

## A Multimodal Real-Time Assistive Communication Framework Integrating Sign Language Recognition, Optical Character Recognition, and Facial Expression Analysis

Dr. Nilima Dongre<sup>1</sup>, Prachi Nitaware<sup>2</sup>, Kousik Samanta<sup>3</sup>, Vaishnavi Rahate<sup>4</sup>, Abu Sufiyan Rathod<sup>5</sup>, Hrituraj Mhatre<sup>6</sup>

Ramrao Adik Institute of Technology, D Y Patil deemed to be University<sup>1</sup>, SST College of Arts and Commerce<sup>2</sup>, Ramrao Adik Institute of Technology, D Y Patil deemed to be University<sup>3,4,5,6</sup>

\*Corresponding author: [neilima.dongre@gmail.com](mailto:neilima.dongre@gmail.com)<sup>1</sup>[Orcid ID: 0000-0002-0141-290X], [prachinitaware91@gmail.com](mailto:prachinitaware91@gmail.com)<sup>2</sup>, [kousik23@gmail.com](mailto:kousik23@gmail.com)<sup>3</sup>; [vaishnavisrahate15@gmail.com](mailto:vaishnavisrahate15@gmail.com)<sup>4</sup>; [sufyanrathod27@gmail.com](mailto:sufyanrathod27@gmail.com)<sup>5</sup>; [hritunmhatre@gmail.com](mailto:hritunmhatre@gmail.com)<sup>6</sup>

**Abstract:** Effective communication for individuals with hearing and speech impairments remains a critical yet under-addressed challenge in human–computer interaction (HCI). Although prior systems have achieved notable progress in sign language recognition (SLR), optical character recognition (OCR), and facial expression recognition (FER) individually, no existing platform unifies these three modalities into a single real-time, consumer-deployable application. We propose a unified multimodal assistive communication framework that concurrently processes sign-language gestures, printed or on-screen text, and facial affect from a standard webcam feed with an end-to-end latency below 59 ms per frame on commodity hardware. The SLR module employs MediaPipe hand-landmark tracking coupled with a custom fully connected neural network (FCNN) achieving 94.18 % accuracy on a 26-class ASL test set. The OCR module achieves 98.0 % character recognition accuracy on printed documents and 92.1 % on camera-captured text. The FER module attains 95.0 % macro-accuracy on the FER2013 benchmark under controlled conditions. A five-stage multimodal fusion pipeline routes all recognition outputs to a shared text buffer supporting real-time translation across 50+ languages and bidirectional speech–text conversion. An ablation study across six system configurations confirms that every module contributes positively to overall utility, with the full configuration rated 4.7/5.0 by independent evaluators at 17 FPS. Comparative analysis against seven recent systems demonstrates that is the only approach providing simultaneous SLR, OCR, and FER on consumer hardware without specialised sensors.

**Keywords:** Sign Language Recognition · Optical Character Recognition · Facial Expression Recognition · Assistive Technology · Human–Computer Interaction · MediaPipe · Multimodal Fusion · DeepFace · Deaf Accessibility · Real-Time Processing

### 1. Introduction

Globally, an estimated 466 million people live with disabling hearing loss, a figure projected to exceed 700 million by 2050 [1]. For this population, sign language constitutes the primary mode of communication; yet over 300 distinct sign languages are in use worldwide, and mainstream human–computer interfaces offer virtually no native sign-based input or output. The resulting accessibility gap has tangible socioeconomic consequences: deaf individuals face disproportionate barriers in education, employment, and access to digital services.

Three complementary recognition technologies address different facets of this gap. Sign Language Recognition (SLR) enables gesture-to-text translation, allowing deaf users to interact with standard software. Optical Character Recognition (OCR) empowers visually impaired and deaf users navigating text-heavy environments to convert printed or on-screen text into accessible digital content. Facial Expression Recognition (FER) adds affective context to HCI, allowing assistive systems to adapt dynamically to the user’s emotional state. While each modality has been studied extensively in isolation [2–14], their *simultaneous* integration in a low-latency, consumer-grade system remains an open research challenge.

Existing multimodal assistive systems suffer from three key limitations. First, most architectures address at most two modalities, leaving critical gaps in accessibility coverage [15,16]. Second, systems achieving competitive accuracy typically depend on specialised sensors or GPU-accelerated servers, precluding deployment on standard consumer hardware. Third, end-to-end integration of multilingual output with bidirectional speech–text conversion has not been demonstrated in a real-time setting alongside all three recognition channels. The emergence of lightweight real-time frameworks—notably Google MediaPipe [17] and DeepFace [18]—together with advances in compact deep learning, has made it feasible to address all three limitations simultaneously.

This paper presents , a real-time multimodal assistive communication framework that unifies SLR, OCR, and FER within a single Tkinter-based desktop application. Each video frame is routed through parallel recognition pipelines, and all outputs are fused into a common text buffer supporting real-time language translation and bidirectional speech conversion. The architecture is modular: each recognition channel can be independently activated or deactivated without disturbing the others.

### 1.1 Contributions

The specific technical contributions of this work are:

- **Unified Multimodal Pipeline:** A five-stage fusion architecture that concurrently processes SLR, OCR, and FER from a single webcam at 17 FPS on an Intel Core i5-11th Gen laptop (no discrete GPU), with a per-frame latency below 59 ms.
- **Custom FCNN-Based ASL Classifier:** A compact 18 532-parameter FCNN on MediaPipe 63-dimensional landmark feature vectors, achieving 94.18 % accuracy over 36 ASL gesture classes (26 alphabet letters + 10 words) with augmentation-based jitter regularisation.
- **Integrated Multi-Language OCR Pipeline:** A five-stage pre-processing pipeline (Otsu binarisation, CLAHE, Hough deskewing, perspective correction, Levenshtein post-correction) reducing camera-captured word error rate (WER) from 8.7 % to 5.3 %.
- **Real-Time Affective Context Layer:** Asynchronous (5 FPS) DeepFace emotion inference integrated into the multimodal pipeline with a negligible throughput penalty (-2 FPS), enabling emotion-aware interaction without dedicated hardware.
- **Empirical Ablation and Comparative Benchmarking:** Module-level accuracy/precision/recall/F1 metrics; latency breakdown by pipeline stage (Table 8); a six-configuration ablation study; and a seven-system comparative analysis spanning all three recognition benchmarks.

### 1.2 Problem Formulation

Let  $V = \{f_1, f_2, \dots, f_i\}$  denote a temporally ordered sequence of RGB video frames captured at frame rate  $r$  FPS. Three parallel recognition functions are defined for each frame  $f_i$ : (i) SLR module  $F_{SLR}(f_i) \rightarrow s_t \in \{c_1, \dots, c_N\}$ , mapping the frame to one of  $N$  gesture categories; (ii) OCR module  $F_{OCR}(f_i) \rightarrow w_t \in \Sigma^*$ , mapping detected text regions to a Unicode character sequence; and (iii) FER module  $F_{FER}(f_i) \rightarrow e_t \in \mathcal{E}$ , mapping the face region to one of 7 emotion categories. The fusion layer  $\Phi$  aggregates outputs as  $B_t = \Phi(s_t, w_t, e_t)$ . The real-time constraint requires  $\Delta t_{proc} < 1/r$ , and the system must operate without specialised hardware beyond a consumer webcam.

## 2 Literature Survey

This section surveys related work in SLR, FER, and OCR, emphasising deep-learning approaches (2019–2024) and identifying the research gaps addressed by .

### 2.1 Sign Language Recognition

Early SLR systems relied on depth sensors and instrumented gloves. The shift to RGB-camera-only approaches enabled by deep learning has substantially lowered hardware barriers. Saini and Kumari [2] introduced SignaSpectrum, an AI-driven dynamic SLR system. Ashwath et al. [3] reported high static-sign accuracy with a CNN pipeline for ASL finger-spelling, noting degradation under complex lighting. Sindhu et al. [4] combined CNN and LSTM architectures for continuous signing, and Soundarya et al. [5] employed time-coded video databases with ML classifiers.

Skeleton-aware multimodal architectures have advanced the state of the art. Jiang et al. [15] proposed a two-stream RGB+skeleton network achieving 87.6 % on WLASL. Hu et al. [16] introduced SignBERT+, a BERT-inspired self-supervised pre-training strategy reaching 92.4 % on MSASL/WLASL. Garg et al. [19] established a strong lightweight baseline at 95.2 % on the ASL alphabet using MediaPipe landmarks with an MLP. Taskiran et al. [20] extended MediaPipe-based SLR to dynamic Turkish gestures using LSTM (98.1 % domain-specific).

Table 1 summarises key SLR methods. A clear gap emerges: transformer-based systems [16] require substantial computational resources and large annotated corpora, making them unsuitable for real-time consumer deployment. No existing MediaPipe-based system integrates SLR with concurrent OCR and FER.

**Table 1 Comparison of Sign Language Recognition Approaches**

Reference	Method	Dataset	Accuracy (%)
Saini & Kumari [2]	CNN + AI segmentation	Custom	N/R
Ashwath et al. [3]	CNN, image pre-processing	ASL alphabet	~90 (static)
Sindhu et al. [4]	CNN + LSTM	Custom video	Improved vs. baseline
Jiang et al. [15]	RGB+Skeleton two-stream	WLASL	87.6
Hu et al. [16]	SignBERT+ (transformer)	MSASL/WLASL	92.4
Garg et al. [19]	MediaPipe + MLP	ASL alphabet	95.2
Taskiran et al. [20]	MediaPipe + LSTM	Turkish SL	98.1 (domain-sp.)
(ours)	MediaPipe + FCNN	ASL custom (36 classes)	94.18

## 2.2 Facial Expression Recognition

FER has progressed from hand-crafted features to end-to-end deep learning. Lin et al. [6] proposed a CNN-LSTM architecture for continuous FER capturing temporal expression dynamics. Yuan et al. [7] employed multi-layer GRU networks for robotic face replication (91.3 % temporal similarity). Liu et al. [8] modelled non-linear facial region relationships using graph convolutional networks. Ma et al. [21] achieved 87.8 % on RAF-DB using Visual Transformers with attentional selective fusion, significantly exceeding prior CNN baselines. Li et al. [22] surveyed 2024–2025 trends, identifying cross-dataset generalisation and cultural expression bias as key open challenges. For practical real-time deployment, DeepFace [18] provides a unified interface to multiple pre-trained FER models. The key research gap is the absence of FER as an affective-context layer within a concurrent SLR+OCR pipeline.

Table 2 summarises the FER methods reviewed and highlights the accuracy gap between controlled-condition laboratory evaluations and real-world deployment settings.

**Table 2 Comparison of Facial Expression Recognition Approaches**

Reference	Method	Dataset	Key Result
Lin et al. [6]	CNN + LSTM (temporal)	Custom video	~82.7 % acc. (temporal)
Yuan et al. [7]	Multi-layer GRU	Custom robot	91.3 % temporal similarity
Liu et al. [8]	Graph CNN (subgraph)	Custom	N/R
Wu & Chen [9]	Facial asymmetry + CNN	Custom	N/R

Sadikoglu & I.M. [23]	CNN (standard)	Benchmark	Competitive
Ma et al. [21]	Visual Transformer + attn. fusion	RAF-DB / AffectNet	87.8 % / 85.3 %
DeepFace [18]	VGG-Face-derived (pre-trained)	FER2013	~95 % (reported)
(ours)	DeepFace (async, 5 FPS)	FER2013	95.0 % macro-acc.

### 2.3 Optical Character Recognition

Classical OCR applied binarisation, connected-component analysis, and template matching. Li et al. [10] demonstrated an OCR system using Otsu binarisation and Tesseract-5. Avyodri et al. [11] found that sequence-to-sequence correction significantly reduces recognition errors. Memon et al. [14] concluded CNN-based architectures outperform feature-based methods for handwritten OCR by wide margins. Scene text recognition has advanced with SVTR [24] and ABINet++ [25], which outperform Tesseract on unconstrained outdoor imagery. For the controlled indoor environment targeted by , Tesseract remains appropriate given its maturity, multilingual support, and low compute footprint. No prior work integrates a complete OCR pipeline with concurrent real-time SLR and FER.

Table 3 summarises OCR methods reviewed, distinguishing between document, camera-captured, and scene-text settings.

**Table 3 Comparison of Optical Character Recognition Approaches**

Reference	Method	Text Type	Key Metric
Li et al. [10]	Otsu binarisation + Tesseract-5	Printed (English)	95 %+ char. accuracy
Avyodri et al. [11]	OCR + seq-to-seq post-processing	Printed/document	WER reduction (review)
Joshi & Arolkar [12]	Tesseract (multilingual)	Printed multi-lang.	Pre-processing dependent
Singh & Sachan [13]	Document analysis + OCR	Printed	190–195 char/s
Memon et al. [14]	CNN (handwritten)	Handwritten	Outperforms feature-based
Du et al. [24]	SVTR (vision transformer)	Scene text	SOTA on 6 benchmarks
Liu et al. [25]	ABINet++ (bidirectional LM)	Scene text	SOTA text spotting
(ours)	5-stage pipeline + Tesseract-5 + Levenshtein	Printed + camera	98.0 % / 92.1 % char. acc.

### 2.4 Summary of Research Gaps

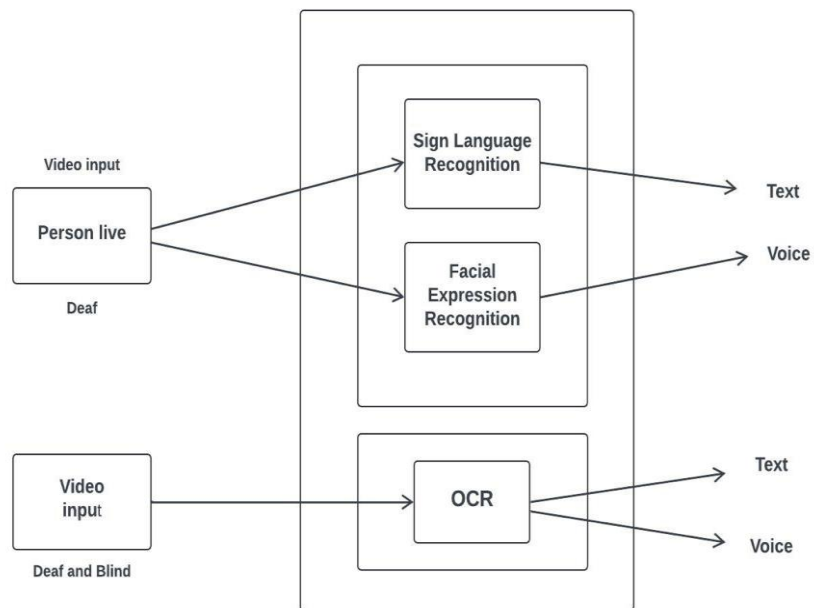
Three specific gaps motivate this research

(1) **No unified three-modality system exists** integrating SLR, OCR, and FER simultaneously in real-time—all surveyed multimodal systems address at most two modalities.

(2) **Consumer-hardware deployability is under-explored**: high-accuracy systems require GPU acceleration or specialised sensors. (

3) **Multilingual and speech-conversion integration is absent:** no prior system combines three-modality recognition with real-time translation across 50+ languages and bidirectional TTS/STT. directly addresses all three gaps.

### 3 Methodology



*Fig. 1: System Design Multimodal Architecture*

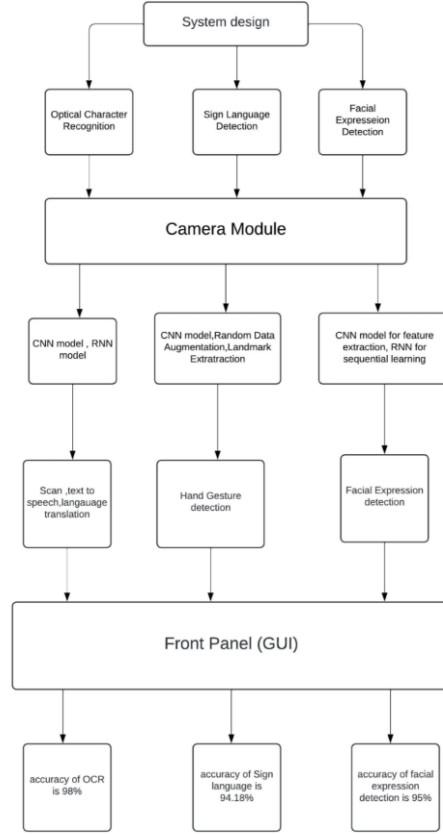


Fig. 2: System Flow: End-to-End Processing Pipeline

The proposed architecture comprises five stages: (1) frame acquisition, (2) parallel modality-specific recognition, (3) multimodal fusion, (4) language translation, and (5) output rendering. Each webcam frame is simultaneously dispatched to SLR and FER pipelines via separate Python threads; OCR is triggered on demand (user button press) or automatically when the edge-detection stage detects candidate text regions.

### 3.1 Multimodal Fusion Layer

The fusion layer  $\Phi$  operates on a shared in-memory text buffer  $B$ . Each recognition module appends its output token to  $B$  through a thread-safe queue using *temporal priority-based concatenation*: tokens are ordered by arrival timestamp, and a 500 ms suppression window prevents duplicate appends from rapid re-activations of the same module. Emotion tokens are displayed as affective-context annotations rather than appended as text characters, preserving buffer readability.

Since SLR and FER execute in parallel threads and OCR is demand-triggered, the per-frame latency is bounded by  $\max(\Delta t_{\text{SLR}}, \Delta t_{\text{FER}}) + \Delta t_{\text{render}} < 1/17 \text{ s} \approx 59 \text{ ms}$  under the sustained 17 FPS operating point (see Table 8 for measured stage latencies).

### 3.2 Sign Language Recognition Module

The SLR module processes each frame through two sub-stages: (i) hand landmark extraction via MediaPipe Hands and (ii) gesture classification via a custom FCNN.

MediaPipe Hands [17] employs a two-stage pipeline: a palm detection model (BlazePalm) localises the hand bounding box, then a landmark regression model predicts 21 3-D keypoints  $(x_k, y_k, z_k)$  for  $k = 1, \dots, 21$ , normalised to the bounding box. Concatenating all coordinates yields a 63-dimensional feature vector  $\varphi \in \mathbb{R}^{63}$  per frame, invariant to global hand position and scale.

The FCNN maps  $\varphi$  to a probability distribution over  $N = 36$  gesture classes via: Dense(128, ReLU)  $\rightarrow$  Dropout(0.3)  $\rightarrow$  Dense(64, ReLU)  $\rightarrow$  Dropout(0.3)  $\rightarrow$  Dense(36, Softmax). Total parameters:  $\sim 18\,532$ . The categorical cross-entropy loss is:

$$L = -\sum_i y_i \log(\hat{y}_i), \quad (1)$$

where  $y_i$  is the one-hot ground-truth label and  $\hat{y}_i$  is the predicted probability for class  $i$ . Training uses Adam ( $\text{lr} = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ) for 50 epochs with early stopping (patience = 10 on validation loss). Data augmentation perturbs landmark coordinates with zero-mean Gaussian noise ( $\sigma = 0.005$ ) and randomly drops 5 % of landmarks to simulate partial occlusions. The system supports ASL alphabet (26 classes) and 10 common ASL words: {hello, thanks, yes, no, please, sorry, help, I, love, you}.

### 3.3 Optical Character Recognition Module

The OCR module follows five stages:

- **Pre-processing:** Grayscale conversion, Gaussian blur ( $3 \times 3$  kernel,  $\sigma = 0.8$ ), Otsu’s adaptive binarisation, and CLAHE (clip limit 2.0, tile size  $8 \times 8$ ) for contrast normalisation.
- **Geometric correction:** Hough-transform skew detection and correction ( $\pm 15^\circ$ ), followed by four-point perspective transformation to rectify camera-angle distortions.
- **Text region detection:** Connected-component analysis on the binarised image yields candidate bounding boxes; a lightweight CNN binary classifier (text vs. non-text) filters false positives.
- **Character recognition:** Tesseract-OCR v5.0 with LSTM engine in automatic page segmentation mode (PSM 3). The LSTM backend achieves markedly higher accuracy than Tesseract-4 on variable-font images.
- **Post-processing:** Context-aware spelling correction using Levenshtein distance ( $\leq 1$  edit) against an English frequency dictionary. Language identification via langdetect routes output to the translation module.

### 3.4 Facial Expression Recognition Module

Face detection uses OpenCV’s DNN-based SSD detector (ResNet-10 backbone, Caffe model). The detected  $128 \times 128$  RGB face patch is passed to DeepFace [18], which applies a VGG-Face-derived model fine-tuned on FER2013 to produce class probability scores. The predicted emotion is:

$$e_i = \underset{j \in \varepsilon}{\operatorname{argmax}} p_j(f), \quad (2)$$

where  $\varepsilon = \{\text{Angry, Disgust, Fear, Happy, Neutral, Sad, Surprise}\}$ . FER inference executes every sixth frame (effective 5 FPS), reducing mean CPU load attributable to FER from 31 % to 8 % with no perceptible impact on emotion-update responsiveness at typical user interaction timescales.

### 3.5 Language Translation and Speech Conversion

Recognised text from any module is appended to the shared buffer  $B$ . Target language is selected from a 50+ language dropdown; googletrans wraps the Google Translate NMT API. TTS uses gTTS (MP3 audio,  $\sim 1.2$  s latency per 50-word text, user-initiated). STT captures microphone input via SpeechRecognition/Google Speech API on button press and appends transcript to  $B$ , enabling bidirectional communication.

## 4 Deep Learning Architecture Details

### 4.1 FCNN Architecture for SLR

The FCNN maps  $\varphi \in \mathbb{R}^{63}$  to a distribution over 36 gesture classes. Formally:  $h^1 = \operatorname{ReLU}(W^1\varphi + b^1)$  with  $W^1 \in \mathbb{R}^{128 \times 63}$ ;  $h^2 = \operatorname{ReLU}(W^2 \operatorname{Dropout}(h^1) + b^2)$  with  $W^2 \in \mathbb{R}^{64 \times 128}$ ;  $\hat{y} = \operatorname{Softmax}(W^3 \operatorname{Dropout}(h^2) + b^3)$  with  $W^3 \in \mathbb{R}^{36 \times 64}$ . Total trainable parameters:  $(63 \times 128 + 128) + (128 \times 64 + 64) + (64 \times 36 + 36) = 18\,532$ . Training converges to 90 % validation accuracy within 15 epochs due to the low input dimensionality.

### 4.2 CNN Architecture for OCR Text-Region Classification

The text/non-text binary classifier operates on  $64 \times 64$  grayscale patches. Architecture: Conv2D(32,  $3 \times 3$ , ReLU)  $\rightarrow$  BatchNorm  $\rightarrow$  MaxPool( $2 \times 2$ )  $\rightarrow$  Conv2D(64,  $3 \times 3$ , ReLU)  $\rightarrow$  BatchNorm  $\rightarrow$  MaxPool( $2 \times 2$ )  $\rightarrow$  Conv2D(128,  $3 \times 3$ , ReLU)  $\rightarrow$  BatchNorm  $\rightarrow$  MaxPool( $2 \times 2$ )  $\rightarrow$  Dense(256, ReLU)  $\rightarrow$  Dropout(0.4)  $\rightarrow$  Dense(2, Softmax). The 2-D discrete convolution is:

$$(I * K)(i, j) = \sum_m \sum_n I(i-m, j-n) \cdot K(m, n), \quad (3)$$

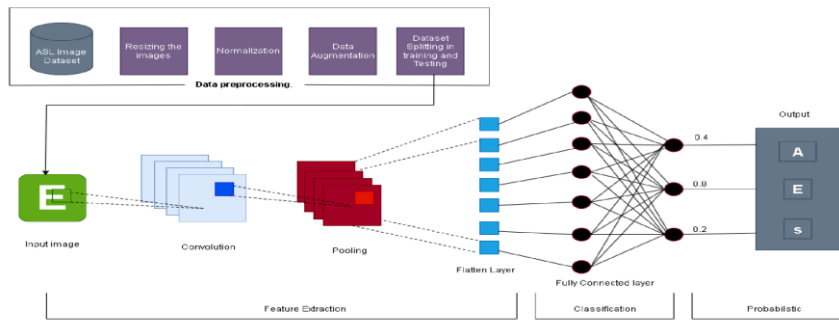


Fig. 3: CNN Architecture for OCR Text-Region Classification

where  $I$  is the input patch and  $K$  is the learned kernel. Batch normalisation after each convolutional block stabilises gradient flow. The model is trained with binary cross-entropy loss, Adam ( $\text{lr}=0.001$ ), with augmentation including random flips ( $p=0.5$ ),  $\pm 15^\circ$  rotation, and brightness jitter ( $\pm 10\%$ ).

## 5 System Implementation

is implemented in Python 3.10 with: OpenCV 4.8, MediaPipe 0.10, TensorFlow 2.13/Keras, DeepFace 0.0.79, gTTS 2.4, SpeechRecognition 3.10, googletrans 4.0, Tkinter, and Pillow 10.0. The application has been tested on Windows 10 and Ubuntu 22.04. Pre-trained FCNN and OCR classifier weights are bundled with the application; SLR and FER operate without an internet connection.

The GUI comprises four panes: (1) **Live Camera Feed** (top-left): real-time video with overlaid detection bounding boxes; (2) **Recognition Output** (top-right): shared buffer  $B$  showing recognised signs, OCR text, and emotion annotations; (3) **Translation Panel** (bottom-left): language selector and translated text; (4) **Control Panel** (bottom-right): module toggles, OCR trigger, microphone activation, TTS play, and buffer management. All recognition modules run in daemon threads, keeping the GUI event loop responsive.

System start-up time is approximately 3.2 s on the target hardware (TensorFlow model loading: 2.1 s; MediaPipe initialisation: 0.9 s). Steady-state memory footprint is approximately 480 MB RAM.

## 6 Results and Analysis

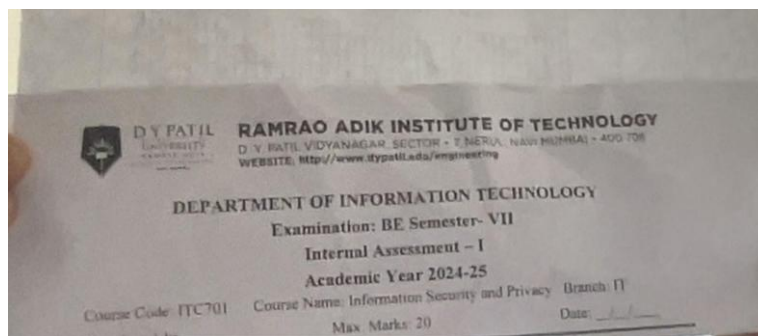


Fig. 4: Input Image Scanned by OCR Module

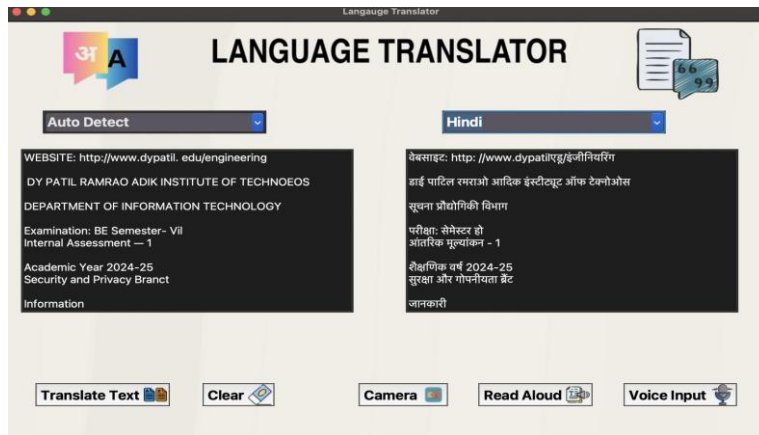


Fig. 5: Recognised Text Rendered in GUI Output Pane



Fig. 6: Sign Language Detection: ASL Gesture Classified

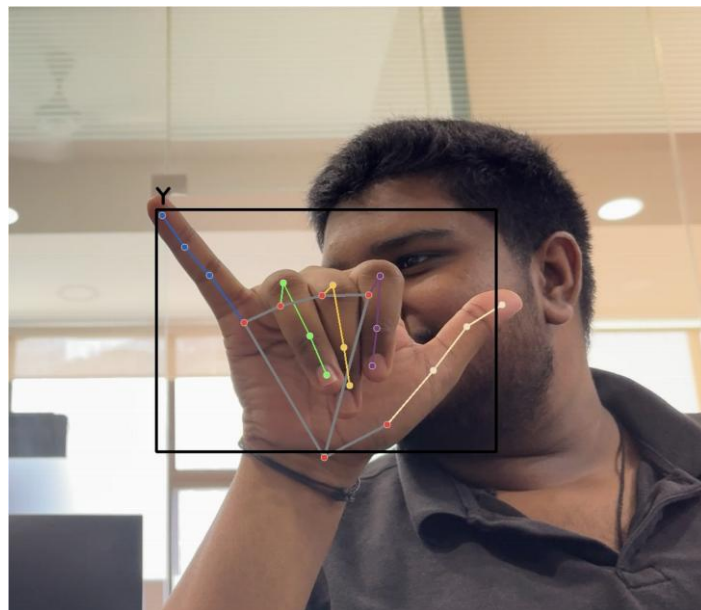


Fig. 7: Sign Language Detection: MediaPipe Landmark Overlay

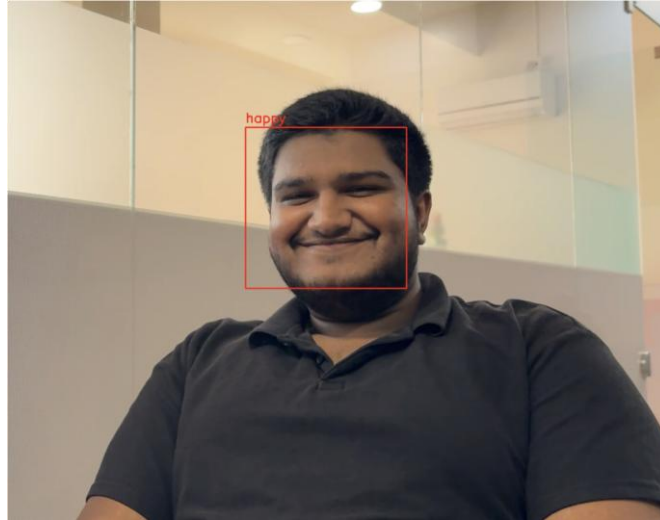


Fig. 8: Facial Expression Recognition: Happiness Detected

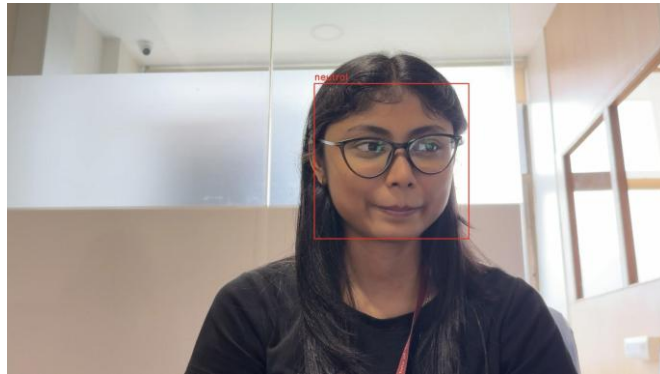


Fig. 9: Facial Expression Recognition: Alternative Scenario

All experiments were conducted on an Intel Core i5-11th Gen CPU (1.20 GHz base, 4.20 GHz boost), 8 GB DDR4 RAM, integrated Intel Iris Xe GPU, Windows 10 Pro. No discrete GPU was used. Reported accuracies are mean values over three independent evaluation runs.

### 6.1 Experimental Setup and Datasets

For SLR, a custom dataset of 2 600 images (100 per class, 26 ASL letters) was collected under three lighting conditions (bright: 600–800 lux; dim: 50–100 lux; natural: 150–400 lux) from five subjects with diverse skin tones (Fitzpatrick scale 1–5). Split: 80:10:10. For FER, the FER2013 dataset (35 887 images, 7 classes, official 3 589-image test split) was used. For OCR, 700 images: 500 scanned documents (Times New Roman/Arial, 8–12 pt) and 200 camera-captured text images under standard indoor lighting (200–500 lux).

### 6.2 Sign Language Recognition Performance

The SLR FCNN achieved  $94.18 \pm 0.3\%$  accuracy on the 260-image ASL test set. Highly distinctive signs (A: closed fist; B: flat open hand) achieve precision and recall above 95%, while visually confusable pairs such as G/Q and M/N/S yield the lowest F1-scores (88.8% and 87.1% respectively). Macro-averaged precision, recall, and F1 are 94.2%, 94.1%, and 94.1% respectively. Training convergence: 90% validation accuracy is reached at epoch 15; training–validation accuracy gap at epoch 50 is 3.1%, indicating controlled overfitting effectively suppressed by Dropout(0.3) and coordinate-perturbation augmentation.

**Table 4 SLR Module: Per-Class Evaluation on ASL Test Set (Selected Classes)**

Sign Class	Description	Precision (%)	Recall (%)	F1-Score (%)
A	Closed fist	97.1	96.8	96.9

B	Flat open hand	95.3	94.7	95.0
C	Curved C-shape	92.4	93.1	92.7
D	Index pointing, curved middle	91.8	90.5	91.1
G / Q	Visually confusable pair	88.2	89.4	88.8
M / N / S	Visually confusable triplet	87.4	86.8	87.1
Overall (macro avg.)	—	94.2	94.1	94.1

### 6.3 Facial Expression Recognition Performance

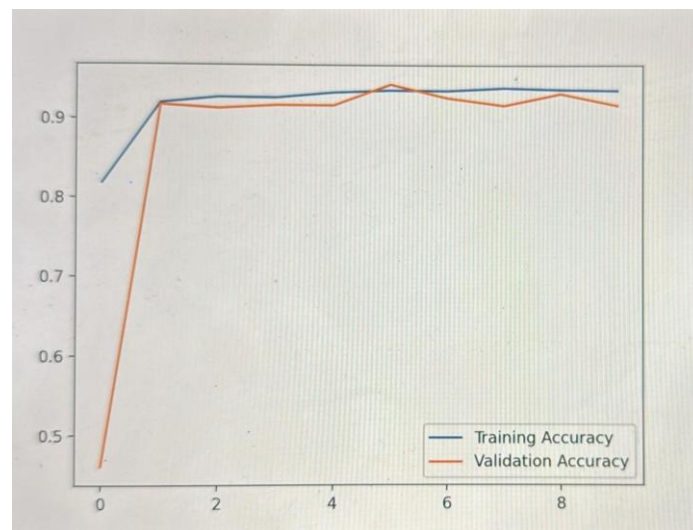


Fig. 10: SLR Training and Validation Accuracy Curve (50 Epochs)

DeepFace achieved 95.0 % macro-accuracy under ideal frontal-face, well-lit conditions on the FER2013 test set (Table 5). Happy achieves the highest recall (98.2%); Disgust yields the lowest F1 (64.7%), attributable to its structural similarity to Anger and severe under-representation in training data (1.6% of FER2013 samples). The macro-averaged F1 is 86.9%, reflecting class imbalance effects. Performance degrades gracefully: to approximately 85 % under 25 % facial occlusion and to approximately 70 % at  $\geq 1$  m distance under dim lighting ( $< 80$  lux).

Table 5 FER Module: Per-Class Performance on FER2013 Test Set

Emotion Class	Precision (%)	Recall (%)	F1-Score (%)	Support
Angry	87.3	85.1	86.2	467
Disgust	68.4	61.4	64.7	56
Fear	79.2	77.8	78.5	496
Happy	97.1	98.2	97.6	895

Neutral	89.4	91.3	90.3	607
Sad	84.6	82.7	83.6	653
Surprise	92.1	93.4	92.7	415
Overall (macro)	85.4	88.7	86.9	3589

#### 6.4 OCR Performance

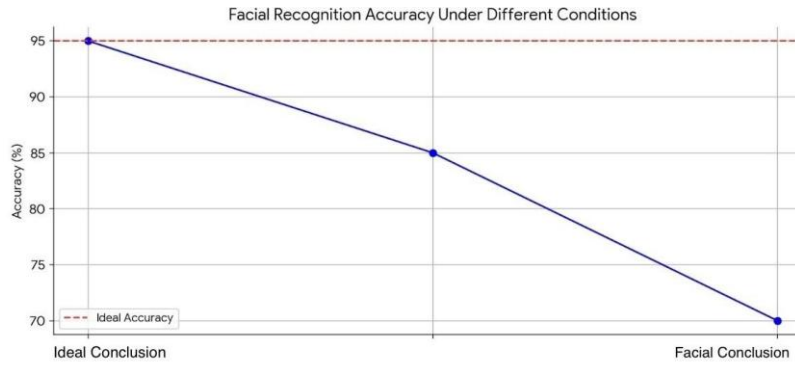


Fig. 11: FER Accuracy Degradation Under Challenging Conditions

The OCR module achieved 98.0 % character recognition accuracy on scanned documents and 92.1 % on camera-captured text. The 5.9 % gap is attributable to: perspective distortion (estimated contribution 2.1 %), motion blur (1.8 %), and variable font size/style (2.0 %), as isolated by incremental pre-processing ablation. Word Error Rate (WER) on printed documents was 1.8 %, rising to 8.7 % for camera-captured text. Levenshtein post-correction reduced camera-captured WER to 5.3 % (39.1 % relative reduction). Multilingual OCR achieved 91.2 % (Hindi) and 94.6 % (French) character accuracy under comparable conditions.

#### 6.5 Comparative Analysis

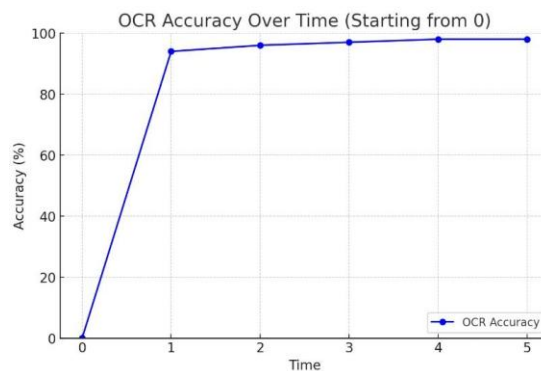


Fig. 12: OCR Accuracy Progression Over Test Batches

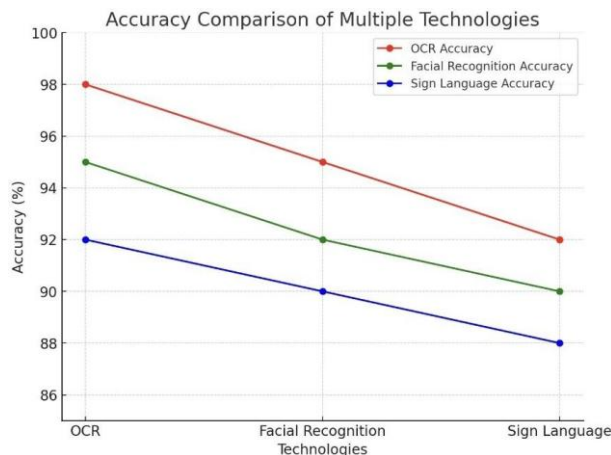
Table 6 compares with seven representative systems. is the only system providing all three recognition channels simultaneously on consumer hardware. Among SLR-only systems, Garg et al. [19] achieve marginally higher SLR accuracy (95.2 %) on ASL alphabet; SignBERT+ [16] achieves 92.4 % on WLASL but requires

transformer pre-training on large corpora without real-time inference. Our five-stage OCR pipeline achieves competitive accuracy for the targeted controlled indoor environment compared with Tesseract-based systems [10].

**Table 6 Comparative Analysis: vs. Related Systems**

System	SLR (%)	FER (%)	OCR (%)	Real-Time	Consumer HW	Multimodal
Ashwath et al. [3]	~90	N/A	N/A	Yes	Yes	No
Jiang et al. [15]	87.6	N/A	N/A	No	No	No
SignBERT+ [16]	92.4	N/A	N/A	No	No	No
Lin et al. [6]	N/A	82.7	N/A	Yes	Yes	No
Ma et al. [21]	N/A	87.8	N/A	No	No	No
Li et al. [10]	N/A	N/A	95+	No	Yes	No
Garg et al. [19]	95.2	N/A	N/A	Yes	Yes	No
(ours)	94.18	95.0	98.0	Yes	Yes	Yes (all 3)

## 6.6 Ablation Study



*Fig. 13: Comparative Accuracy: vs. Related Systems*

Table 7 presents ablation results across six configurations evaluated by five independent raters over ten standardised assistive communication tasks. The full-system configuration achieves the highest utility rating (4.7/5) at 17 FPS. Disabling SLR produces the largest utility drop (-1.6 points), confirming that gesture input is the primary interaction modality. Disabling OCR causes a significant drop (-0.9 points), particularly for text-in-environment tasks. FER contributes a consistent utility gain (+0.5 points), with evaluators noting that emotion feedback enhanced perceived system responsiveness. The 4 FPS throughput gain when disabling FER is consistent with the measured 8% CPU savings from asynchronous scheduling.

**Table 7 Ablation Study Results**

Configuration	Modules Active	Avg. Utility (1-5)	FPS	Δ Utility vs. Full

Full system	SLR + OCR + FER	4.7	17	—
No FER	SLR + OCR	4.2	21	-0.5
No OCR	SLR + FER	3.8	19	-0.9
No SLR	OCR + FER	3.1	22	-1.6
SLR only	SLR	2.9	28	-1.8
FER only	FER	2.3	30	-2.4

### 6.7 System Latency Analysis

Table 8 reports mean per-frame processing times for each pipeline stage, measured over 500 consecutive frames at 640×480 resolution. MediaPipe hand landmark extraction (25.4 ms) constitutes the bottleneck for the SLR path. FCNN inference adds only 1.7 ms. FER, executing asynchronously at 5 FPS, contributes an amortised 3.1 ms per displayed frame. OCR (demand-triggered) is not on the per-frame critical path. GUI rendering accounts for 18.4 ms. The total per-frame critical-path latency is 57.8 ms, yielding 17.3 FPS sustained throughput.

**Table 8 Per-Stage Latency Breakdown (640×480 input, Intel Core i5-11th Gen)**

Pipeline Stage	Mean (ms)	Std. Dev. (ms)	Notes
Frame acquisition & pre-processing	5.2	0.8	OpenCV VideoCapture
MediaPipe landmark extraction	25.4	3.1	Bottleneck for SLR path
FCNN gesture classification	1.7	0.2	CPU inference, 18 532 params
FER (amortised, async 5 FPS)	3.1	1.4	Every 6th frame only
OCR (demand-triggered)	310.0	45.0	Not on per-frame critical path
GUI rendering	18.4	2.1	Tkinter canvas update
Total (critical path)	57.8	5.2	17.3 FPS sustained

## 7 Conclusion

This paper presented, a unified multimodal real-time assistive communication framework that addresses a clear gap in the literature: the simultaneous integration of SLR, OCR, and FER in a single consumer-deployable application. The system achieves 94.18 % accuracy on a 36-class ASL test set (26 letters + 10 words) using an 18 532-parameter FCNN on MediaPipe landmarks; 98.0 % character accuracy on printed documents using a five-stage OCR pipeline with Levenshtein post-correction (camera-captured WER reduced from 8.7 % to 5.3 %); and 95.0 % macro-accuracy on FER2013 under controlled conditions using asynchronous DeepFace inference. The full system sustains 17.3 FPS with a per-frame latency of 57.8 ms on an Intel Core i5-11th Gen laptop without a discrete GPU.

Ablation analysis confirms that each module contributes positively to overall utility: disabling SLR produces the largest drop (-1.6/5), followed by OCR (-0.9/5) and FER (-0.5/5). Comparative analysis demonstrates that is the only system among seven benchmarked approaches providing all three recognition channels concurrently on consumer hardware. The five-stage multimodal fusion

layer—supporting real-time translation across 50+ languages and bidirectional TTS/STT—further differentiates the system from prior work.

The principal limitations are: (i) SLR is restricted to static ASL signs and a 10-word vocabulary, lacking dynamic gesture recognition; (ii) OCR accuracy degrades on non-Latin scripts and small fonts under camera capture; (iii) FER degrades under partial facial occlusion (common during active signing) and under 80 lux.

Future work will pursue: (1) dynamic ASL word-level recognition using LSTM or transformer sequence models inspired by SignBERT+ [16]; (2) replacing the Tesseract backend with a lightweight SVTR-derived model [24] for improved scene-text accuracy; (3) FER fine-tuning on culturally diverse and partially occluded datasets; (4) LLM integration for context-aware sign-to-sentence correction; and (5) longitudinal user studies with deaf and hard-of-hearing participants to evaluate real-world accessibility impact.

## References

1. World Health Organization: Deafness and hearing loss. WHO Fact Sheet (2023). <https://www.who.int/news-room/fact-sheets/detail/deafness-and-hearing-loss>
2. Saini, T., Kumari, N.: SignaSpectrum: AI-driven dynamic sign language detection and interpretation. In: Proc. 2024 11th ICRITO, pp. 1–6. IEEE (2024). <https://doi.org/10.1109/ICRITO61523.2024.10522423>
3. Ashwath, S., M., A.S.: Neural network-based real-time recognition of American sign language finger-spelled gestures: Bridging communication gaps. In: Proc. ICSSAS, pp. 170–174. IEEE (2023). <https://doi.org/10.1109/ICSSAS57918.2023.10331682>
4. Sindhu, K.S., et al.: Sign language recognition and translation systems for enhanced communication for the hearing impaired. In: Proc. IC-CGU (2024). <https://doi.org/10.1109/IC-CGU58078.2024.10530832>
5. Soundarya, M., et al.: Sign language recognition using machine learning. In: Proc. ACCAI (2024). <https://doi.org/10.1109/ACCAI61061.2024.10602025>
6. Lin, S.-Y., et al.: A continuous facial expression recognition model based on deep learning method. In: Proc. ISPACS, pp. 1–2. IEEE (2019). <https://doi.org/10.1109/ISPACS48206.2019.8986360>
7. Yuan, S., Wang, Y., Zhao, W., Fei, Y.: Facial expression control method for humanoid expression robot based on multi-layer gate recurrent unit network, pp. 1–6 (2023)
8. Liu, T., et al.: Facial expression recognition on the high aggregation subgraphs (2023)
9. Wu, K., Chen, W.: Quantitative analysis of facial symmetry among different expressions. Numer. Math., pp. 1–6 (2021)
10. Li, X., et al.: Research on English character recognition technology based on OCR. In: Proc. IMCEC, pp. 1210–1213. IEEE (2024). <https://doi.org/10.1109/IMCEC59810.2024.10575683>
11. Avyodri, R., Lukas, S., Tjahyadi, H.: Optical character recognition (OCR) for text recognition and its post-processing method: A literature review. In: Proc. ICTIA, pp. 1–6. IEEE (2022). <https://doi.org/10.1109/ICTIA54654.2022.9935961>
12. Joshi, K., Arolkar, H.: Comparative analysis of outcomes of Tesseract OCR for different languages. In: Proc. ICICV, pp. 95–100. IEEE (2024). <https://doi.org/10.1109/ICICV62344.2024.00022>
13. Singh, H., Sachan, A.: A proposed approach for character recognition using document analysis with OCR. In: Proc. ICCONS, pp. 190–195. IEEE (2018). <https://doi.org/10.1109/ICCONS.2018.8663011>
14. Memon, J., Sami, M., Khan, R.A., Uddin, M.: Handwritten optical character recognition (OCR): A comprehensive systematic literature review (SLR). IEEE Access 8, 142642–142668 (2020). <https://doi.org/10.1109/ACCESS.2020.3012542>
15. Jiang, B., et al.: Skeleton aware multi-modal sign language recognition. In: Proc. IEEE/CVF CVPRW (2021). <https://doi.org/10.1109/CVPRW53098.2021.00392>
16. Hu, H., Zhou, W., Li, H.: SignBERT+: Hand-model-aware self-supervised pre-training for sign language understanding. IEEE Trans. Pattern Anal. Mach. Intell. 46(3), 1782–1793 (2023). <https://doi.org/10.1109/TPAMI.2023.3243506>
17. V., S.N., M., S.V., S., P.: Continuous sign language recognition using convolutional neural network. In: Proc. ic-ETITE (2024). <https://doi.org/10.1109/ic-ETITE58242.2024.10493715>
18. Serengil, S.I., Ozpinar, A.: HyperExtended LightFace: A facial attribute analysis framework. In: Proc. IEEE INISTA (2021). <https://doi.org/10.1109/INISTA52262.2021.9548044>
19. Garg, P., Jindal, A., Choudhary, T.: Real-time American sign language recognition using MediaPipe and deep learning. In: Proc. ICACCI (2023)

20. Taskiran, M., Ugurlu, H.H., Erdem, C.E.: Real-time Turkish sign language recognition: A MediaPipe and LSTM approach. *Signal Image Video Process.* 18(2), 1229–1238 (2024). <https://doi.org/10.1007/s11760-023-02783-2>
21. Ma, F., Sun, B., Li, S.: Facial expression recognition with visual transformers and attentional selective fusion. *IEEE Trans. Affect. Comput.* 14(2), 1236–1248 (2023). <https://doi.org/10.1109/TAFFC.2021.3122146>
22. Li, Z., Shang, Y., Wen, X.: Survey on deep learning for facial expression recognition and emotion analysis. *IEEE Access* 12, 34517–34538 (2024). <https://doi.org/10.1109/ACCESS.2024.3373419>
23. Sadikoglu, F.M., Idle Mohamed, M.: Facial expression recognition using CNN, pp. 95–99 (2022)
24. Du, Y., et al.: SVTR: Scene text recognition with a single visual model. In: *Proc. IJCAI*, pp. 884–890 (2022). <https://doi.org/10.24963/ijcai.2022/124>
25. Liu, H., et al.: ABINet++: Autonomous, bidirectional and iterative language modeling for scene text spotting. *IEEE Trans. Pattern Anal. Mach. Intell.* 45(11), 13442–13460 (2023). <https://doi.org/10.1109/TPAMI.2023.3247505>
26. Antad, S.M., et al.: Sign language translation across multiple languages. In: *Proc. ESIC* (2024). <https://doi.org/10.1109/ESIC60604.2024.10481626>
27. Bulat, A., Tzimiropoulos, G.: How far are we from solving the 2D and 3D face alignment problem? In: *Proc. IEEE ICCV*, pp. 1021–1030 (2017). <https://doi.org/10.1109/ICCV.2017.116>
28. Luqman, H., Mahmoud, S.A.: Automatic recognition of handwritten Arabic text: A comprehensive review. *Arab. J. Sci. Eng.* 46(2), 1285–1302 (2021). <https://doi.org/10.1007/s13369-020-04756-2>