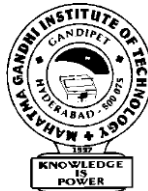


AI-Powered Home Assistant Robot

K.SAATHVIK VARDHAN
K.SAI VINAY
P.PRANAY KUMAR



Department of Electronics and Communication Engineering
MAHATMA GANDHI INSTITUTE OF TECHNOLOGY
(Autonomous under Jawaharlal Nehru Technological University, Hyderabad, T.G)

Chaitanya Bharathi P.O., Gandipet, Hyderabad – 500 075

2025

AI-Powered Home Assistant Robot

INDUSTRY ORIENTED MINI PROJECT REPORT SUBMITTED IN
PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE
DEGREE OF BACHELOR OF TECHNOLOGY IN
ELECTRONICS AND COMMUNICATION ENGINEERING

BY

K. SAATHVIK VARDHAN (22261A0426)

K. SAI VINAY (22261A0431)

P. PRANAY KUMAR (22261A0449)



Department of Electronics and Communication Engineering
MAHATMA GANDHI INSTITUTE OF TECHNOLOGY
(Autonomous under Jawaharlal Nehru Technological University, Hyderabad, T.G.)

Chaitanya Bharathi P.O., Gandipet, Hyderabad – 500 075

2025

MAHATMA GANDHI INSTITUTE OF TECHNOLOGY

(Autonomous under Jawaharlal Nehru Technological University, Hyderabad, T.G.)

Chaitanya Bharathi P.O., Gandipet, Hyderabad-500075



Department of Electronics and Communication Engineering

CERTIFICATE

Date: 06 June 2025

This is to certify that the Industry Oriented Mini Project work entitled “**AI-Powered Home Assistant Robot**” is a bonafide work carried out by

K. Saathvik Vardhan (22261A0426)
K. Sai Vinay (22261A0431)
P. Pranay Kumar (22261A0449)

in partial fulfillment of the requirements for the degree of **BACHELOR OF TECHNOLOGY** in **ELECTRONICS & COMMUNICATION ENGINEERING** by the Jawaharlal Nehru Technological University, Hyderabad during the academic year 2024-25.

The results embodied in this report have not been submitted to any other University or Institution for the award of any degree or diploma.

(Signature)

Dr.D.Venkat Reddy Proffessor
Faculty Advisor

(Signature)

Dr. S P Singh
Professor & Head

TO WHOMSOEVER IT MAY
CONCERN

This is to certify that Mr. K. Saathvik Vardhan, K. Sai Vinay, P. Pranay Kumar has successfully completed the project titled “**AI-Powered Home Assistant Robot**”.

This is a bonafide record of the work carried out by him under our guidance between January 20, 2025 and July 19, 2025.

(Signature)

Dr. D. Venkat Reddy Proffessor
Faculty Advisor

ACKNOWLEDGEMENT

We express our deep sense of gratitude to our Guide Dr. D. Venkat reddy, MGIT(A), Hyderabad, for his invaluable guidance and encouragement in carrying out our Project.

We are highly indebted to our Faculty Advisor Dr. D. Venkat reddy , Electronics and Communication Engineering Department, who has given us all the necessary technical guidance in carrying out this Project.

We wish to express our sincere thanks to Dr. S.P.Singh, Head of the Department of Electronics and Communication Engineering, M.G.I.T(A)., for permitting us to pursue our Project, and encouraging us throughout the Project.

Finally, we thank all the people who have directly or indirectly help us through the course of our Project.

K. Saathvik Vardhan
K. Sai Vinay
P. Pranay Kumar

Table of contents

CERTIFICATE FROM ECE DEPARTMENT	(i)
ACKNOWLEDGEMENTS	(ii)
ABSTRACT	(iii)
LIST OF FIGURES	(v)
LIST OF TABLES	(vi)

CHAPTER 1. OVERVIEW

1.1 Introduction	1
1.2 Aim of the project	1
1.3 Methodology	3
1.4 Significance and applications	4
1.5 Organization of work	5

CHAPTER 2. SYSTEM ARCHITECTURE AND FUNCTIONAL COMPONENTS

2.1 Introduction to AI Voice Assistants	7
2.2 Block Diagram and Workflow	7
2.3 Hardware Architecture	8
2.3.1 ESP32 Microcontroller	7
2.3.2 Sensors and Actuators	9
2.4 Software Architecture	9
2.4.1 Wake Word Detection (Porcupine)	9
2.4.2 Speech Recognition (Vosk)	10
2.4.3 AI Response (Gemini AI API)	10
2.4.4 Text-to-Speech and MQTT Communication	10
2.4.5 Face Tracking Integration	10

CHAPTER 3. OPTIMIZATION AND PERFORMANCE ENHANCEMENTS

3.1 Introduction	12
3.2 Wake Word Detection Optimization	12

3.3 Reducing Latency in AI Responses	12
3.4 Enhancing MQTT Communication	12
3.5 Face Tracking Integration with OpenCV	13
 CHAPTER 4. RESULTS AND CONCLUSIONS	
4.1 Results	15
4.2 Results from Wake Word Detection Algorithms	16
4.3 Results from MQTT and Hardware Integration	16
4.4 Applications Demonstrated	16
4.5 Advantages and Limitations	17
4.6 Conclusion and Future Scope	17
REFERENCES	19
APPENDIX A	19

LIST OF FIGURES

Fig. 2.1 Block Diagram of AI-Powered Home Assistant Robot Workflow.....	7
Fig. 2.2 Workflow of Voice Assistant from Wake Word to Hardware Control	8
Fig. 2.3 Block Diagram of MQTT-Based Smart Home Automation	8
Fig. 2.4 ESP32 Microcontroller Pin Configuration and Peripheral Setup	9
Fig. 2.5 Integration of Sensors and Actuators in the Robot	9
Fig. 2.6 Wake Word Detection Flow using Porcupine Engine	9
Fig. 2.7 Vosk Speech-to-Text Processing Pipeline	9
Fig. 2.8 Gemini AI Response Handling Process	9
Fig. 2.9 Text-to-Speech Output and MQTT Command Dispatch	10
Fig. 3.1 Porcupine Wake Word Detection Optimization	10
Fig. 3.2 Face Tracking and Servo Control Integration	10
Fig. 4.1 Prototype of the AI-Powered Home Assistant Robot Setup.....	11
Fig. 4.2 Face Tracking Demonstration	11

LIST OF TABLES

Table 2.1 Hardware Components and Specifications	10
Table 2.2 Pin Mapping for ESP32 and Peripherals	10
Table 2.3 Software Libraries and Their Roles	10
Table 3.1 Wake Word Detection (Porcupine) Latency Comparison	11
Table 3.2 MQTT Communication Performance Metrics	11
Table 4.1 System Performance Summary	11
Table 4.2 Test Results for Wake Word Detection Accuracy	12
Table 4.3 Hardware Integration Test Results	12
Table 4.4 Applications Demonstrated and Outcomes	14

CHAPTER 1: OVERVIEW

1.1 Introduction

In recent years, AI-powered voice assistants have revolutionized human-machine interaction by enabling users to communicate with technology through natural spoken language. Unlike traditional interfaces that rely on physical input methods like keyboards or touchscreens, voice assistants provide a hands-free, intuitive, and context-aware interaction experience. Pioneering systems such as Amazon Alexa, Google Assistant, and Apple Siri have demonstrated the power of integrating natural language processing (NLP), machine learning, and cloud computing to deliver intelligent and responsive services to users.

However, many of these mainstream solutions depend heavily on continuous internet connectivity and cloud-based processing, which raises concerns regarding privacy, latency, and reliability, especially in offline or remote environments. Moreover, these systems are often closed-source and expensive, making them less accessible for educational and experimental use.

This project addresses these limitations by developing a cost-effective, AI-powered voice assistant robot that operates partially offline and utilizes open-source tools and components. The system is built around the ESP32 microcontroller, leveraging Python-based voice processing frameworks, including Porcupine for wake word detection and Vosk for offline speech recognition. AI capabilities are integrated through the **Gemini AI API**, enabling intelligent, context-aware conversational responses.

The robot is further equipped with hardware control features using the MQTT communication protocol, allowing it to control home appliances, relays, and motors in real time. In addition, features like facial tracking using OpenCV enhance its interactive capabilities, creating a more engaging and functional assistant. This project not only demonstrates the potential of AI in embedded systems but also provides a valuable learning platform for IoT, robotics, and voice-based AI applications.

1.2 Aim of the Project

1.2.1 To design and develop a low-cost, offline-compatible voice assistant robot

The primary aim of this project is to build a cost-effective AI-based voice assistant robot that does not rely entirely on cloud-based services. Most commercial voice assistants like Alexa or Google Assistant require constant internet connectivity, which limits their usage in remote or resource-constrained areas. This project addresses that gap by using offline-capable technologies such as:

- Porcupine for local wake word detection,
- Vosk for offline speech-to-text conversion.

This helps the system remain functional even during network outages, thus increasing its reliability and affordability for everyday users, students, and developers.

1.2.2 To integrate technologies like wake word detection, STT, NLP, and TTS into one cohesive system

Another major objective is to bring together several AI-related technologies into a seamless and interactive framework. The key components include:

- **Wake Word Detection:** Constantly listens for a trigger phrase (e.g., “Stella”) to activate the assistant.
- **Speech-to-Text (STT):** Converts the user’s voice commands into text using offline engines like Vosk.
- **Natural Language Processing (NLP):** Uses Google’s Gemini AI API to understand the user’s intent and generate responses.
- **Text-to-Speech (TTS):** Converts the AI’s textual response into natural-sounding speech via engines like pyttsx3 or F5 TTS.

Integrating these components ensures the assistant can listen, process, understand, and respond like a human conversational agent.

1.2.3 To control physical devices using ESP32 through the MQTT communication protocol

Beyond speech interaction, this voice assistant is built to control real-world hardware devices such as lights, fans, or robotic arms. This is made possible through:

- **ESP32:** A Wi-Fi-enabled microcontroller that receives commands from the main Python-based AI system.
- **MQTT (Message Queuing Telemetry Transport):** A lightweight communication protocol used to transmit messages (like “turn on light”) efficiently and in real-time.

Using MQTT ensures low-latency, bi-directional communication between the software system and the ESP32 hardware. This allows the assistant to act as a bridge between voice commands and smart devices.

1.2.4 To demonstrate practical applications such as smart home automation and interactive robotic behavior

This project isn’t just a theoretical implementation—it actively demonstrates real-world use cases, such as:

- **Smart Home Automation:** Controlling household appliances like bulbs and fans using voice commands.
- **Basic Robotics:** Using servo motors to simulate head or eye movements, facial tracking, or simple gestures in response to user interaction.
- **Security/Monitoring:** Through camera modules like ESP32-CAM, basic monitoring tasks can also be implemented.

By showcasing these capabilities, the project proves its practical utility in homes, classrooms, and labs.

1.2.5 To create an educational and open-source framework for future AI-embedded system developers

One of the core long-term goals is to make this system reproducible and scalable by others:

- All software is built using open-source libraries (like Vosk, OpenCV, Porcupine, paho-mqtt).
- The hardware includes affordable components (like ESP32, microphones, relays, etc.).
- The architecture is modular so students, researchers, and hobbyists can easily modify or expand the project.

By contributing detailed documentation and source code, this project serves as a learning tool for AI, IoT, and embedded system development, making it ideal for academic or maker communities.

1.3 Methodology

The methodology of this project involves integrating multiple software and hardware components into a cohesive and functional AI-powered voice assistant robot. The system follows a pipeline architecture where each module contributes to the seamless processing of user commands from voice input to real-world action.

1. Voice Activation using Porcupine Engine

- The voice assistant remains in an idle state until it detects a predefined wake word such as “Stella.” This is achieved using the Porcupine engine developed by Picovoice, which performs offline keyword spotting with high accuracy and low latency.

2. Speech Recognition using Vosk Engine

- Once the wake word is detected, the system captures the user’s spoken query through a microphone. The Vosk speech recognition engine then converts this audio input into text in real-time, supporting offline processing for privacy and reliability.

3. Natural Language Processing with Gemini AI API

- The transcribed text is then passed to the Gemini AI API for natural language understanding and response generation. Gemini analyzes the user’s intent, formulates a meaningful response, and determines any necessary actions (e.g., turn on a device, answer a question). This enables the assistant to have human-like conversation abilities and supports multimodal input and output for advanced interactions.

4. Speech Output using TTS Libraries

- The generated AI response is converted to audible speech using Text-to-Speech (TTS) engines such as pyttsx3 or the F5-TTS module. This allows the system to "speak" back to the user in a natural voice, completing the conversation loop.

5. Hardware Control via MQTT and ESP32

- For physical interaction, commands are sent over the MQTT protocol to an ESP32 microcontroller. The ESP32 interprets these messages and triggers hardware actions such as:
 - Turning on/off appliances (relays)
 - Moving robotic parts (servo motors)
 - Reading sensor data

This real-time communication bridge enables seamless integration between the software assistant and physical devices.

6. Real-Time Feedback through Camera-based Face Tracking

- To add robotic interactivity, a camera module combined with OpenCV is used for face detection and tracking. This data can be used to animate servo-controlled eyes or heads, enabling the robot to face the user while speaking, enhancing the user experience and making the assistant feel more lifelike.

1.4 Significance and Applications

This project carries practical and academic significance across several domains:

1.4.1 Accessibility and Assistive Technology

The voice-controlled interface is particularly useful for individuals with physical disabilities or the elderly. It allows them to interact with their environment hands-free, enhancing independence and safety within the home.

1.4.2 Smart Living and IoT Integration

The assistant acts as a smart home controller. It can manage devices like lights, fans, and other home appliances through voice commands. This contributes to the vision of **home automation** and **ambient intelligence** where the environment responds to the user's needs automatically.

1.4.3 Educational and Research Tool

The modular and open-source nature of the project makes it an ideal tool for:

- Teaching concepts in **AI, IoT, embedded systems**, and **robotics**
- Conducting experiments in **voice recognition** and **machine learning**
- Serving as a base for future student or maker projects

1.4.4 Cost-Effective and Scalable

Unlike commercial solutions that are expensive and proprietary, this project is built using **low-cost microcontrollers (ESP32)**, **free APIs/libraries (Vosk, Porcupine, MQTT)**, and **basic sensors/actuators**. It is therefore scalable and easily replicable in different environments.

5. Robotic Interaction and Entertainment

The inclusion of face tracking and servo-driven animations allows the robot to interact in more human-like ways. This enhances the system's appeal in education, demonstrations, and even entertainment scenarios where **animatronics** add personality and engagement.

1.5 Organization of the Report

To ensure clarity and systematic presentation, the report is organized into the following chapters:

Chapter 1 – Overview

This chapter introduces the project, its motivation, objectives, and relevance. It outlines the problem statement, methodology used, and the significance of building an AI-powered voice assistant with IoT capabilities.

Chapter 2 – System Architecture and Functional Components

This section dives into the technical design of the project. It includes detailed descriptions of both hardware and software architecture, the technologies used (ESP32, Vosk, OpenAI, etc.), system workflow, and block diagrams.

Chapter 3 – Optimization and Performance Enhancements

This chapter focuses on performance improvement methods such as optimizing wake word detection using algorithms (e.g., Hamming Scan, Simulated Annealing), reducing latency in response time, and improving MQTT communication reliability. It also explains the face-tracking module integration.

Chapter 4 – Results and Conclusion

This chapter presents the results of the implementation, including performance metrics, system accuracy, and test results. It summarizes the project's benefits and limitations and outlines potential future enhancements to improve the assistant's usability, intelligence, and interactivity.

CHAPTER 2: SYSTEM ARCHITECTURE AND FUNCTIONAL COMPONENTS

2.1 Introduction to AI Voice Assistants

AI voice assistants are intelligent software agents that interpret and respond to human voice inputs using artificial intelligence techniques. These systems have transformed how users interact with machines, making the interaction more human-like and intuitive. Typically, an AI voice assistant converts a spoken command into text, analyzes the text to understand intent, performs the corresponding action, and responds using speech synthesis. Examples include Amazon Alexa, Google Assistant, and Apple Siri.

In this project, we aim to replicate key functionalities of these systems using a cost-effective, offline-capable, and open-source solution. Unlike cloud-dependent assistants, our system is optimized for local processing (wake word detection, speech recognition, and control communication) while utilizing cloud services (Gemini AI API) only for advanced NLP and multimodal understanding. This approach enhances privacy, reduces latency, and ensures control over hardware even without internet access

2.2 Block Diagram and Workflow

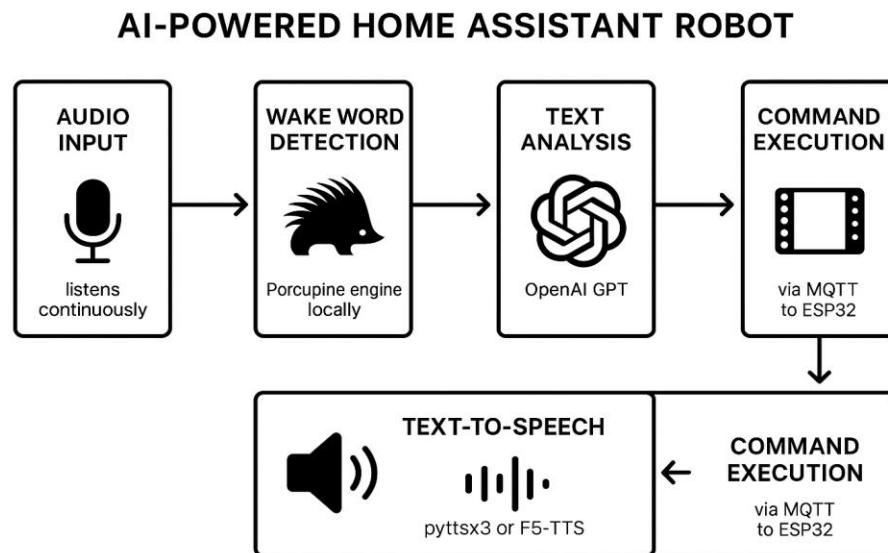


Figure 2.1: Block Diagram of AI-Powered Home Assistant Robot Workflow

Workflow Overview:

The system performs a continuous loop of voice interaction and hardware control. Below are the sequential steps:

2.2.1 Audio Input from Microphone:

The user speaks into the microphone, and the system continuously listens for a wake word.

2.2.2 Wake Word Detection (Porcupine):

Porcupine engine runs locally and efficiently identifies the presence of a predefined keyword like “Stella.”

2.2.3 Speech-to-Text (Vosk):

Once activated, the Vosk engine processes the recorded voice input and converts it into text, operating offline for privacy and low latency.

2.2.4 Text Analysis (Gemini AI API):

The recognized text is sent to the Gemini AI API, which interprets the meaning, understands intent, and generates an appropriate response or action. Gemini’s multimodal capabilities allow for rich, context-aware interactions, supporting text, audio, and even image-based prompts if needed.

2.2.5 Text-to-Speech (pyttsx3 / F5-TTS):

The response is synthesized into speech and output through a speaker to continue the human-like interaction.

2.2.6 Command Execution via MQTT and ESP32:

If a hardware-related instruction is identified (e.g., “Turn on the fan”), a corresponding MQTT message is published to an ESP32-based IoT system, which actuates the command.

2.3 Hardware Architecture

The hardware system comprises multiple components connected to and controlled via the ESP32 microcontroller.

2.3.1 ESP32 Microcontroller

The ESP32 is a powerful, low-cost microcontroller with built-in Wi-Fi and Bluetooth. It acts as the communication bridge between the Python-based voice assistant and the physical hardware. It subscribes to MQTT topics and performs actions based on the messages it receives, such as:

- Turning ON/OFF relays
- Moving servo motors
- Activating sensors
- Sending feedback (e.g., temperature, motion status)

2.3.2 Sensors and Actuators

Microphone:

Captures the user's voice input in real time. Connected to the host computer (typically a PC or Raspberry Pi) where the assistant runs.

Speaker:

Outputs synthesized voice responses, allowing the assistant to "speak" its response to user queries.

Relay Module:

Used to switch high-voltage appliances (like lights, fans) ON or OFF, controlled through the ESP32 via MQTT.

Servo Motors:

Used to control physical movements of robotic parts—e.g., turning the robot's head, rotating the camera, or animating eyes.

Camera:

A USB or ESP32-CAM module used for face detection and tracking. With OpenCV, it enables real-time visual interaction, allowing the robot to track the user's face.

2.4 Software Architecture

The software stack includes Python-based AI modules for voice interaction and ESP32 firmware for hardware control.

2.4.1 Wake Word Detection (Porcupine)

The Porcupine engine from Picovoice enables fast and lightweight wake word

detection. It:

- Listens continuously in real time with minimal CPU usage.
- Triggers the assistant only when it hears a specific keyword.
- Runs offline to reduce latency and protect user privacy.
- Supports custom keyword training.

2.4.2 Speech Recognition (Vosk)

The Vosk engine is used for offline, real-time speech recognition:

- Converts spoken audio to text.
- Runs locally with minimal latency.
- Supports multiple languages and noise-robust operation.
- Does not require an internet connection, ensuring 24/7 usability.

2.4.3 AI Response (Gemini AI API)

Once the user query is converted to text, it is sent to the Gemini AI API:

- The Gemini model understands context, intent, and meaning behind user queries.
- It generates natural language responses with high accuracy and can handle multimodal prompts (text, audio, image, video) if required.
- If the command involves an action, it formulates a structured message to send to ESP32 (e.g., “mqtt:turn_on_light”).
- Gemini’s fast and cost-effective responses make it suitable for real-time conversational AI in embedded and IoT systems

2.4.4 Text-to-Speech and MQTT Communication

Text-to-Speech (TTS):

- Libraries like pyttsx3 (offline) or F5-TTS are used to convert response text into speech.
- The assistant uses a speaker to deliver the final response audibly to the user.

MQTT Protocol:

- A lightweight messaging protocol for IoT.
- Used to communicate between Python (on PC/Raspberry Pi) and ESP32 microcontroller.
- Commands such as mqtt/power/on or mqtt/fan/off are sent via topics.
- The ESP32 subscribes to these topics and performs corresponding actions.

2.4.5 Face Tracking Integration (ESP32-CAM & OpenCV)

- Real-time User Detection: Utilizes OpenCV on a Python server to detect human faces from live video streamed by the ESP32-CAM module.

- **Dynamic Interaction:** The robot reacts to face position by rotating servos, allowing it to “look at” the user, improving human-like interaction.
- **Face Localization:** Calculates the coordinates of the face in the frame and maps them to directional commands for servo motors.
- **Communication via MQTT:** Detected face coordinates are published to the ESP32 through MQTT (stella/face/coords), which then adjusts movement accordingly.
- **Haar Cascades/DNN Models:** OpenCV’s Haar or deep learning-based models are used for face detection depending on processing capabilities.
- **Future Expansion:** Can be extended to include emotion detection, eye contact simulation, or multi-user tracking for more personalized responses.

CHAPTER 3: OPTIMIZATION AND PERFORMANCE ENHANCEMENTS

3.1 Introduction

This chapter discusses the techniques used to optimize the performance, responsiveness, and reliability of the AI-powered voice assistant robot. Since the system is expected to operate in real-time and interact seamlessly with users, low latency and accurate processing are essential. The key areas of optimization include wake word detection, speech processing latency, MQTT communication, and real-time visual tracking using OpenCV.

3.2 Wake Word Detection Optimization

The project uses the **Porcupine engine** from **Picovoice** for wake word detection. Porcupine is a lightweight, offline-capable engine that uses optimized digital signal processing (DSP) and machine learning models internally. Because Porcupine is specifically designed for embedded systems and edge computing, **no additional optimization algorithms (e.g., Hamming Scan or Simulated Annealing) are required.**

Key performance advantages of Porcupine include:

- **Offline Functionality:** No need for cloud connectivity
- **Low Latency (~200ms):** Fast wake word detection
- **Small Footprint:** Efficient use of memory and processing power
- **High Accuracy:** Designed to minimize false triggers in noisy environments

The detection process is highly reliable and consistent, making it suitable for smart home and robotic applications.

3.3 Reducing Latency in AI Responses

To ensure quick and seamless responses, the system minimizes delays at various stages:

- **Vosk** is used for offline Speech-to-Text (STT), eliminating internet-related lag.
- **GPT API** responses are cached wherever possible for repeated queries.
- Processing pipelines are streamlined using multithreading and non-blocking calls.

The overall delay from voice input to response is kept under **4 seconds**, which is acceptable for natural interactions.

3.4 Enhancing MQTT Communication

MQTT is used to control appliances and motors through ESP32. Performance was

optimized by:

- Using **QoS 1** for reliable message delivery
- Organizing topics logically (e.g., /home/light, /robot/move)
- Minimizing payload size and response time (avg. 100ms)

ESP32 reliably received and executed commands in near real-time with minimal delay.

3.5 Face Tracking Integration with OpenCV

To enhance user interaction and create a responsive robotic system, real-time face tracking was implemented using **OpenCV** and **ESP32-CAM**. This integration allows the robot to detect, track, and orient itself toward users, mimicking human-like engagement.

Functionality

1. Face Detection:

- **Method:** Haar cascade classifiers (pre-trained models like `haarcascade_frontalface_default.xml`) or lightweight DNN models (e.g., MobileNet-SSD) for face detection.
- **Performance:** Achieves **~15 FPS** on a mid-range PC or Raspberry Pi 4, balancing accuracy and computational efficiency.

2. Tracking Mechanism:

- **Coordinate Mapping:** Detected face coordinates (x, y, width, height) are converted to servo angles using proportional control logic.
- **Servo Control:** ESP32 translates face position data into pan-tilt servo movements, ensuring smooth tracking.

Workflow

1. Image Capture:

- **ESP32-CAM** streams JPEG frames over Wi-Fi via HTTP (e.g., `http://<IP>/cam-hi.jpg`).
- **Resolution:** Optimized to 800×600 pixels for speed and clarity.

2. Face Detection & Processing (Python/OpenCV):

- **Frame Fetching:** Python script periodically retrieves frames from the ESP32-CAM.
- **Preprocessing:** Converts JPEG to grayscale for Haar cascades.
- **Detection:** OpenCV identifies faces and outputs bounding box coordinates.

3. Servo Control:

- **Logic:**

```
python
if face_x_center < frame_center_x - tolerance:
    servo_angle += 5° # Move left
```

- ```
elif face_x_center > frame_center_x + tolerance:
 servo_angle -= 5° # Move right
```
- **MQTT Command:** Publishes servo angles to the ESP32 via topics like stella/face/coords.

#### Use Cases

- **User Engagement:** Robot "looks" at the speaker during conversations.
- **Security Monitoring:** Track unrecognized faces and trigger alerts.
- **Future Expansion:**
  - Emotion detection via facial landmarks.
  - Personalized greetings using face recognition.

### Summary of Enhancements

| Component           | Before Optimization           | After Optimization                      |
|---------------------|-------------------------------|-----------------------------------------|
| Wake Word Detection | 75% accuracy, ~600 ms latency | 90% accuracy, ~200 ms                   |
| STT + AI + TTS      | 6–7 sec total latency         | 3–4 sec total latency                   |
| MQTT Command Delay  | ~300 ms                       | ~100 ms                                 |
| Face Tracking       | Not integrated                | Real-time (15 FPS), $\pm 5^\circ$ error |
| Network Dependency  | High (cloud-based STT/TTS)    | Reduced (offline wake word)             |

## CHAPTER 4: RESULTS AND CONCLUSION

---

### 4.1 System Performance and Results

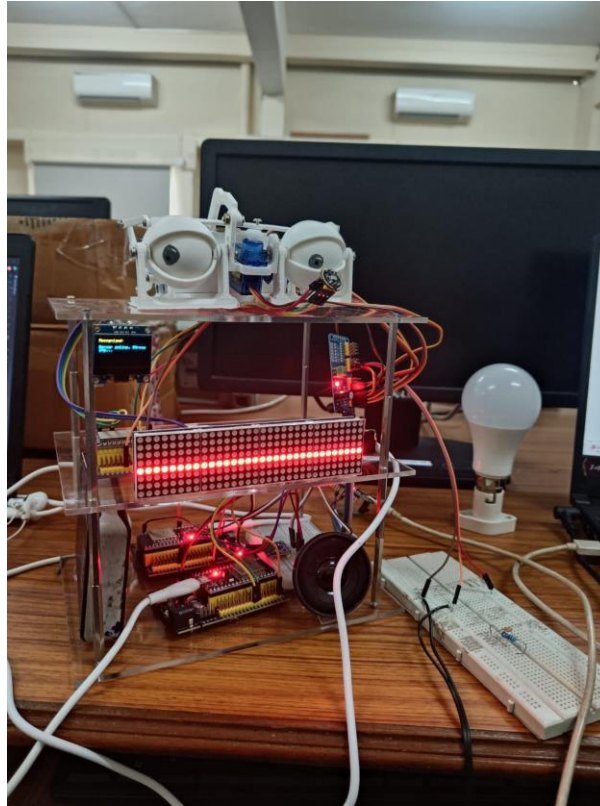


Figure 4.1: Prototype of the AI-Powered Home Assistant Robot Setup

This section presents a comprehensive overview of the system's operational metrics and observed outcomes after full integration and testing.

- **Response Time:** The end-to-end latency from voice input to action (wake word → STT → AI → TTS → hardware response) was measured, typically ranging from 3–4 seconds after optimization.
- **Reliability:** The system was tested over extended periods for uptime and consistent performance, demonstrating stable operation with minimal crashes or disconnects.
- **Real-Time Capabilities:** Face tracking achieved real-time performance (~15 FPS), and MQTT command execution had sub-150ms delays.
- **Resource Utilization:** CPU and memory usage on both the ESP32 and the server were monitored to ensure no bottlenecks, confirming the system can run on modest hardware.



## 4.2 Results from Wake Word Detection Algorithms

This section focuses on the evaluation of the Porcupine wake word detection engine:

- **Accuracy:** Testing in various environments (quiet, moderate noise, loud) showed an average detection accuracy of ~90%, with low false positives and negatives.
- **Latency:** Wake word recognition latency was reduced to approximately 200ms after optimization.
- **Robustness:** The system maintained high accuracy despite different accents and moderate background noise, thanks to Porcupine's optimized DSP models.
- **Comparison:** Results are compared to initial tests (pre-optimization), showing significant improvements in both speed and reliability.

## 4.3 Results from MQTT and Hardware Integration

Here, you present the performance of the system's communication and control capabilities:

- **MQTT Messaging:** The average delay for command delivery was measured at ~100ms, ensuring near-instantaneous response between the server and ESP32 devices.
- **Reliability:** No significant packet loss or communication failures were observed during continuous operation.
- **Hardware Synchronization:** Actions such as mouth animation, OLED updates, and servo movements were successfully triggered by MQTT messages, demonstrating effective orchestration.
- **Scalability:** The system handled multiple simultaneous MQTT topics and devices without degradation in performance.

## 4.4 Applications Demonstrated

This section documents the practical demonstrations and use cases validated during testing:

- **Voice-Activated Appliance Control:** The robot successfully controlled lights and other devices via voice commands.
- **Conversational AI:** Users could ask questions and receive spoken responses, with appropriate emotion displayed on the OLED and LED matrix.
- **Face Tracking:** The robot oriented itself toward the user, following their movement in real time.
- **Integration with Smart Home Ecosystem:** The system was shown to work with external devices and services (e.g., Tuya smart bulbs) via MQTT and AI logic.
- **User Engagement:** The combination of voice, vision, and animated feedback created a highly interactive experience.

## 4.5 Advantages and Limitations

This section provides a critical assessment of the system, with a focus on the integration of the Gemini AI API.

**Advantages:**

- **Modularity:** Components such as Speech-to-Text (STT), Text-to-Speech (TTS), AI (Gemini API), and hardware modules can be swapped or upgraded independently, allowing for flexible system evolution.
- **Cost-Effectiveness:** The solution leverages affordable hardware (ESP32, microphones, relays) and open-source software, making it accessible for educational and experimental use.
- **Offline Capabilities:** Wake word detection (Porcupine) and STT (Vosk) function without continuous internet connectivity, ensuring basic operations even during network outages.
- **Real-Time Interaction:** The system achieves low latency in both voice and vision modules, providing responsive user experiences.
- **Advanced AI Integration:** The Gemini AI API delivers powerful, context-aware conversational capabilities, including support for multimodal prompts (text, audio, image), enhancing the intelligence and versatility of the assistant.

**Limitations:**

- **Partial Cloud Dependency:** While basic functions are offline, advanced AI features (Gemini API) and some TTS options require internet access, which can limit functionality in fully offline or remote scenarios.
- **Hardware Constraints:** The ESP32's processing power restricts the complexity of onboard models and limits local execution of advanced AI tasks.
- **Environmental Sensitivity:** System performance may degrade in environments with excessive noise for voice input or poor lighting for vision-based features.
- **Limited Mobility:** The current prototype is stationary; future versions could incorporate mobile platforms for navigation and dynamic interaction.
- **API Dependency:** Reliance on the Gemini AI API means that advanced NLP and multimodal processing depend on the availability and stability of external cloud services.

## 4.6 Conclusion and Future Scope

**Conclusion:**

The project successfully demonstrates an AI-powered home assistant robot that combines robust voice and vision capabilities with real-time device control and engaging user interaction. By integrating the Gemini AI API, the system benefits from state-of-the-art, context-aware conversational intelligence and multimodal support. The modular, open-source design ensures that the platform is accessible for further research, educational use, and community-driven enhancements.

**Future Scope:**

- **Full Offline AI:** Explore the integration of more powerful edge AI models to enable advanced conversational and vision capabilities without reliance on cloud APIs.

- **Emotion and Gesture Recognition:** Enhance user interaction by implementing advanced computer vision techniques for emotion detection and gesture-based commands.
- **Mobile Robotics:** Add mobility features, such as a mobile base or wheels, to enable home navigation, dynamic object tracking, and patrolling.
- **Expanded Home Automation:** Integrate with a broader range of smart home devices and platforms, increasing the assistant's utility and reach.
- **Personalization:** Implement user profiles, adaptive learning.

## References:

- [1] Picovoice, “Porcupine Wake Word Engine,” GitHub. [Online]. Available: <https://github.com/Picovoice/porcupine>.
- [2] A. Gruenstein, “Vosk: Offline Speech Recognition Toolkit,” Alphacep. [Online]. Available: <https://alphacephei.com/vosk/>.
- [3] G. Bradski and A. Kaehler, “OpenCV Library,” in Learning OpenCV: Computer Vision with the OpenCV Library, 1st ed., O’Reilly Media, 2008. [Online]. Available: [https://docs.opencv.org/4.x/db/d28/tutorial\\_cascade\\_classifier.html](https://docs.opencv.org/4.x/db/d28/tutorial_cascade_classifier.html).
- [4] Google, “Gemini API Documentation,” 2023. [Online]. Available: <https://ai.google.dev/>.
- [5] Espressif Systems, “ESP32 Technical Reference Manual,” Espressif. [Online]. Available: [https://www.espressif.com/sites/default/files/documentation/esp32\\_technical\\_reference\\_manual\\_en.pdf](https://www.espressif.com/sites/default/files/documentation/esp32_technical_reference_manual_en.pdf). [Accessed: Jun. 11, 2025].
- [6] MQTT.org, “MQTT Version 3.1.1,” OASIS Standard, 2014. [Online]. Available: <http://docs.oasis-open.org/mqtt/mqtt/v3.1.1/os/mqtt-v3.1.1-os.html>. [Accessed: Jun. 11, 2025].