

GRADUATION PROJECT PROGRESS REPORT

A Real-Time, Dual-Function Sign Language Translation and Education System

Prepared by:

Salah Mohamed Salah Ibrahim Ahmed Adel Sayed Goda Ahmed
Alaa Osama Mohammed Shehab Hadeel Gamal Eldin Mohamed Ibrahim

Supervised by:

Dr. Nada Mostafa Abdelaziem Mostafa
Dr. Reham Taher Abdelmajeed Salem Al Magraby

December 16, 2025

Contents

1	Executive Summary	1
2	1. Introduction and Background	1
3	2. Project Goals and Objectives: Review of Progress	1
4	3. Phase 1 Progress: Model Development	2
4.1	3.1. Hardware and GPU Configuration	2
4.2	3.2. Architectural Model Development	2
5	4. Technical Results and Performance Metrics	2
6	5. Phase 2 Progress: Real-Time Integration	3
6.1	5.1. Data Curation and Pre-Processing	3
6.2	5.2. Real-Time Application Development	3
7	6. Challenges Encountered and Remedial Strategies	4
8	7. Forthcoming Stages of Development (Timeline Review)	4
9	8. Conclusion	5

1 Executive Summary

This report documents the developmental progress achieved during Term 9, confirming the successful conclusion of the research and initial proposal phase and the definitive transition to the system implementation stage. The foundational technical milestone—the establishment of a functional, real-time recognition engine for Arabic Sign Language (ArSL)—has been validated. This was achieved through the successful implementation and training of two distinct Artificial Intelligence paradigms: the image-based MobileNetV2 architecture and the geometry-centric MediaPipe model. The system currently demonstrates real-time translation capabilities, thereby fulfilling the core technical deliverable prerequisite for the subsequent phases involving application enclosure and user interface development.

2 1. Introduction and Background

The project aims to address significant communication barriers by developing a dual-function system that acts as both a real-time ArSL translator and an interactive educational platform. The initial stage of the project focused on selecting and validating the optimal computer vision and deep learning architectures capable of high-accuracy and low-latency performance on consumer-grade hardware. The algorithmic foundation of the recognition engine is predicated upon a dual-component strategy: the MobileNetV2 architecture, a highly optimized convolutional neural network employed for image-based classification via transfer learning, and the MediaPipe framework, utilized for the robust extraction of key skeletal features and three-dimensional hand landmarks for subsequent classification by a multilayer perceptron (MLP). The successful implementation of this dual-model approach provides the requisite technical redundancy to ensure deployment flexibility, regardless of the target machine's computational constraints.

3 2. Project Goals and Objectives: Review of Progress

The primary objectives outlined in the original proposal remain the guiding mandates of the project. A review of the status against the established phases is provided below:

Goal 1: Develop a High-Accuracy Recognition Model: Achieved via the successful training of both MobileNetV2 and MediaPipe MLP models, yielding validation accuracies exceeding 90%.

Goal 2: Implement Real-Time Inference: Achieved through the development of the `Combined_Architecture` script, which processes live video and provides continuous predicted text output.

Goal 3: Integrate Bidirectional Communication: Pending. Requires the integration of external Speech-to-Text and Text-to-Speech components (Stage II).

Goal 4: Develop an Educational Module: In progress. The foundational geometry comparison logic (Stage I) is being planned for implementation.

Goal 5: Deploy in a Dedicated Application: Pending. Requires the transition from the prototype environment (Jupyter/OpenCV) to a formal GUI wrapper (Stage III & IV).

4 3. Phase 1 Progress: Model Development

The entirety of the Model Development phase has been concluded, encompassing the establishment of the foundational infrastructure and the training of the primary recognition architectures.

4.1 3.1. Hardware and GPU Configuration

Formal documentation confirms the successful institution of the foundational technical ecosystem. The dedicated Graphics Processing Unit (GPU) was configured for optimal operation within the TensorFlow framework, necessitating specific kernel adjustments for stable performance under load. Critical to this stage was the execution of dynamic memory growth strategies, a necessary measure to rigorously preclude the manifestation of Out of Memory (OOM) exceptions during the intensive model training regimen, given the hardware's VRAM limitations. The core computational libraries (OpenCV, MediaPipe, Keras) have been integrated into a unified software architecture, guaranteeing system integrity for real-time operations.

4.2 3.2. Architectural Model Development

The strategic decision to pursue two disparate deep learning methodologies was executed to ensure technical robustness.

3.2.2. Approach A: MobileNetV2 (Transfer Learning Implementation)

The pre-trained MobileNetV2 architecture was employed as the foundational core, exploiting its deep feature-extraction capabilities. The stability of the network was augmented through the incorporation of BatchNormalization layers within the classification head, and the mandatory execution of the fine-tuning procedure on the proprietary, image-based ArSL dataset was completed. The resultant model weights are archived as `mobilenet_arabic_best_final`.

3.2.3. Approach B: MediaPipe and Multi-Layer Perceptron (MLP) Integration

The geometry-centric feature extraction pipeline utilizes the MediaPipe framework for the precise extraction of twenty-one distinct 3D hand landmark coordinates, successfully reducing the input dimensionality. This feature vector extraction precedes the training of a specialized lightweight MLP classifier, detailing the resultant low-latency performance profile that renders this model exceptionally appropriate for deployment within CPU-constrained environments.

5 4. Technical Results and Performance Metrics

A quantitative, empirical assessment of the performance efficacy exhibited by the two constructed models is provided, presented separately for each architectural methodology.

Table 1: MobileNetV2 (Image-Based) Performance Metrics Table 2: MediaPipe + MLP (Landmark-Based) Performance Metrics

Metric	Value	Metric	Value
Training Accuracy	95.2%	Training Accuracy	94.5%
Validation Accuracy	89.1%	Validation Accuracy	91.0%
Inference Speed	15 – 20 Frames per Second	Inference Speed	\geq 30 Frames per Second
Lighting Sensitivity	High	Lighting Sensitivity	Low

Detailed Performance Observation. A critical comparative analysis meticulously contrasts the superior proficiency of the MobileNetV2 architecture in the recognition of complex static imagery (signs involving subtle shading or internal hand texture) against the MediaPipe methodology's empirically verified computational efficiency (achieving frame rates above 30 FPS). This dual implementation provides the fundamental rationale supporting the strategic necessity of adopting a dual-model deployment paradigm, allowing the system to leverage the respective strengths of each model dynamically based on run-time conditions.

6 5. Phase 2 Progress: Real-Time Integration

The Real-Time Integration phase was initiated with the creation of the core communication and smoothing components, necessary for transitioning the trained models into a functional, user-facing application.

6.1 5.1. Data Curation and Pre-Processing

Quality assurance protocols were established for the ArSL dataset to guarantee semantic diversity across sign classes, encompassing variations in camera angle and background clutter. Furthermore, the dataset was augmented with specialized classes, including a NULL or 'NOTHING' class to suppress predictions when no sign is intended, and a dedicated SPACE class to enable segmentation between recognized letters for sentence construction.

6.2 5.2. Real-Time Application Development

The unified Combined_Architecture control script was developed to manage the live inference loop. This script facilitates concurrent video stream acquisition and implements the necessary graphic UI overlay for user feedback. Crucially, three technical components were integrated to ensure a robust user experience:

- **Temporal Smoothing Algorithm:** A confidence-weighted prediction history queue (`collections.deque`) was utilized to actively mitigate the output text instability, specifically the detrimental phenomenon of character flickering in the predictive display, by requiring multiple consecutive frames to concur on a single prediction.
- **Sentence Construction Logic:** The system incorporates logic to manage the predicted sequence of signs, appending classified signs to a sentence buffer and utilizing the designated SPACE class to insert word boundaries.

- **Hand Constraint Enforcement:** The current system architecture prioritizes the detection and classification of single-hand signs. Hand detection logic was implemented to monitor the presence of a second hand and output a visual warning or suppress the prediction if two hands are simultaneously detected, thus preventing misclassification in the current scope.

7 6. Challenges Encountered and Remedial Strategies

Formal documentation of significant technical obstacles encountered during the development cycle, alongside the precise, implemented strategies utilized for mitigation.

6.1. Intra-Class Sign Confusion

Issue: A difficulty was evidenced in distinguishing between visually homologous Arabic signs (e.g., the exemplars 'K' and 'X'), which share nearly identical high-level kinematic features.

Mitigation: Remedial adoption of Test-Time Augmentation (TTA), utilizing simulated rotations and slight scale variations, was executed. This was concurrent with the systematic execution of targeted dataset expansion to increase the marginal separation between these classes in the feature space, thereby resolving the ambiguity.

6.2. Hardware Resource Constraints

Issue: The system instability and operational restrictions were directly attributable to the inherent limitations of the GPU's VRAM capacity, which restricted the model complexity and training parameters.

Mitigation: This was addressed through the implementation of memory growth protocols and the rigorous optimization of batch size selection parameters (e.g., maximum size 32). These measures strictly ensured reliable resource utilization without precipitating system failure.

8 7. Forthcoming Stages of Development (Timeline Review)

The following objectives are provisioned for execution during the subsequent reporting period (Term 10), providing a clear trajectory for the project's final application development and system deployment phase.

- **Stage I: Refinement of the Educational Module (Term 10 Kickoff)** Implementation of the Euclidean distance metric within the 3D landmark space, which is required for the quantitative comparison of a learner's instantaneous gesture against the established "Golden Reference" signs.
- **Stage II: Bidirectional Communication Integration (Term 10)** Integration of external Speech-to-Text (STT) and Text-to-Speech (TTS) components to finalize the system's core objective: the establishment of a complete, two-way, multimodal communication loop.
- **Stage III & IV: Application Enclosure and Deployment Strategy (Term 10)** Development of the final Graphical User Interface (GUI) and the integration of the "Combined Architecture" within this enclosure, facilitating the dynamic, run-time selection of the optimal classification model to guarantee maximal system performance.

9 8. Conclusion

A definitive concluding statement formally affirming the project's absolute adherence to the estimated schedule and established technical specifications for the completion of Phase 1 and the initial phase of Phase 2. The core recognition engine functionality is demonstrably robust, and the team is positioned for the final stages of application enclosure, comprehensive system integration, and the development of the educational features, thereby setting the stage for the final project submission.