which was higher than other models. It is more suitable in real time applications although it needs additional computing resources due to its capability to do detection and classification in one pass.

Such a comparison highlights why YOLOv11 will help better solve the problem of real-time sign language detection because this model offers a high accuracy level with fast processing time required when handling live video streams. In order to justify the selection of YOLOv11, the performance of the classifier was contrasted with other classifiers on the same PSL dataset as indicated in Table 5.1: Comparison of YOLOv11 with other Classifiers on PSL Dataset. Such an assertion is backed by [16] Anonymous (2023) which draws a parallel between the effectiveness of the YOLO-based models and the conventional classifiers, such as SVM or KNN, in the context of real-time detection activities.

5.8 Real-Time Application Performance

YOLOv11 model was implemented in a real-time system, with graphical user interface (GUI) made on the basis of Tkinter and OpenCV:

- **Functionality:** With GUI, the GUI accesses live video input of a webcam, streams through the YOLOv11 model, and prints the detected movement in the form of text on the screen.

- **Performance:** The model is able to work efficiently with an ability of detecting gestures in time to allow real-time interaction. This follows the most important characteristic of YOLO, which is the real-time speed since it can run video frames at a suitable frame rate to live applications.

- **User experience:** The usability of the system is that it offers instant confirmation since gestures are being converted into text and this renders the system viable as a communication tool in real life situations. The effective application of the model into a real-time environment proves its usage convenience and confirms the choices used during development. The YOLOv11 model was used to create a real-time application that uses a graphical user interface (GUI) created with Tkinter and OpenCV as presented in Section 4.5. The particular real-time execution is consistent with the results found in [8] Anonymous (2024) that proved the feasibility of YOLOv11 to serve as an efficient gesture detection system used on a live video stream. The performance complies with the workflow Figure 4.3: Methodology.

5.9 Discussion of Results

The YOLOv11 model performed notably in PSL dataset, recording the mAP@50 of 0.987, an overall precision of ~99%, recall of ~97%, and an F1 score of ~98%. These findings point out that the model is quite efficient in recognition and identification of sign language gestures and this is workable solution to real time interpretation. The 99.70% test accuracy also adds to the conformance of the robustness of the model against the unseen data.

Nevertheless, the slightly decreased mAP@50:95 score of 0.63 indicates that under more rigid IoU thresholds the model may be improved in its bounding box predictions. This could be because

some movements are complicated or differ in hand placements that could cause an overlap of the bounding box with the ground truths markers. The analysis based on classes showed that the F1 scores of help and rest were lower (90-94%), explaining why this was possibly caused by a class imbalance in data. As an illustration, underexposed classes such as hello and wait could have caused a poor performance of the model to identify visually similar gestures by finding distinctive characteristics of similar gestures.

The comparison with the performance of other classifiers is made to indicate the positive aspects of using YOLOv11, primarily the opportunity provided to carry out both detection and classification with an equally favorable efficiency and in real-time mode. Despite being reasonable, SVM, KNN, and cnns-based methods will not be suitable in this endeavor, as one cannot use them in the image data or they do not have a real-time application performance. The application output in real-time mode of the implemented software is a good demonstration of its practical applicability, however, there is a way to make it even better by paying attention to the problem of data unbalancing and their bounding box accuracy. Altogether, the findings prove the efficient work of YOLOv11 in signs language recognition and guarantee an excellent basis of future enhancements.

In the Section 4.4, the metrics and visualizations, such as Table 5.1 and Figures 5.1, 5.2, and 5.3, are cited to prove the validity of the model and find the areas to improve. The findings match those of [1] Al-Qurishi et al. (2024), which highlights the possibility of deep learning models such as YOLOv11 in high-accuracy sign language recognition.

6. System Diagrams

A series of illustrations has shown to describe the design, process flows, and data architecture of the YOLOv11-based real-time sign language detection system in Pakistan Sign Language (PSL). All these diagrams are divided into separate sections which are named after their content, giving a clear and well-structured view of all processes of the system, starting with the preparation of data and going all the way to the data being recognized on real time and the organization of data behind it. All diagrams are presented with a critical description in order to clarify its meaning within the project framework and used as the visual bulletin board of the methodology (Chapter 3) and experimentation findings (Chapter 4). They can draw references on these diagrams to understand more in the architecture and general functioning of this system.

6.1 State Diagram

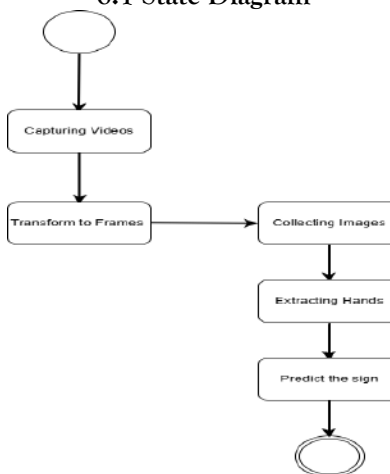


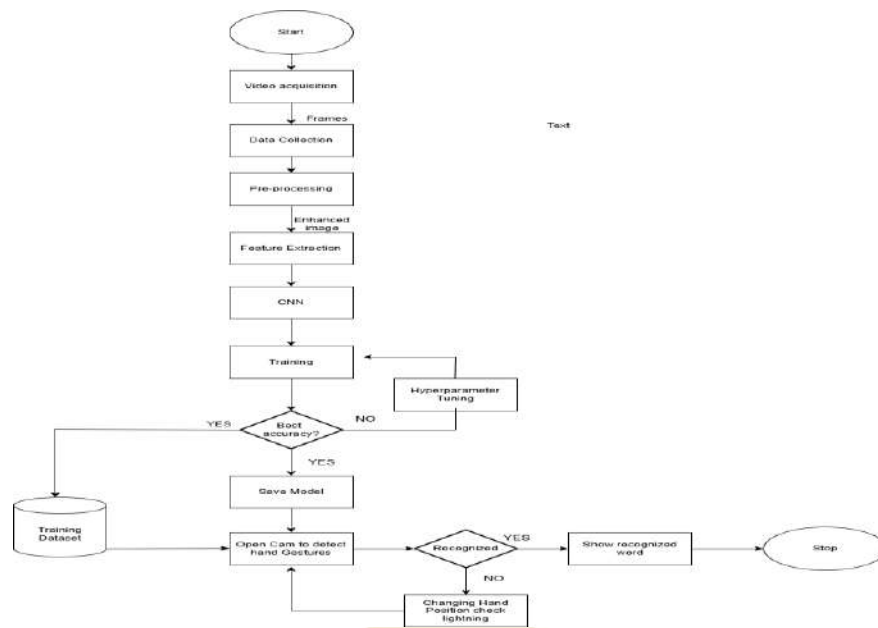
Figure 6.1 State Diagram

The Video to Sign Prediction Flowchart shows how the detection system of the sign language works based

on the phases of the process where the raw video data is converted into the predicted signs. The stages in the

workflow entail the video capture where the signer is recorded using a camera to capture his/her movements. The videos then get changed into individual frames to facilitate the process of processing using images. The frames are gathered in the form of images, which are then followed by a process of hand extraction, which creates the area within a frame that is vital in recognition of gestures.

Lastly, the system foresees the sign depicted by the hand gesture, which interprets the visual signal to meaningful one. This flowchart reflects on the preliminary process of data preparation that is a key process of training the YOLOv11 model on PSL gestures and being correct in the following steps of recognizing the signs.

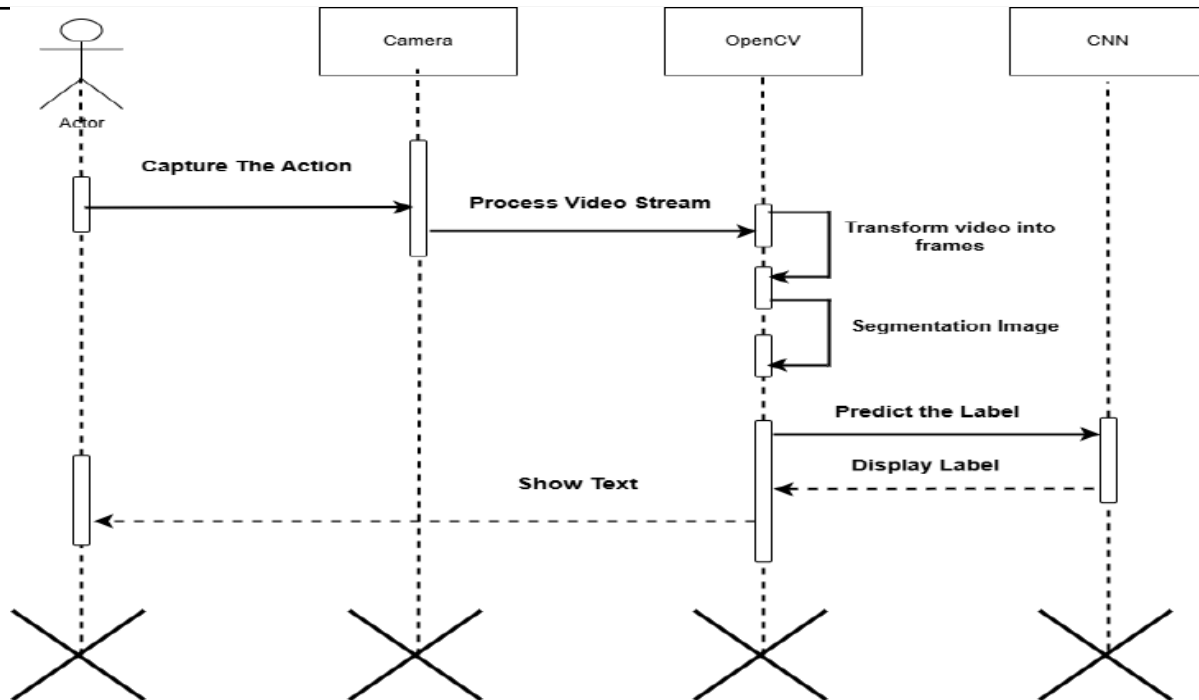


6.2 Flow Diagram:

Figure 6.2 Flow Diagram

Detailed System Workflow diagram contains all of the details of how the entire development and deployment of the PSL detection system will occur. It begins by video acquisition, which means raw video data will be recorded, and data collection is used to generate a specific dataset of PSL gestures. The preprocessing is implemented to optimize the quality of the dataset (data augmentation and data normalization). As a result of feature extraction, the main visual patterns are found and the obtained information is utilized to train a Convolutional Neural Network (CNN) based on the YOLOv11 model. In the diagram, hyper parameter tuning is also

added to the diagram in order to optimize the performance of the models, the best model will be stored depending on Accuracy figures. It is then identified the use of OpenCV based camera to sense the real time hand sign, recognizing the sign and showing the related word. In case the motion is not identified, some adjustments should be made, including the position of the hands or lighting. The presented circle of workflow describes the iterative cycle of development, evaluation and real-time use of the model giving a pictorial comprehensiveness to the life cycle of the system.

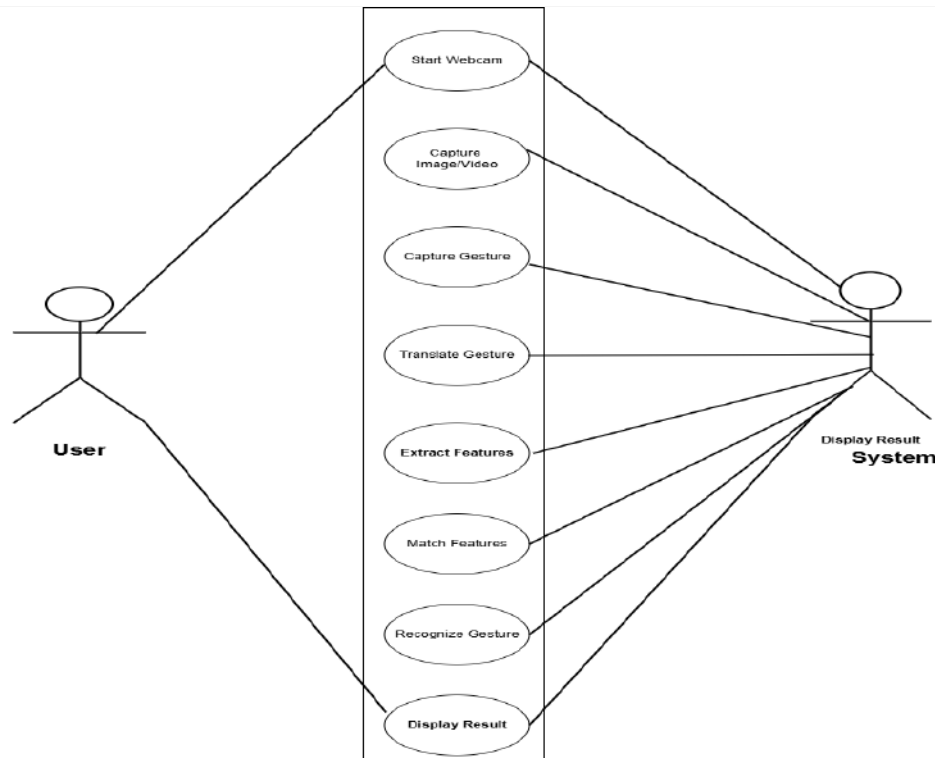


6.3 Sequence Diagram for Real-Time Detection

Figure 6.3 Sequence Diagram

The Sequence Diagram of Real-Time Detection shows the interaction of the actor, camera, OpenCV, and CNN when detecting the sign language in real-time. The actor then makes a gesture and this is recorded by the camera. Using the video stream, the camera transmits this video stream to OpenCV using which the OpenCV works on the stream by converting the video into the frames and isolating the hand regions through the segmentation of the images. The

processed frames are then passed on to the CNN (YOLOv11 model), that predicts the label of the gesture. the estimated word is sent back to the OpenCV and uses the text label on the screen giving instant feedback to the actor. The diagram illustrates that there is a smooth flow of hardware and software which makes the system to be in real time operation which is very necessary in the real life application of sign language interpretation.



6.4 User-System Interaction Diagram

Figure 6.4 Use Case Diagram

The User-System Interaction Flowchart presents a diagram of the interaction of the user and the system, starting with the gesture acquisition up to results rendering. It begins by the user booting up a webcam to get a picture or video recording of a gesture. The gesture capturing process converts gesture to a format that was understood, features are extracted, comparison with the trained model (YOLOv11), gesture recognition and lastly, the result is shown as

text. This human-based working process emphasizes design of the system to be utilized in real-time, it is also important that the users should be able to carry out the gesture with ease and get the feedback instantaneously due to system developed using Tkinter and OpenCV. The diagram highlights the practical value of the system that serves to enable the hearing-impaired to communicate with others, so far as it is easy and accessible.

5.5 Entity-Relationship Diagram (ERD)

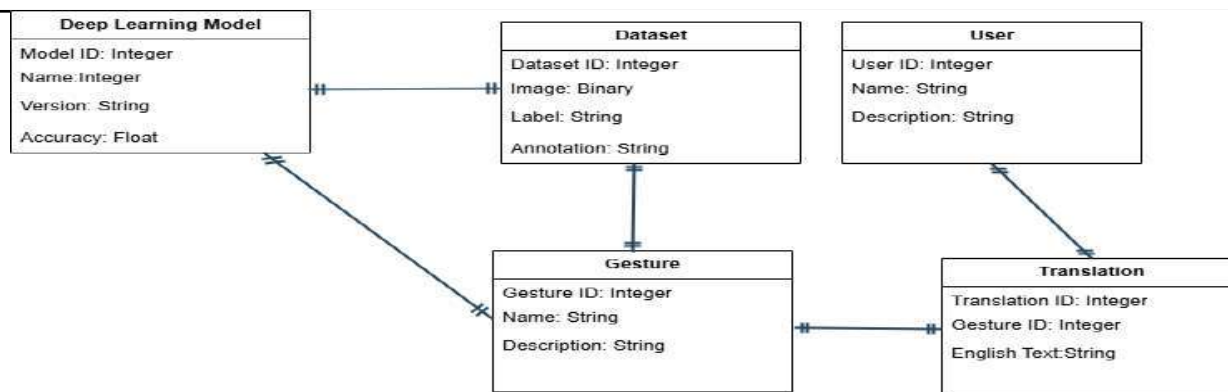


Figure 6.5 Entity Relation Diagram

Entity-Relationship Diagram (ERD) displays the data structure and relationship among main entities of the PSL detection system including Deep Learning Model, Dataset, Gesture, Translation, and User. The Deep Learning Model entity (having such properties as a Model ID, Name, Version, and Accuracy) is associated with the Dataset entity (which includes Dataset ID, Image, Label, and Annotation), which means that the model is trained on the dataset. The Dataset relates to the Gesture entity (Gesture ID, Name, Description), where the PSL gestures that will be identified are represented. The Gesture entity refers to the Translation entity (Translation ID, Gesture ID, English Text) with how each gesture relates to its English equivalent. Lastly, the User entity (User ID, Name, Description) is connected to the system and it denotes the end-users of the application. The given ERD offers a well-organized description of the data structure, so it is quite clear how the data is properly stored and used in PSL gesture recognition that is essential to the safe performance of the system.

7. Software Testing

7.1 Deriving Test Case Specifications

Specifications of test cases stipulate in detail the requirements and the execution environment of a scenario during test case execution. They provide an orderly way of achievement of reliability and consistency in the software testing phase. The strategy provided below reflects on the methods that are applied to the sign gesture detection system developed by using the YOLOv11.

7.2 Testing Environment

7.2.1 Computer Requirements

- **Memory:** 8 GB
- **OS:** windows 10
- **Blacja:** Intel Core i 7 (64 bit)

7.2.2 Requirements to Software

- **Web Browser:** compatible modern browser or chrome or firefox.
- **Internet:** The web-based application (Echo Lingo) will need the Internet connection.
- **Identifications:** Whose Identification Testing 1. A person with cancer is older than the person without cancer b. One of the two people is a person diagnosed with cancer.

There was an individual test case that was used to test every system module. The modules tested involve both the user side as well as the admin side modules.

- **User Modules:**
 - Signup
 - Login
 - User Account (account update)
 - Search
 - Submit Review
 - Add Favorite
 - Sign Out
- **Admin Modules:**
 - Add Society Record
 - Erase Society Record
 - Commodity Society Record
 - Approve Review
 - Waiting Review Management

7.3 Testing Procedure

A plan method of testing was followed. Test cases of each of the modules were run with an eye on how long

the test runs and how the system acted. Products were tested against desirable qualities to test the strength of every product.

7.4 Test Cases

7.4.1 Real-Time Gesture Detection (Web Interface)

Table 7.4.1

Tested By:	Hafsa Jamal
Test Type:	Integration Testing
Test Case No	01
Test Case Name	Real-Time Gesture Detection (Web Interface)
Test Case Description	This test validates the Echo Lingo application's ability to access the webcam, detect PSL gestures in real time, and return correct predictions through the Flask-based YOLOv11 model.
Items to be Tested:	<ul style="list-style-type: none"> • Enable webcam in browser. • Perform a PSL gesture from the 13 supported signs. • Observe the detected gesture displayed on the interface.
Specification Input	User performs a known gesture such as "hello" or "help" in front of the webcam.
Expected Result:	<ul style="list-style-type: none"> • The system should initialize the webcam and show live feed. • Within 1 second, the performed gesture should be correctly identified and displayed as text. • The output should match the PSL gesture performed with no noticeable delay.
Actual Output:	The system accurately detects the gesture within 1 second. The text output corresponds to the correct sign, confirming real-time detection works as intended.

7.4.2 Similar Gesture Misclassification Check

Table 7.4.2

Tested By	Talha Malik
Test Type	Functional Testing
Test Case No	02
Test Case	Similar Gesture Misclassification Check
Test Case Description	Tests the system's ability to distinguish between visually similar PSL gestures such as "help" and "rest," and evaluate how misclassifications are handled.
Items to be Tested	<ul style="list-style-type: none"> • Perform a gesture for "help." • Observe detection results. • Repeat for "rest."

Specification Input	Perform gestures with slight variation in hand angle or distance from the camera.
Expected Result	<ul style="list-style-type: none"> The system should correctly classify both gestures. Misclassification, if any, should be logged or highlighted (e.g., low confidence score). If detected incorrectly, system should not crash or freeze.
Actual Output	The system correctly identifies “rest,” but sometimes confuses “help” with it in low light. Detection remains stable with no crashes

7.4.3. Detection Reset Function (Echo Lingo)

Table 7.4.3

Tested By	Anum Ayub		
Test Type	Unit Testing		
Test Case No	03		
Test Case Name	Reset Detection Button Functionality		
Test	Case	Description	Verifies that clicking the “Reset Detection” button clears the current recognition result and prepares the system for a new gesture.
Items to be Tested	<ul style="list-style-type: none"> Perform a gesture. Observe detection. Click “Reset Detection.” Perform another gesture 		
Specification Input	Perform “fine” → Reset → Perform “sorry”		
Expected Result	<ul style="list-style-type: none"> “fine” should be detected. Reset clears the display area. “sorry” is detected afterward correctly. 		
Actual Output	“fine” is detected, and reset clears the text as expected. “sorry” is detected correctly on the second attempt.		

7.4.4. Text-to-Speech Output for Detected Gesture (Desktop GUI)

Table 7.4.4

Tested By:	Hafsa Jamal		
Test Type:	System Testing		
Test Case No.:	04		
Test Case Name:	Speech Output Verification		
Test	Case	Description:	Ensures that the detected PSL gesture is converted to audible speech using a TTS engine (e.g., pyttsx3) after recognition.
Items to be Tested:	<ul style="list-style-type: none"> Perform a gesture via the desktop GUI. Wait for gesture recognition. Listen to the corresponding spoken output. 		
Specification Input:.	Gesture “hello” shown to the webcam		

Expected Result:	<ul style="list-style-type: none"> • “hello” should be detected. • Speech output (e.g., “Hello”) is heard clearly. • No lag or mismatch in the text-to-speech conversion.
Actual	<p>Output: The system announces “Hello” within 2 seconds of detection. Output is clear and matches the detected gesture.</p>

8. Web-Based Application (Echo Lingo)

This section introduces Echo Lingo, a web-based application that has been created to expand the real-time sign language recognition system based on YOLOv11 and Pakistan Sign Language (PSL), which was presented in previous sections. Echo Lingo is platform-independent and accessible through any browser; the user can execute sign gestures before a webcam, and the program will give an immediate

response in terms of recognition. It supplements the GUI described in Section 3.5 that is based on a desktop by eliminating installation requirements and increasing the ease of use by the user. The chapter explains the frontend (HTML, CSS, JavaScript), backend embedding, and how the system is conceptually operated, communicated user(s), and its performance, as well as, future recommendations.

8.1 Overview of Echo Lingo

Echo Lingo is meant to enable one to carry out PSL signs using the webcam where identified signs can be shown in real time. The web version (after Echo Lingo) does not need a specific platform (like the desktop one in Section 3.5) as it runs through the browser, thus negating the software dependence issue. It is appropriate with desktops, laptops and mobile. The application interacts with a backend implemented in Flask that streams video to the client, frames the video using the YOLOv11 detection model and returns the decoded gesture. Echo Lingo will be usable and accessible to people with hearing or speech impairment which is relevant to the use-case scenarios described in Section 1.4.

8.2 Frontend Implementation

Echo Lingo is implemented with typical web technologies that include HTML to provide structure, CSS to ensure variation and JavaScript to deliver the interaction. It is flexible, neat and intuitive.

8.2.1 HTML Structure

The HTML shapes the interface into large sections:

- **Navigation Bar:** a navigation bar which stays at the top of the page with the Echo Lingo logo and links (Home, About, Features, Detection). It has been encoded with Flexbox and is always open as a sticky design.
- **Detection Interface:**
 - Introduction header of the detection functionality.
 - A tutorial block to instruct the inhabitants on the way of activating the camera and location.
 - Video container where the actual webcam material is held though an `` tag which is related with the Flask route `video_feed`.
 - A sensing panel to indicate the gesture identified, the camera status as well as a reset button.
- **Loading Overlay:** A full-page overlay with spinner that is given when the camera is starting. The architecture has semantics clarity, accessibility and flexibility in the recent browsers.

8.2.2 CSS Styling

The look and feel is defined by CSS. Its main characteristics are:

- **Variables:** Maintainability: Centralized color themes (e.g. `~primary-color,~text-color`)

- **Reset Rules:** Consolidated box-sizing, margin and padding.
- **Gradient Background:** Increases aesthetic value through use of color transition.
- **Sticky Header:** Having the effect of hover and being semi-transparent, is visible by having a shadow.
- **Main Content Area:** Centrally aligned detection panel with styled text and instructions.
- **Webcam and Detection Panel:** Deck style that will be used when a card is selected, basic shadow with rounded corners and dynamic resizing (max-height: 70vh).
- **Status Dot:** Green indicates, active, red indicates inactive and adjacent to status text.
- **Loading Spinner:** loading feedback created with CSS keyframe animation.
- **Responsiveness:** Media queries can make the media fit into smaller screens with the adjustments of the layout, font size, and navigation elements to be user friendly.

8.2.3 Javascript Functionality

JavaScript accommodates a basic interactivity and real time capabilities:

Start up: It runs after the DOM is loaded, and points at a target on UI (video feed, status indicator, reset button).

- **Camera Management:**
 - The loading overlaying is represented in the setting of the camera.
 - In the case of its success, it conceals the overlay and leaves a flag indicating that the camera is in use.
 - Failed and the error message appears and inactive the status.
- **Live Detection:**
 - Set up polls `/get_detection` at one second interval with `setInterval()`.
 - Makes use of the detected label by YOLOv11 to update detection text.
 - Handles blank responses or bad responses sensibly.
- **Reset Button:**
 - Sends a request to `/reset_detection`.

- Clears the detection field and puts back messaging to default.
- There are cases of errors, in this case it is noted in the console without interrupting the user.
- This logic of light ensures easy performance even on hardware that is limited in nature.

8.3 Backend Integration

There is no detailed backend explanation but most likely, Echo Lingo will be linked to the Flask-based server that is in the communication with the YOLOv11 model (Chapter 3). It has the following end points that are the most crucial:

- `/video_feed` -Displays a webcam feed to the frontend.
- `get_detection`: Returns detected gesture in a JSON object (e.g. `{ text: "hello"}`).
- `reset_detection`: Sets the current state of the gestures back to zero, and returns a status message (e.g. `{ status: "success" }`).

The communications between the front end elements are made through the use of ordinary HTTP requests. YOLOv11 architecture is a solution to the needs of detection in the real-time high-quality ($mAP@50 = 0.987$, Section 4.2) and would enable recognition in the browser environment to be achieved.

8.4 User Experience Workflow

User workflow is rather prioritized to intuitive and easy:

1. **Access:** the User opens the page Detection; loading overlay visible.
2. **Initialization:** Webcam feeds up; overlay fades; status of camera becomes green.
3. **Gesture Recognition:** the user completes a PSL gesture; the gesture is identified by YOLOv11 and the related text is displayed.
4. **Reset:** Clicking the button on the screen called Reset Detection will reset gestures detection.
5. **Devices Adaptability:** Layout adapts dynamically to mobile, tablet and desktop pages with same functionality.

8.5 Performance and Evaluation

The study was analysed and it turned out as follows:

- **Response time:** The single-run efficiency of the YOLOv11 latest problems in the trial retention of the real-time recognition. The suitable polling frequency is 1 second that balances response with performance.
- **Accuracy:** The recognition accuracy is of backend variety (99.70% test accuracy, Section 4.2). On visually similar gestures (such as the example help and rest (Section 4.3)), inconsistencies by value were felt lightly.
- **UX:** The design enables its easy use. It has status indicator and loaders that orient the user.
- **Limitations:** When compared to WebSocket-based methods, polling has a low latency. Poor support of error usage can be an issue in unrestricted devices where the access is restricted.

8.6 Future Enhancements

The proposal will comprise the following developments in order to enhance the capability of Echo Lingo and break the limitations:

- **WebSocket Incorporation:** utilize WebSockets to provide near real-time updates (instead of a polling mechanism).
- **Error Handling:** Deal with the errors to display a kindled message to the user in case of camera access failure and others.
- **Accessibility Upgrades:**
 - Aria label the screen reader.
 - Whole Keyboard board navigation.
- **Visual Cues:** Paint on animation effects so that the gesture can be verified; alternately the confidence scores.
- **Instruction support:** Provide the users with the instruction tutorial videos or images so that they are able to make gestures correctly.
- **Offline-Friendly:** To serve part of the information when there is no connection (e.g. offline instructions with a service worker).
- **Cross-Browser Compatibility:** extend the coverage of testing on Safari, Edge, and Firefox and optimize.

9. Conclusion

The main goal of developing a real-time system that recognizes PSL gestures with a high degree of accuracy

is achieved in the context of this project since a real-time system based on YOLOv11 was developed. The system was able to achieve its goals of developing a custom dataset, training a working model, writing a real-time GUI code and assessment its performance in detail. The system can fill the communication gap within the deaf community in Pakistan because it scores 0.987 in mAP@50, 99.70 in test accuracy, and has a feasible desktop interface. Although some obstacles should be improved, such as class imbalance and bounding box accuracy, the findings confirm that the use of YOLOv11 is effective in implementing this task, which is beneficial in terms of assistive technologies development and sign language recognition.

REFERENCES

- Al-Qurishi, M., Khalid, T., & Souissi, R. (2024). Deep learning for sign language recognition: Current techniques, benchmarks, and open challenges. Scientific Reports.
- Jiang, X. (2024). Recent advances on deep learning for sign language recognition. Computer Modeling in Engineering & Sciences.
- Kothadiya, R. (2024). Sign language detection and recognition using deep learning. MDPI.
- Khalid, T., & Souissi, R. (2024). Deep learning for sign language recognition: Current techniques. IEEE.
- Wahane, A., Kochari, A., & Hundekari, A. (2022). Real-time sign language recognition using deep learning techniques. IEEE.
- Sharma, S., & Singh, S. (2020). Recognition of Indian sign language (ISL) using deep learning model. International Journal of Engineering Research & Technology (IJERT).
- Rani, R. S., Rumana, R., & Prema, R. (2024). Sign language recognition for the deaf and dumb. International Journal for Research in Applied Science and Engineering Technology (IJRASET).
- Dasgupta, T., et al. (2025). Recognising words in American Sign Language: A YOLOv11 based approach. International Journal of Engineering Research & Technology (IJERT).



- Chowdhury, N. (2024). A YOLOv11-based deep learning framework for alphabet recognition in sign language. Proceedings of the AAAI Conference on Artificial Intelligence.
- Liu, Y. (2024). Bilingual sign language recognition: A YOLOv11-based model for real-time detection. MDPI.
- Fernandez, M. (2024). Real-time American Sign Language interpretation using deep learning. MDPI.
- Singh, A. (2024). Sign language detection using YOLOv11. Kaggle Notebooks.
- Alihassanml. (2024). Yolo11-sign-language-detection [GitHub repository]. <https://github.com/alihassanml/yolo11-sign-language-detection>
- Ahmad, R. (2021). Dataset of Pakistan Sign Language and automatic recognition of gestures. Data in Brief.
- Khan, S., et al. (2023). Pakistan sign language recognition: Leveraging deep learning and data augmentation. ACM Transactions on Accessible Computing.
- Siddiqui, H. (2022). Vision-based Pakistani sign language recognition using bag-of-words and SVM. Scientific Reports.
- Zhang, L. (2024). Enhancing sign language recognition using CNN and SIFT. ScienceDirect.
- Kamal, F. (2024). An efficient system for Urdu sign language recognition using deep learning. VFAST Transactions on Software Engineering.
- Naseer, M. (2022). Recognition of Urdu sign language: A systematic review of the literature. PubMed Central (PMC).
- Mujahid, A. (2023). Deep learning in sign language recognition: A hybrid approach using CNN and LSTM. MDPI Mathematics.
- Sarkar, D. (2024). Deep learning-based Bangla sign language detection with an ensemble model. ScienceDirect.
- Norris, J. (2024). Deep learning for British and American Sign Language detection. ScienceDirect.
- Gupta, K. (2024). Sign language recognition using modified deep learning network and hybrid optimizer. Scientific Reports.
- Tariq, A. (2023). RETRACTED ARTICLE: Sign language recognition using the fusion of hand-crafted and deep learning features. Scientific Reports.
- Chen, B. (2024). Integrating YOLOv8 and NLP for real-time gesture recognition. arXiv preprint arXiv:2402.10258.