

**A**  
**PROJECT REPORT**  
**ON**  
**SeeU: Real-Time Object Detection and Distance**  
**Estimation App**

*Submitted by*  
**Jaivik Jamnadas Rathod(21IT552)**

# Abstract

SeeU is an AI-powered mobile application designed to assist visually impaired individuals by detecting objects in real-time using their smartphone camera. The application identifies objects, estimates their distance, and provides audio feedback, enabling users to navigate safely in their surroundings.

The key components of SeeU include:

1. Real-Time Object Detection - Utilizing YOLO, the app can detect objects live through the camera.
2. Depth Estimation for Distance Calculation - Using MiDaS, the app estimates the distance of detected objects.
3. Audio Feedback System - The detected object and distance are announced via Text-to-Speech (TTS).
4. Optimized Mobile Performance - The backend is developed using FastAPI, processing frames and sending results efficiently.

The application enhances mobility for the visually impaired by providing real-time object awareness and surroundings estimation. Future improvements include gesture-based controls, multi-language support, and better AI accuracy.

## List of Figures

|                                  |    |
|----------------------------------|----|
| Figure 1. Gantt chart.....       | 09 |
| Figure 2. Use Case Diagram ..... | 14 |
| Figure 3. DFD Level 0 .....      | 15 |
| Figure 4. DFD Level 1 .....      | 15 |
| Figure 5. Class Diagram .....    | 16 |
| Figure 6. ER Diagram .....       | 17 |
| Figure 7. Flow Chart .....       | 18 |
| Figure 8. Simple UI.....         | 19 |
| Figure 9. Listening UI.....      | 20 |
| Figure 10. Laptop Detected.....  | 21 |
| Figure 11. Bottle Detected ..... | 21 |
| Figure 12. Question Asked .....  | 22 |
| Figure 13. LLM Response .....    | 22 |
| Figure 14. Asked question .....  | 23 |
| Figure 15. Log of backend .....  | 24 |

## List of Abbreviations

| Abbreviation | Full Form                           |
|--------------|-------------------------------------|
| AI           | Artificial Intelligence             |
| API          | Application Programming Interface   |
| CNN          | Convolutional Neural Network        |
| CUDA         | Compute Unified Device Architecture |
| YOLO         | You Only Look Once                  |
| DFD          | Data Flow Diagram                   |
| LLM          | Large Language Model                |
| HTTP         | Hypertext Transfer Protocol         |
| LiDAR        | Light Detection and Ranging         |
| MiDaS        | Multi-scale Depth Estimation        |
| STT          | Speech-to-Text                      |
| TTS          | Text-to-Speech                      |
| UI           | User Interface                      |

## Table of Contents

|  |    |
|--|----|
| Chapter 1 : Introduction .....   | 1  |
| 1.1 Brief Overview of the Work.....  | 1  |
| 1.2 Objective .....  | 1  |
| 1.3 Scope .....  | 1  |
| 1.4 Project Modules.....   | 1  |
| 1.5 Project Hardware/Software Requirements .....                                       | 2  |
| 1.5.1 Hardware Requirements.....   | 2  |
| 1.5.2 Software Requirements .....  | 2  |
| Chapter 2: Literature Review .....   | 3  |
| 2.1 Object Detection Models .....  | 3  |
| 2.2 Depth Estimation Techniques .....  | 3  |
| 2.3 Assistive Technologies for Visually Impaired Individuals.....                      | 4  |
| 2.4 Challenges and Limitations.....  | 5  |
| Chapter 3: System Analysis & Design.....   | 6  |
| 3.1 Comparison of Existing Applications with my project with merits and demerits ..... | 6  |
| 3.2 Project Feasibility Study .....  | 7  |
| 3.2.1 Technical Feasibility .....  | 7  |
| 3.2.2 Operational Feasibility.....   | 7  |
| 3.2.3 Economic Feasibility .....   | 7  |
| 3.3 Requirement Gathering.....   | 8  |
| 3.4 Project Timeline chart.....  | 9  |
| 3.5 Detailed Modules Description .....   | 9  |
| 3.5.1 User Interface (UI) .....  | 9  |
| 3.5.2 Object Detection Module .....  | 10 |
| 3.5.3 Distance Estimation Module .....   | 11 |

|   |    |
|---|----|
| 3.5.4 Audio Feedback System .....           | 13 |
| 3.5.5 LLM Processing Module .....           | 14 |
| 3.6 Project SRS .....                       | 16 |
| 3.6.1 Use Case Diagram.....                 | 16 |
| 3.6.2 Data Flow Diagrams .....              | 17 |
| 3.6.3 Class Diagram.....                    | 18 |
| 3.6.4 Entity Relationship Diagrams .....    | 19 |
| 3.6.5 Flow Chart .....                      | 20 |
| Chapter 4: Implementation and Testing.....  | 21 |
| 4.1 User Interface Snapshots.....           | 21 |
| 4.2 Testing.....                            | 23 |
| Chapter 5: Conclusion and Future Work ..... | 27 |
| 5.1 Conclusion .....                        | 27 |
| 5.2 Limitations .....                       | 27 |
| 5.2 Future Work.....                        | 28 |
| Chapter 6. References .....                 | 29 |

# Chapter 1 : Introduction

## 1.1 Brief Overview of the Work

SeeU is an AI-powered mobile application designed to assist visually impaired individuals by detecting objects in real-time using their smartphone camera. The app analyzes live camera frames, estimates the distance of detected objects, and provides audio feedback to help users navigate safely. The app leverages advanced technologies such as YOLOv8 for object detection, MiDaS for depth estimation, and FastAPI for backend processing.

## 1.2 Objective

- To develop a real-time object detection and distance estimation app for visually impaired individuals.
- To provide accurate and timely audio feedback to users about their surroundings.
- To ensure the app is optimized for mobile devices with fast inference and low latency.
- To create a seamless communication system between the mobile app and the backend server.

## 1.3 Scope

- Implement real-time object detection using YOLOv8.
- Estimate the distance of detected objects using MiDaS depth estimation.
- Provide audio feedback using Android Text-to-Speech (TTS).
- Develop a FastAPI backend for processing frames and returning object and distance information.
- Ensure the app is optimized for mobile devices using TensorFlow Lite and PyTorch Mobile.

## 1.4 Project Modules

- User Interface (UI)
- Object Detection Module
- Distance Estimation Module
- Audio Feedback System
- Backend Processing Module

## **1.5 Project Hardware/Software Requirements**

### **1.5.1 Hardware Requirements**

- Mobile Device: Smartphone with a camera.
- Server: Backend server with adequate processing power for real-time frame processing.

### **1.5.2 Software Requirements**

- Frontend (Mobile App):
  - Language: Dart
  - Framework: Flutter
  - Camera API: Camera package for Flutter
  - Networking: HTTP client package for Flutter (e.g., http)
  - Image Processing: Bitmap & Base64 Encoding
  - Audio Feedback: Coqui TTS XTTS-v2
- Backend (FastAPI Server):
  - Framework: FastAPI
  - Object Detection: YOLOv8s
  - Depth Estimation: MiDaS (PyTorch)
  - Image Handling: OpenCV & NumPy
  - Model Deployment: PyTorch



## Chapter 2: Literature Review

The development of the SeeU application builds upon existing research and technologies in the fields of object detection, depth estimation, and assistive technologies for visually impaired individuals. This section reviews key studies, methodologies, and tools that have influenced the design and implementation of the application.

### 2.1 Object Detection Models

Object detection is a critical component of the SeeU application, enabling the identification of objects in real-time. Several state-of-the-art models have been developed for object detection, each with its strengths and limitations:

- **YOLO (You Only Look Once)** : YOLO is a popular real-time object detection system known for its speed and accuracy. The latest version, YOLOv8, offers improved performance and is capable of detecting objects in real-time with high precision [2]. However, it requires more computational resources compared to MobileNet SSD.
- **MobileNet SSD** : MobileNet is a lightweight convolutional neural network (CNN) designed for mobile and embedded devices. When combined with Single Shot MultiBox Detector (SSD), it provides a balance between accuracy and speed, making it suitable for real-time applications [1]. MobileNet SSD is widely used in mobile applications due to its efficiency and low computational requirements.
- **Faster R-CNN** : Faster Region-based CNN (Faster R-CNN) is another widely used object detection model. It provides high accuracy but is computationally intensive, making it less suitable for real-time applications on mobile devices [3].

### 2.2 Depth Estimation Techniques

Depth estimation is essential for calculating the distance between the user and detected objects. Several techniques and models have been developed for depth estimation:

- **MiDaS (Multi-scale Depth Estimation)** : MiDaS is a state-of-the-art model for depth estimation that works across different scales and environments. It is capable

of estimating depth from a single image, making it highly suitable for real-time applications [4]. MiDaS is known for its robustness and accuracy in various lighting conditions.

- **Stereo Vision** : Stereo vision techniques use two cameras to capture images from slightly different angles, allowing the system to calculate depth based on the disparity between the images. While effective, stereo vision requires specialized hardware and is computationally expensive [5].
- **LiDAR-based Depth Estimation** : LiDAR (Light Detection and Ranging) uses laser pulses to measure distances. It is highly accurate but requires expensive hardware, making it less feasible for mobile applications [6].

## 2.3 Assistive Technologies for Visually Impaired Individuals

Several assistive technologies have been developed to help visually impaired individuals navigate their surroundings. These technologies range from simple tools to advanced AI-based systems:

- **Smart Canes** : Traditional white canes have been enhanced with sensors and vibration feedback to detect obstacles. Some smart canes also incorporate GPS for navigation assistance [7].
- **Wearable Devices** : Smart glasses and wearable devices equipped with cameras and sensors can provide real-time feedback to visually impaired users. These devices often use object detection and depth estimation to help users avoid obstacles [8].
- **Mobile Applications** : Several mobile applications have been developed to assist visually impaired individuals. Examples include:
  - **Seeing AI** : Developed by Microsoft, Seeing AI uses AI to describe the world to visually impaired users. It can recognize text, objects, and people [9].
  - **Be My Eyes** : This app connects visually impaired users with volunteers who provide visual assistance through live video calls [10].

## 2.4 Challenges and Limitations

While existing technologies have made significant progress, several challenges remain:

- **Real-time Performance** : Achieving real-time performance on mobile devices with limited computational resources is a significant challenge [11].
- **Accuracy** : Improving the accuracy of object detection and depth estimation models, especially in complex environments, remains an ongoing area of research [12].
- **Accessibility** : Ensuring that assistive technologies are user-friendly and accessible to individuals with varying levels of visual impairment is crucial [13].

## Chapter 3: System Analysis & Design

### 3.1 Comparison of Existing Applications with my project with merits and demerits

| Feature/Aspect        | SeeU  | Seeing AI (Microsoft)   | Be My Eyes  |
|-----------------------|---|---|---|
| Core Functionality    | Real-time Object Detection, Distance Estimation, Voice Querying, LLM-based Description, TTS, Wake Word Activation             | Object Detection, Scene Description, Short/Document Reading, Product ID, Person/Currency/Color/Light Detection, TTS | Live Video Assistance by Volunteers, Specialized Help   |
| AI Capabilities       | Object Detection (YOLO), Distance Estimation (MiDaS), LLM for Description   | Object Recognition (Proprietary), Scene Understanding, OCR  | Limited (Relies on human vision)  |
| Voice Interaction     | Wake Word ("Jarvis"), Voice Querying for Information  | Primarily for describing detected elements, less conversational   | Manual initiation of calls  |
| Distance Estimation   | Explicit estimation of distance to detected objects   | Implicit in scene descriptions, less direct/detailed  | No direct distance estimation   |
| Human Assistance      | None (Autonomous AI-driven)   | None (Autonomous AI-driven)   | Core feature: relies on sighted volunteers  |
| Range of Features     | Currently focused, potential for expansion  | Comprehensive suite of diverse features   | Broad range of tasks through human assistance   |
| Platform Availability | (Potentially Cross-Platform via Flutter)  | Primarily iOS   | iOS and Android   |
| Offline Capabilities  | Potential for some (depending on model choices)   | Some features may work offline  | Requires network for live calls   |
| Learning/Improvement  | Through model training and development  | Continuous development by Microsoft   | Relies on volunteer network growth and training   |
| Strengths/Merits      | Integrated advanced AI, Conversational interaction, Customizable wake word, Potential for personalized assistance             | Wide range of features, Mature and reliable, Strong accessibility focus   | Real-time human assistance, Versatile use cases, Global network                                       |
| Weaknesses/Demerits   | Early development stage, Potential performance issues, Accuracy/reliability dependent on models, Backend dependency (current) | Primarily iOS, Less emphasis on conversational interaction, Distance estimation might be less direct                | Reliance on volunteer availability, Privacy considerations, Not fully autonomous, No direct AI vision |

## 3.2 Project Feasibility Study

A feasibility study assesses the viability of the SeeU project from technical, operational, and economic perspectives.

### 3.2.1 Technical Feasibility

- **Technologies Used:** The app leverages well-established technologies such as TensorFlow Lite, PyTorch Mobile, and FastAPI, ensuring robust technical feasibility.
- **Real-Time Processing:** The combination of YOLOv8 for object detection and MiDaS for depth estimation ensures accurate and real-time processing of camera frames.
- **Scalability:** The backend is designed to handle multiple users simultaneously, with the potential to scale using cloud-based solutions.

### 3.2.2 Operational Feasibility

- **User-Friendly Interface:** The app is designed with accessibility in mind, ensuring that visually impaired users can easily navigate and use the app.
- **Training and Support:** Comprehensive documentation and training materials will be provided to users and administrators for seamless adoption.
- **Low Learning Curve:** The app's intuitive design ensures minimal learning curve for new users.

### 3.2.3 Economic Feasibility

- **Cost-Effective Development:** The use of open-source technologies reduces development costs.
- **Monetization Strategies:** The app can be monetized through partnerships with organizations supporting visually impaired individuals or through premium features such as advanced object recognition.
- **Long-Term Benefits:** The app provides significant long-term benefits by improving the quality of life for visually impaired users, making it a worthwhile investment.

### 3.3 Requirement Gathering

**Functional Requirements:**

- Real-time object detection using the smartphone camera.
- Distance estimation for detected objects.
- Audio feedback for detected objects and their distances.
- Integration with a backend server for processing.

**Non-Functional Requirements:**

- Low-latency performance for real-time feedback.
- Compatibility with Android 9+ devices and iOS. User-friendly interface for visually impaired users

### 3.4 Project Timeline chart

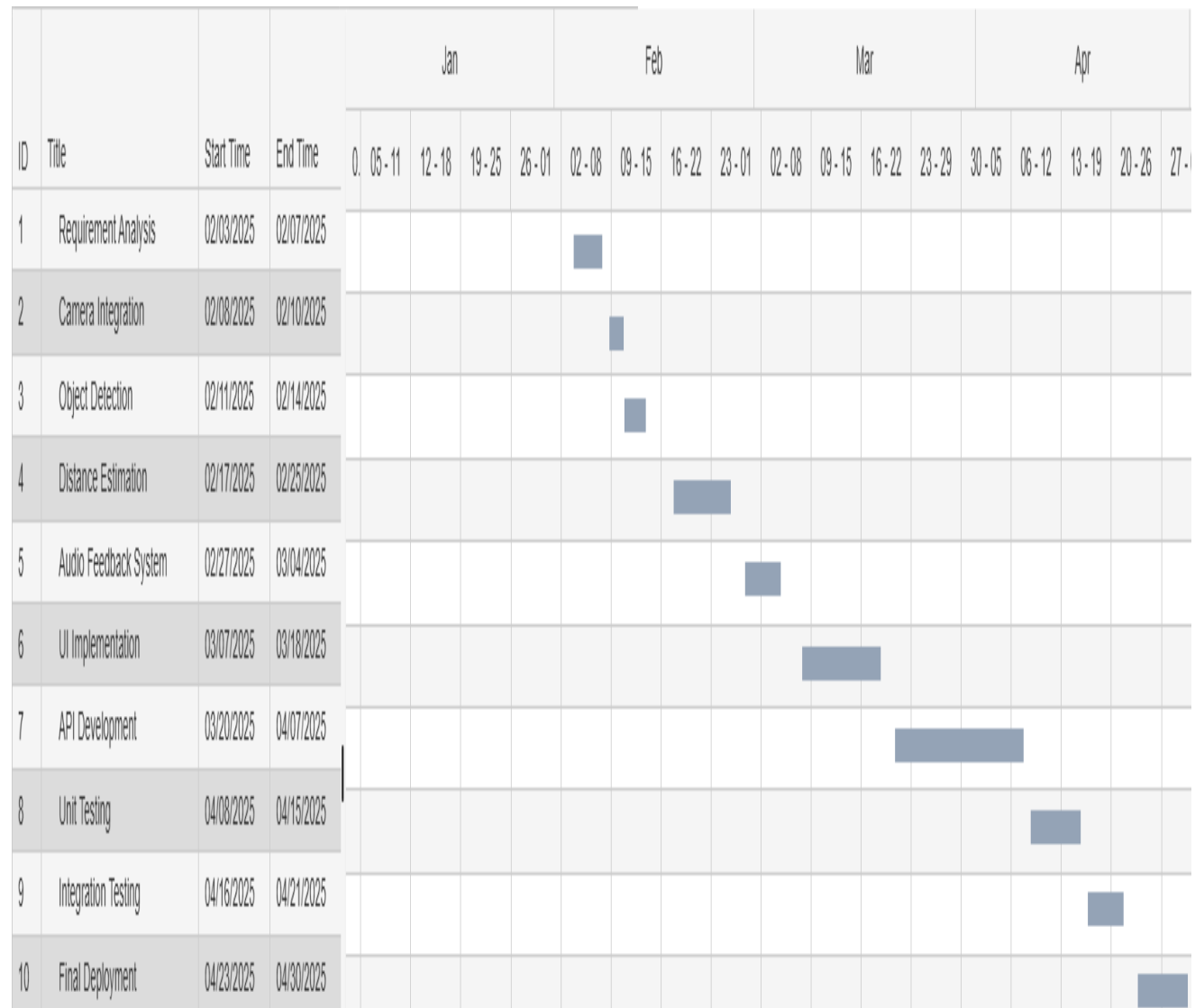


Figure 1. Gantt Chart

### 3.5 Detailed Modules Description

#### 3.5.1 User Interface (UI)

- Camera View
  - Description: Real-time camera feed displaying the user's surroundings.
  - Features:
    - Live preview from the device's rear camera.
    - Overlay-free design to avoid visual clutter
- Wake Word Activation Panel

- Description: Visual indicator for voice command readiness.
- Features:
  - Status bar showing wake word detection state (e.g., "Ready for 'Jarvis'").
  - Color-coded feedback
    - Green: Listening for wake word ("Jarvis").
    - Red: Wake word inactive (with a restart button).
- Detected Objects List
  - Description: Scrollable list of detected objects and distances.
  - Features:
    - Card-based layout showing:
      - Object name (e.g., chair).
      - Distance (e.g., 2.5m away) or unknown if depth estimation fails.
    - Auto-updates every second.

### 3.5.2 Object Detection Module

The object detection module in SeeU is responsible for detecting objects in real-time using a **YOLOv8** model trained on a custom dataset.

YOLOv8 (You Only Look Once, Version 8) is an anchor-free, fully convolutional neural network that predicts bounding boxes and class probabilities in a single forward pass through the network.

#### Model Details:

- **Input:** Preprocessed RGB frames (640x480 resolution).
- **Architecture:** YOLOv8 leverages advanced backbone structures and a decoupled head for classification and regression tasks.
- **Parameters Tuned:**
  - **Confidence threshold:** 0.3 (Only detections with probability above 30% are kept. Lower threshold increases detections but may introduce false positives.)
  - **IoU threshold for NMS:** Default 0.5 (Controls how overlapping boxes are merged — higher IoU reduces multiple detections of the same object).



- **Performance:** Achieved mAP@0.5 of 92% on validation dataset.

**Working:**

1. Camera captures the frame.
2. Frame is encoded in Base64 and sent to FastAPI backend.
3. YOLOv8 model processes the frame and outputs bounding boxes, confidence scores, and class labels.
4. Low-confidence detections ( $<0.3$ ) are filtered out.
5. Post-processing (Non-Maximum Suppression) is applied to eliminate duplicate boxes.

**Influence of Parameters:**

- Increasing confidence threshold improves precision but reduces recall.
- Lowering IoU threshold in NMS increases sensitivity but may lead to more duplicate detections.

**YOLO Model Integration**

- Model Used: Custom-trained YOLOv8 (besttrain.pt).
- Input : Base64-encoded image from the camera feed.
- **Output:**

```
{  
  "object": "chair",  
  "confidence": 0.92, # Filtered to >0.3  
  "bbox": [x1, y1, x2, y2] # Bounding box coordinates  
}
```

**3.5.3 Distance Estimation Module**

The distance estimation module uses **MiDaS** (Monocular Depth Estimation) to predict the relative depth of objects within a 2D image.

**Technical Working:**

- **Model:** MiDaS Small v2.1 from Intel ISL, a lightweight transformer-based depth estimation model optimized for mobile devices.
- **Input:** RGB frame of the captured image.
- **Processing:**
  - Frame is converted to RGB format using OpenCV.
  - Frame is transformed using MiDaS preprocessing pipeline (resizing, normalization).
  - The frame is passed through the MiDaS model to produce a dense depth map.
- **Depth Extraction:**
  - For each detected object, the bounding box region is extracted from the depth map.
  - The average depth value in this region is calculated.

**Parameters Configured:**

- **Calibration Factor:** 114
  - Used to map the average depth values into approximate real-world distances.
  - Distance is calculated as:

```
distance = CALIBRATION_FACTOR / avg_depth # CALIBRATION_FACTOR = 114
```

**Impact of Parameters:**

- A higher calibration factor results in larger estimated distances.
- The resolution of the input frame affects the sharpness of the depth prediction — lower resolution may result in noisier depth maps.

**Sample output :**

```
"detections": [  
  {  
    "object": "chair",  
    "confidence": 0.92,  
    "bbox": [120, 200, 300, 400],  
    "distance": 2.5  
  },  
  {  
    "object": "person",  
    "confidence": 0.87,  
    "bbox": [50, 150, 180, 350],  
    "distance": 1.8  
  }  
]
```

### 3.5.4 Audio Feedback System

The audio feedback module uses the **Coqui TTS XTTS-v2** engine to generate real-time voice outputs that describe detected objects and their estimated distances.

#### Technical Working:

- **Model:** Coqui TTS XTTS-v2
- **Input:** Text describing detected objects and distances.
- **Processing:**
  - Text and optional speaker embedding are passed to the TTS engine.
  - Output is generated as a WAV audio file.
- **Transmission:** Audio is encoded in Base64 and sent to the Flutter frontend, where it is decoded and played.

#### Parameters Configured:

- **Language:** 'en' (English)
- **Speaker WAV:** Optional Base64 encoded audio sample for voice cloning.

**Impact of Parameters:**

- Selecting the language affects pronunciation and naturalness of the audio.
- Providing a custom speaker WAV personalizes the voice but slightly increases generation time.

**Wake Word Detection**

- Library : Porcupine SDK (porcupine\_flutter) with "Jarvis" keyword.
- Activation Flow : 1. "Jarvis" detected → 2. Stops YOLO timer → 3. Starts speech recognition.

**3.5.5 LLM Processing Module**

The LLM (Large Language Model) processing module uses a local deployment of **LLaVA-v1.5-7b** to intelligently answer user queries related to the surroundings.

**Technical Working:**

- **Model:** LLaVA v1.5 7B — a multimodal model combining vision and language understanding.
- **Input:**
  - Detected object labels and distances (text).
  - Captured image (Base64 encoded).
  - User query (text).
- **Prompt Construction:**
  - A system prompt is built, describing the detections and context.
  - The user query is appended.
- **Inference:**
  - Both image and text are passed to the model.
  - Model generates a short, warm, casual response.
- **Output:**
  - The response text is sent back to the frontend for speech generation.

**Parameters Configured:**

- **Max Tokens: 150**
  - Limits the length of the LLM's response to keep it short and relevant.

**Impact of Parameters:**

- Increasing max tokens would generate longer answers but would increase latency.
- Using both image and text inputs improves the relevance and accuracy of the answer.

**Workflow**

1. User Query:
  - Voice → STT → Text ("What's near me?").
2. Context Preparation:
  - Combines YOLO detections ("chair (2.5m)") + query → LLM payload.
3. LLM Inference:

```
llm_payload = {  
  "model": "llava-v1.5-7b",  
  "messages": [  
    {"role": "system", "content": "You are a vision assistant..."},  
    {"role": "user", "content": ["Detections: chair (2.5m)", "What's near me?"]}  
  ]  
}
```

4. Output:
  - Short, actionable response → TTS (e.g., "A chair is 2 meters to your left").

### 3.6 Project SRS

#### 3.6.1 Use Case Diagram

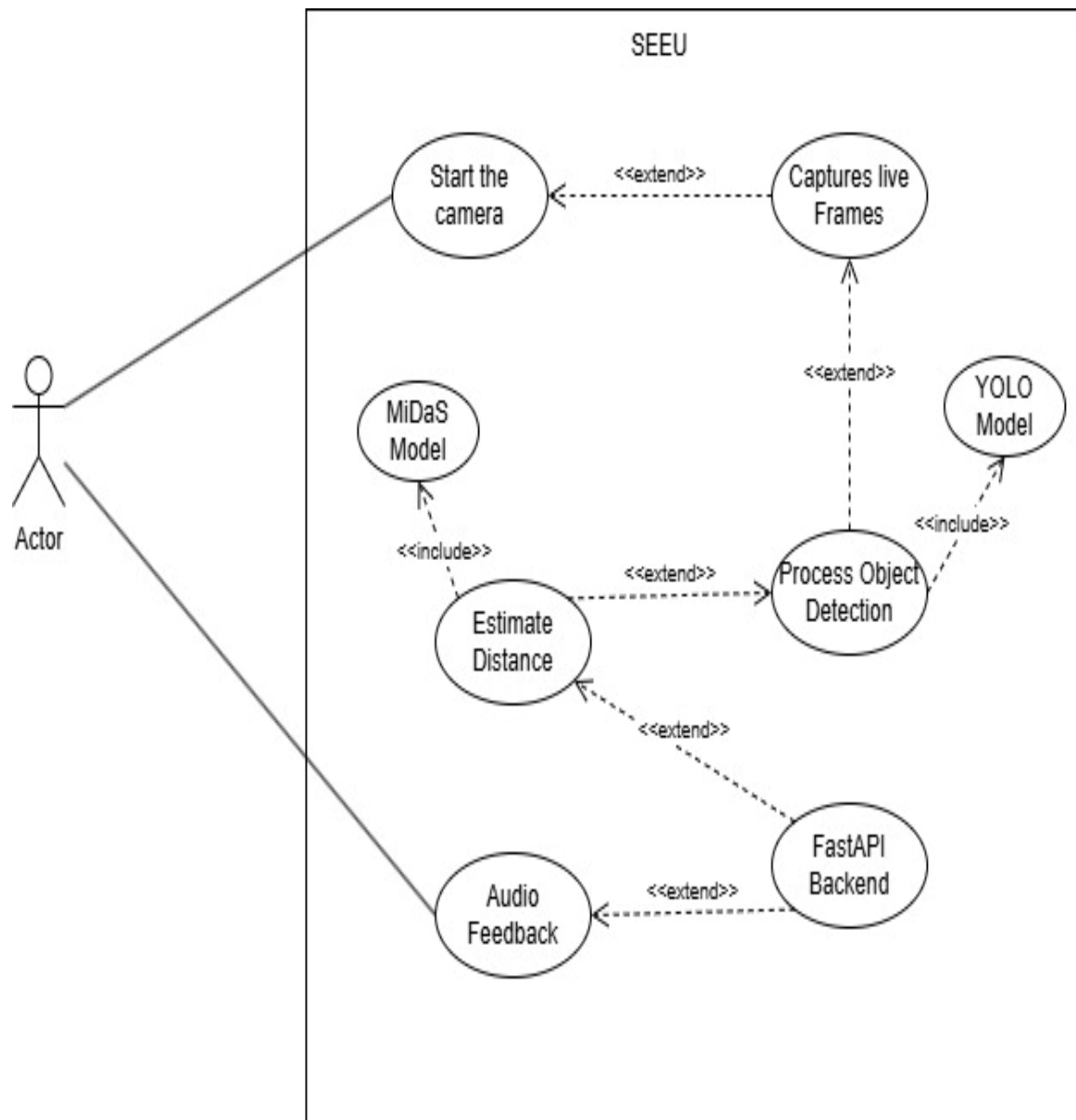


Figure 2. Use case diagram

### 3.6.2 Data Flow Diagrams

- DFD Level 0

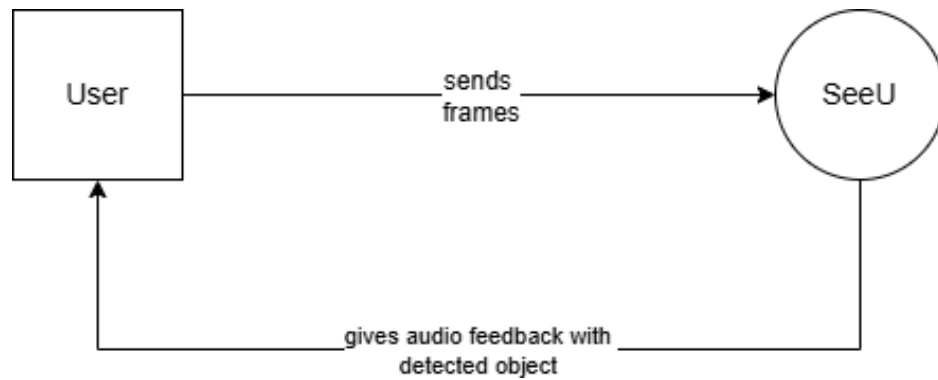


Figure 3. DFD level 0

- DFD Level 1

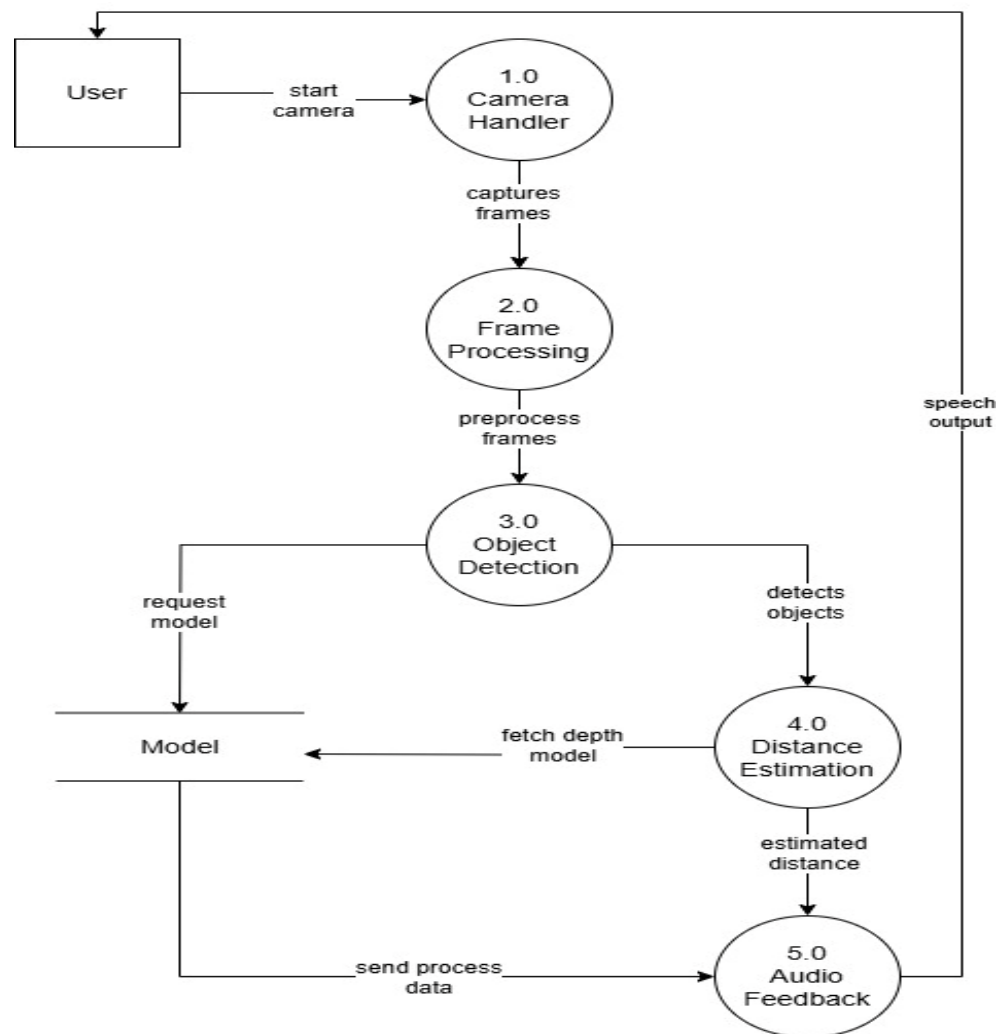


Figure 4. DFD level 1

### 3.6.3 Class Diagram

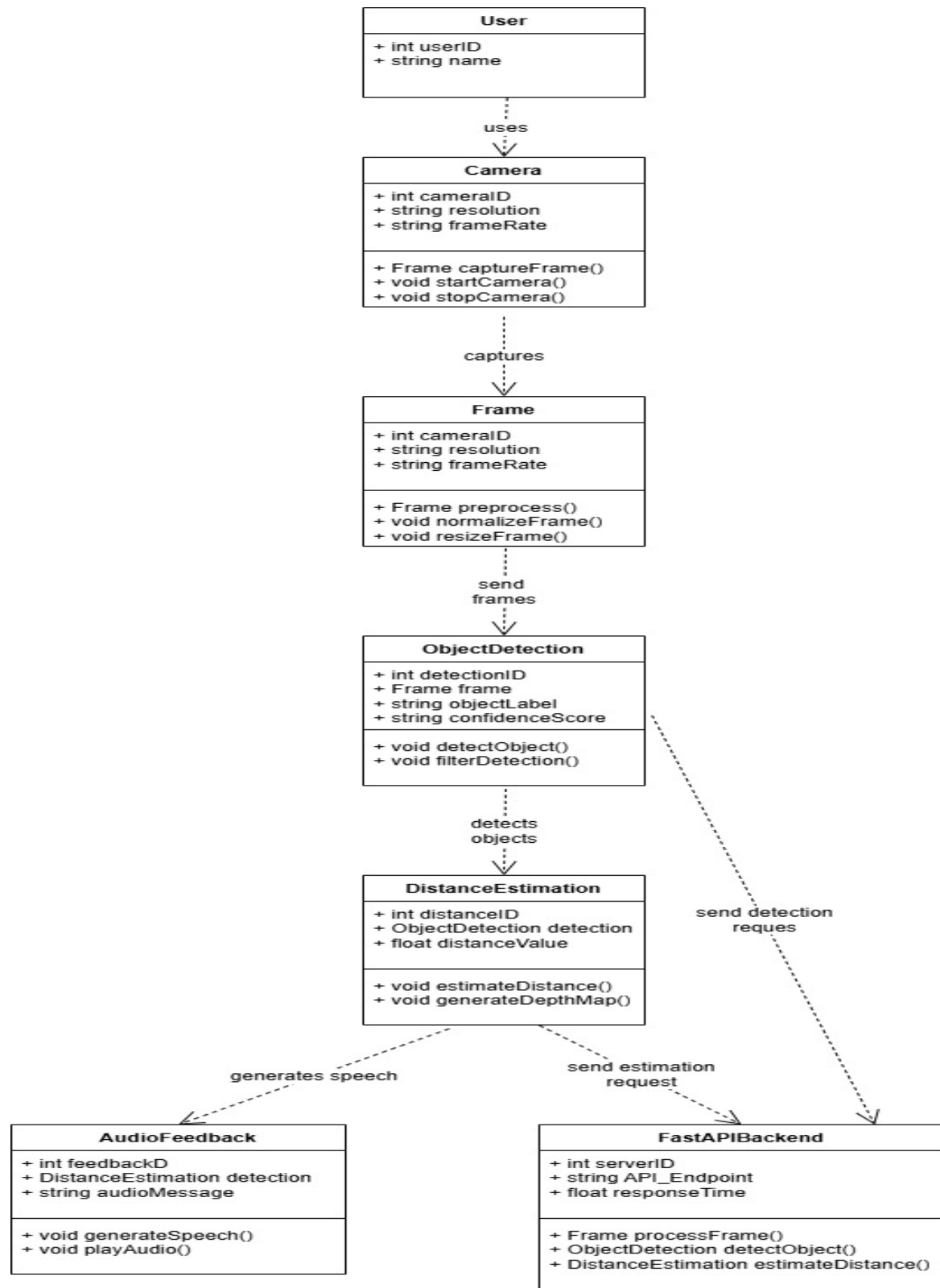


Figure 5. Class Diagram



### 3.6.4 Entity Relationship Diagrams

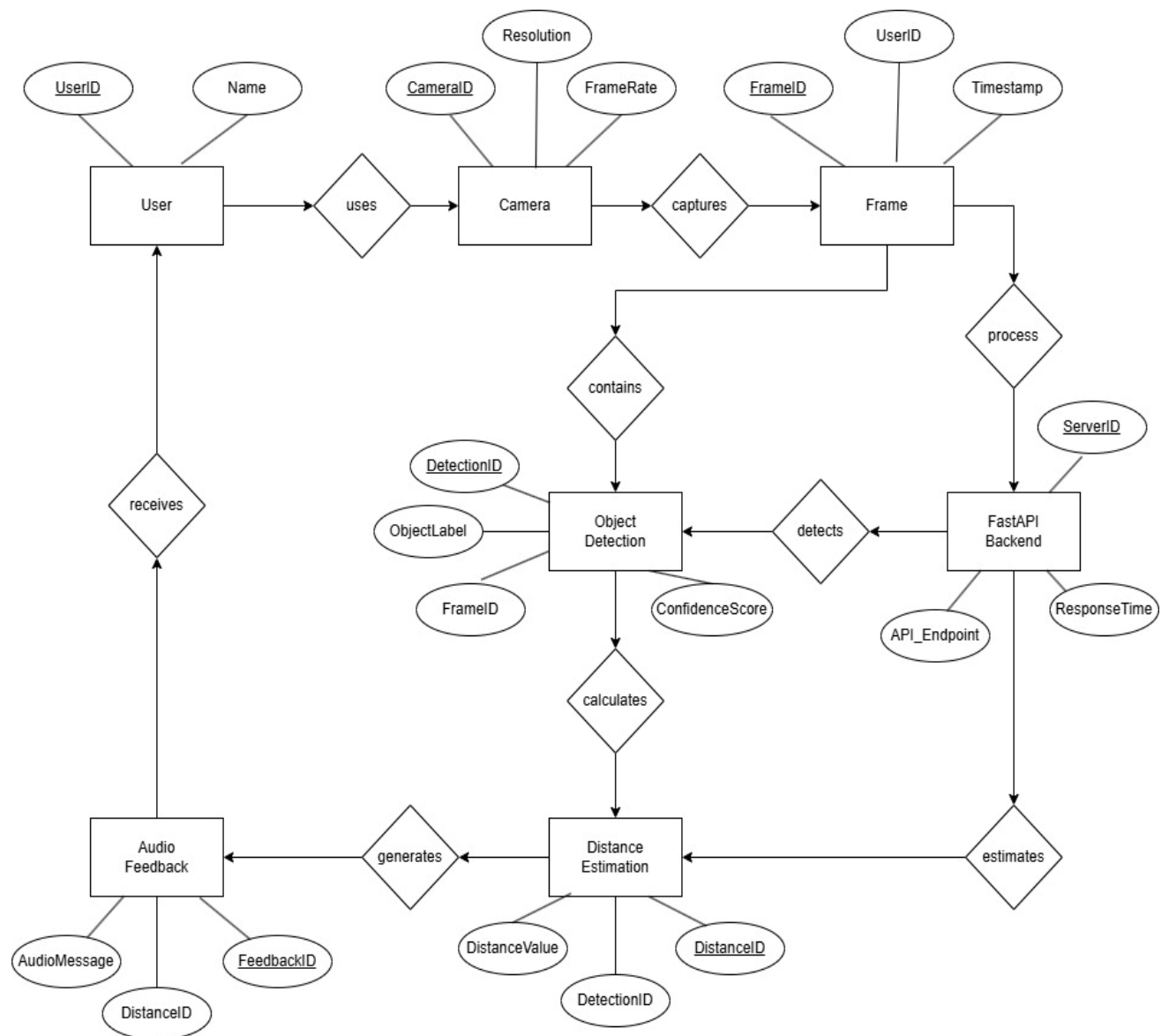


Figure 6. ER Diagram

### 3.6.5 Flow Chart

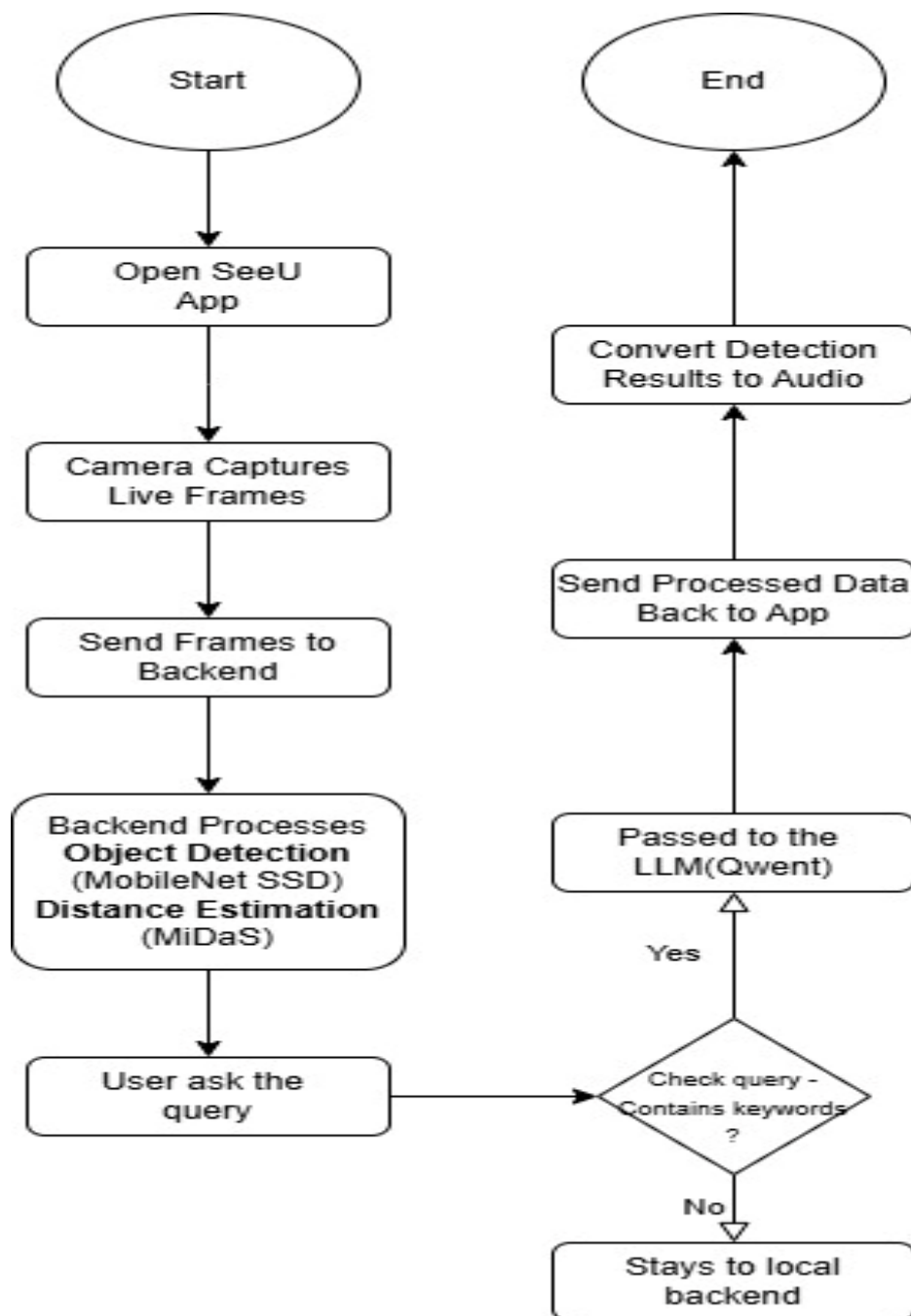


Figure 7. Flow Chart

## Chapter 4: Implementation and Testing

### 4.1 User Interface Snapshots

A Simple UI with Camera preview and the “Detected Object” label where the details of the detected objects displays with its distance. And the symbol of the “Ready for Jarvis” which shows that you can say Jarvis to open the voice recorder.

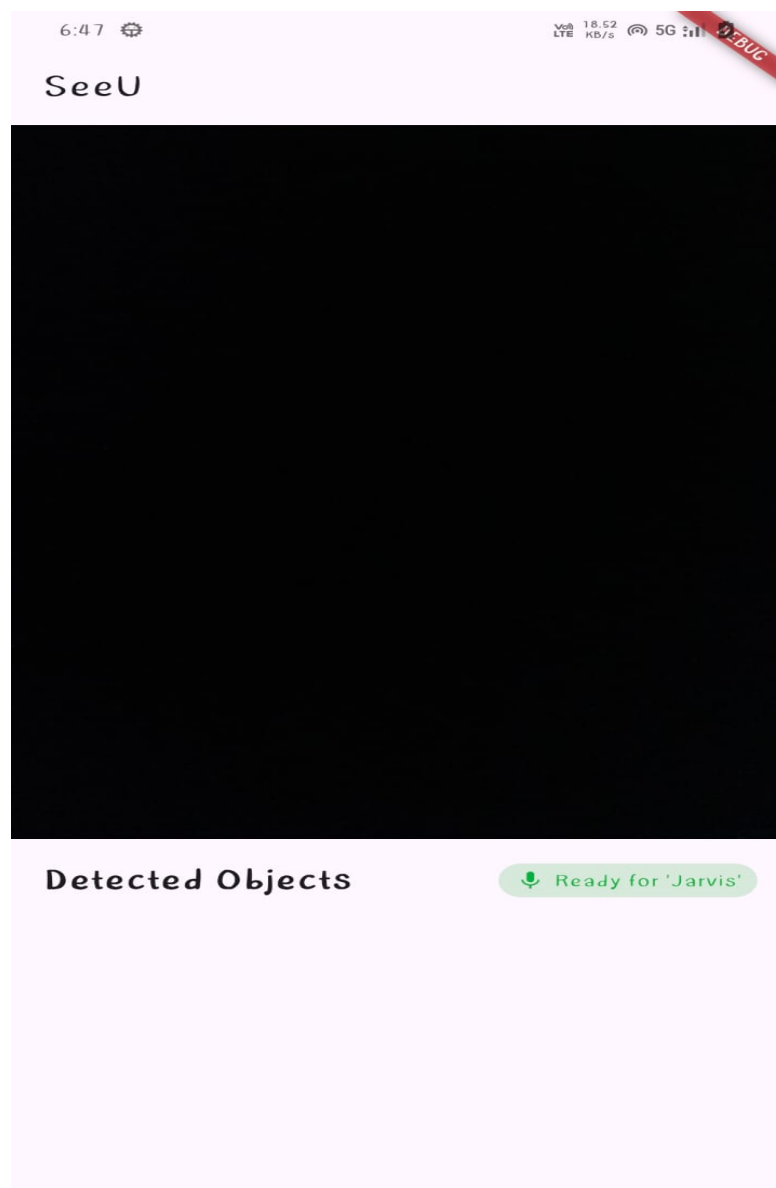


Figure 8. Simple UI

When you say “Jarvis” it will open the voice recorder where you can ask the question.

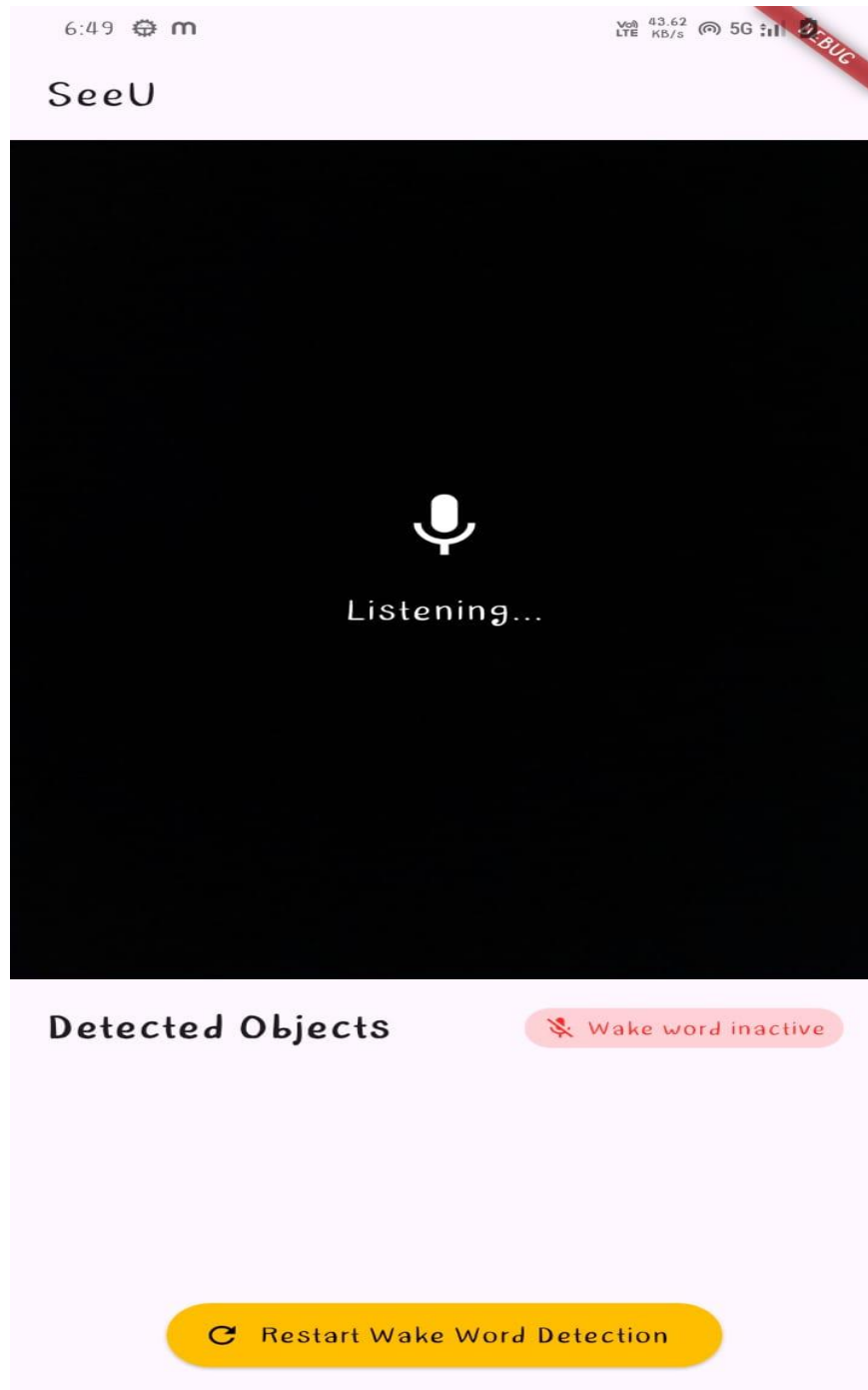
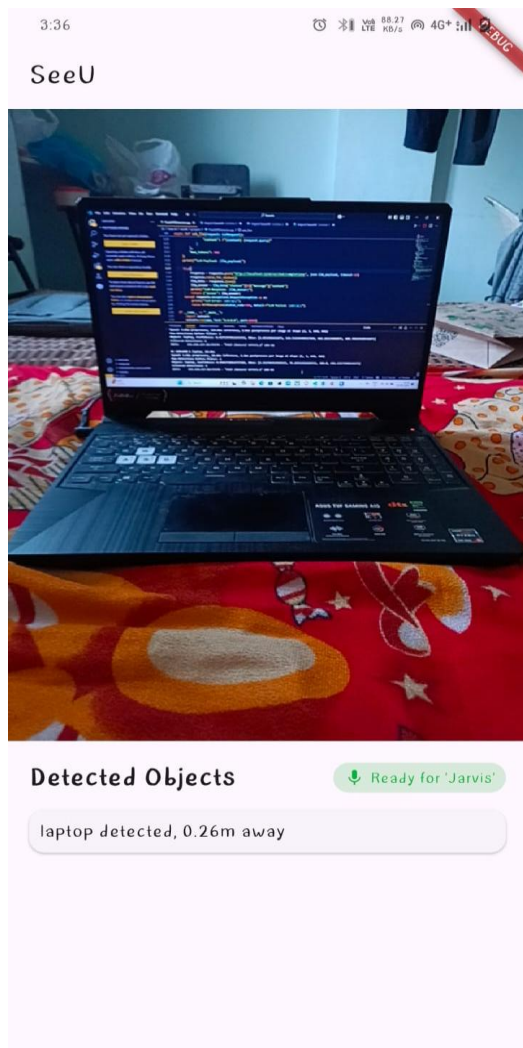


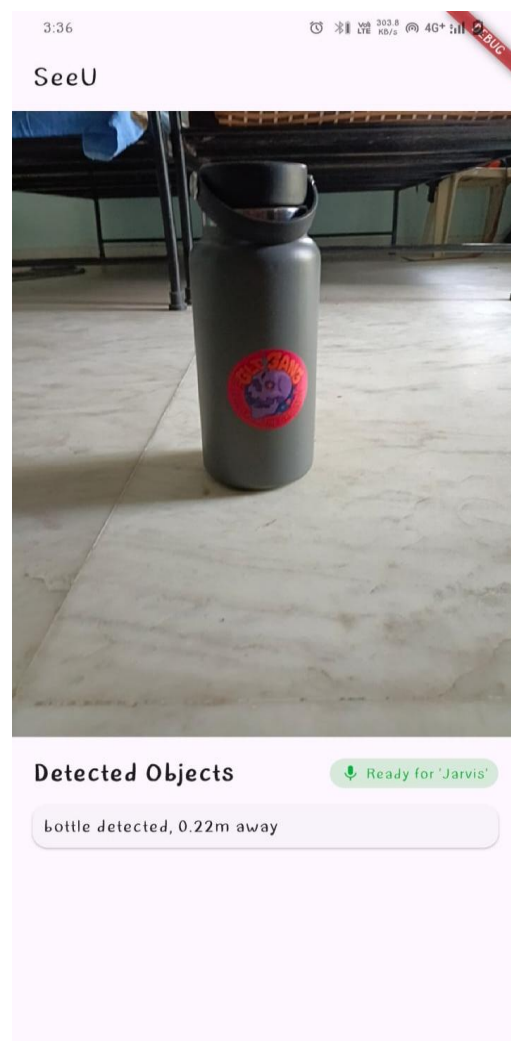
Figure 9. Listening UI

## 4.2 Testing

- Object Detection with Distance Estimation Testing



**Figure 10. Laptop Detected**



**Figure 11. Bottle Detected**

- In Figure 9, we can see that the laptop is detected with its distance i.e 0.26m away.
- In Figure 10, we can see that the detected object is water bottle and its 0.22m away.

- LLM Response Testing

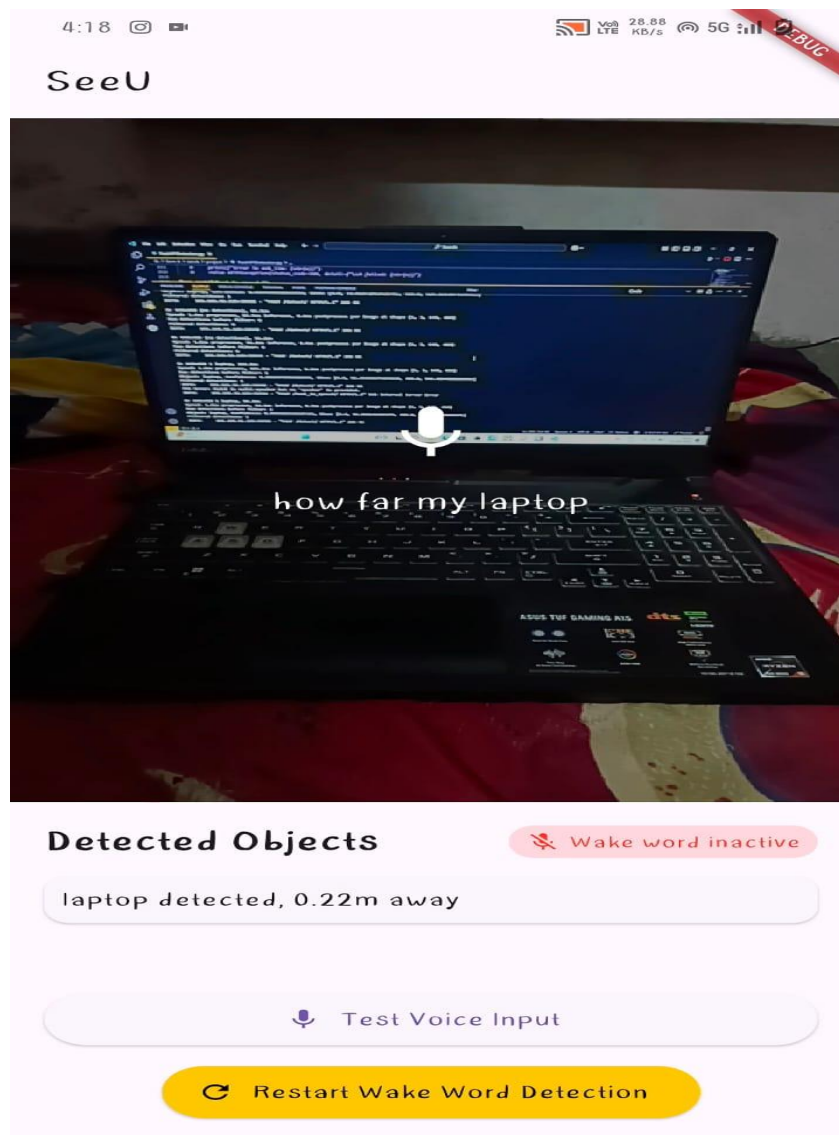


Figure 12. Question asked

- In Figure 11, I asked the “How far my laptop” and it is passed to the llava-v1.5-7b model with the image and the details of the detected object with its distance

```
"content": "Hey there! I see a laptop sitting on a bed, and it's displaying some code or text. It looks like you might be working on a project or studying computer science. The laptop is positioned at about 0.24 meters away from the camera, which gives us a sense of scale for the image. If you need help with anything related to programming or technology, feel free to ask!"
}
```

Figure 13. LLM Response

- In Figure 12, is the response the LLM model that is llava-v1.5-7b is generated

- Full app testing

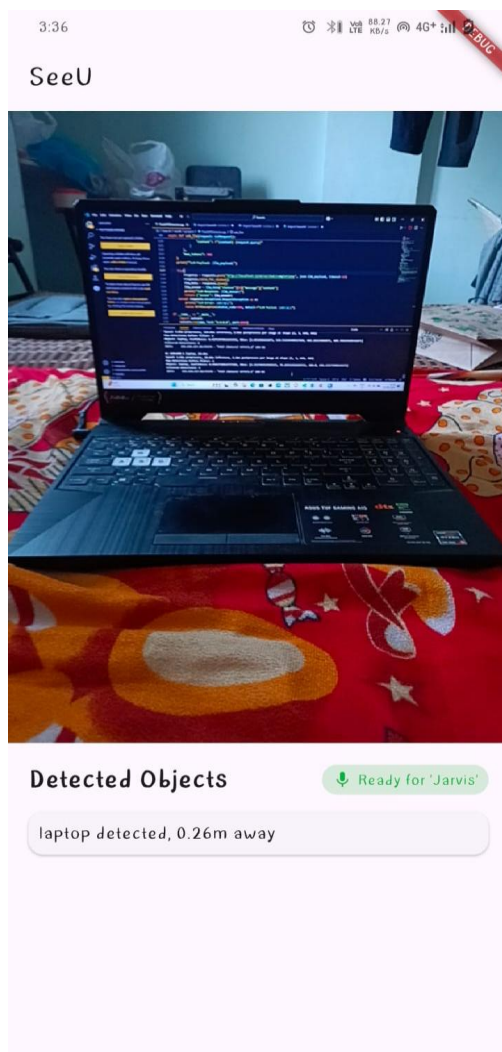


Figure 10. Laptop Detected

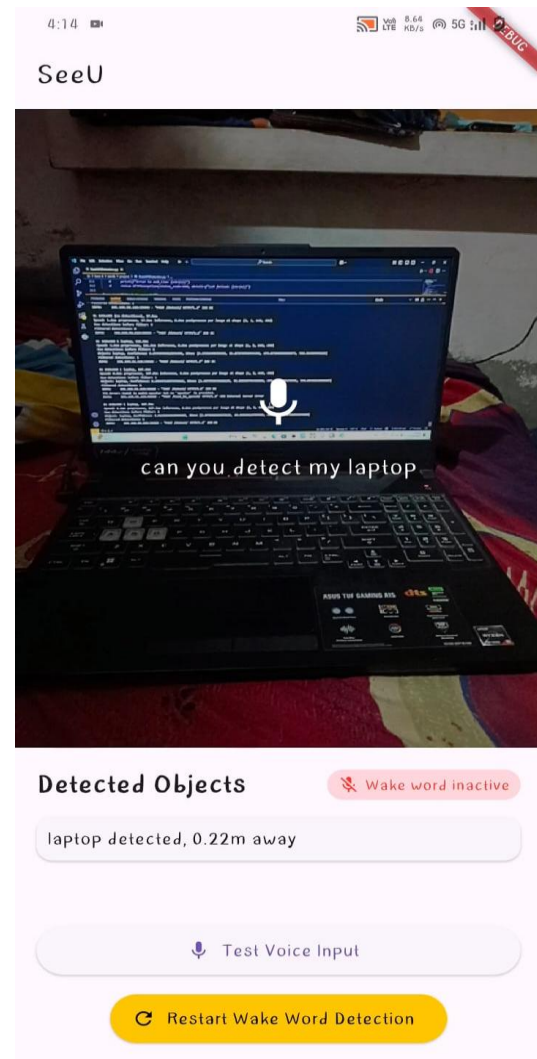


Figure 14. Asked question

- In Figure 9, we can see the laptop is detected and it is 0.26m away from it.
- When I say “Jarvis” the voice recorder is opened while speaking “I am listening”
- In Figure 10, I asked the question “Can you detect my laptop”

```
0: 640x480 1 laptop, 94.2ms
Speed: 0.9ms preprocess, 94.2ms inference, 0.7ms postprocess per image at shape (1, 3, 640, 480)
Raw detections before filter: 1
Object: laptop, Confidence: 0.8976629376411438, BBox: [0.8494873046875, 92.94847106933594, 480.0, 551.473388671875]
Filtered detections: 1
INFO: 192.168.52.115:39924 - "POST /detect/ HTTP/1.1" 200 OK
Sending to LLM: Detections: laptop (0.24m). can you detect my laptop
LLM Response:
Hey there! It looks like your laptop is open to a screen with lots of text on it, possibly some kind of program or document. The keyboard is visible in the foreground and appears to be black. I'm not sure what exactly you are working on, but it seems to be an important task that requires focus and attention. If you need any help or have any questions, feel free to ask!
INFO: 192.168.52.115:39926 - "POST /ask_llm/ HTTP/1.1" 200 OK
```

**Figure 15. Log of backend**

- In Figure 14, is the log of the backend. Where we can see the content of LLM Response which the output spoke and the user will listen it.



## Chapter 5: Conclusion and Future Work

### 5.1 Conclusion

The project SeeU: An Object Detection and Distance Estimation App for Assisting Visually Impaired Individuals successfully demonstrates the integration of deep learning, computer vision, and speech synthesis to provide real-time environmental feedback. By combining wake-word activation ("Jarvis"), multilingual TTS, and an LLaVA-based LLM, SeeU provides an accessible, privacy-preserving alternative to cloud-dependent solutions.

By combining object detection, monocular depth estimation, large language models, and advanced text-to-speech systems, SeeU offers a comprehensive and personalized experience for visually impaired users to navigate their surroundings more safely and independently.

### 5.2 Limitations

While promising, SeeU has areas for growth:

- Latency: 4–5 seconds is functional but not ideal for real-time navigation.
  - Solution : Model quantization (e.g., TensorRT) could reduce inference time.
- Hardware Constraints: Performance drops on low-end smartphones.
  - Solution : Offer a "Lite Mode" with smaller models (e.g., YOLO-Nano).
- Wake Word False Positives: Background noise sometimes triggers accidental activations.
  - Solution : Integrate noise suppression (e.g., RNNoise).

## 5.2 Future Work

While **SeeU** already offers a robust set of features, there are several areas for future enhancement that can further improve its functionality and user experience.

### 1. Performance Optimization:

- Implement TensorRT quantization to reduce latency to <2 seconds.
- Develop a "Lite Mode" using YOLO-Nano for low-end devices.

### 2. Feature Expansion:

- Contextual Awareness: Use LLM fine-tuning to describe object relationships (e.g., "Chair blocking the door").
- Offline Mode: Bundle quantized models to enable functionality without internet.

### 3. User-Centric Refinements:

- Add other languages TTS support to serve regional users.
- Integrate RNNoise to improve wake-word accuracy in noisy environments.

## Chapter 6. References

1. Howard, A. G., et al. (2017). "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications."
2. Redmon, J., & Farhadi, A. (2018). "YOLOv8: An Incremental Improvement." arXiv preprint arXiv:1804.02767
3. Ren, S., He, K., Girshick, R., & Sun, J. (2015). "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks." Advances in Neural Information Processing Systems.
4. Ranftl, R., Bochkovskiy, A., & Koltun, V. (2021). "Vision Transformers for Dense Prediction." arXiv preprint arXiv:2103.13413.
5. Szeliski, R. (2010). "Computer Vision: Algorithms and Applications." Springer Science & Business Media.
6. Chen, Y., et al. (2020). "LiDAR-based 3D Object Detection: A Survey." arXiv preprint arXiv:2007.13132.
7. Khan, M., et al. (2019). "Smart Cane: A Depth-Based Solution for the Visually Impaired." Sensors.
8. Bai, J., et al. (2017). "Wearable Assistive Devices for Visually Impaired: A Review." IEEE Sensors Journal.
9. Microsoft. (2017). "Seeing AI: An Intelligent Assistant for the Visually Impaired." Microsoft Research.
10. Be My Eyes. (2015). "Be My Eyes: Connecting Blind and Low-Vision People with Sighted Volunteers." Mobile Application Review.
11. Zhang, Z., et al. (2021). "Optimizing Real-Time Object Detection for Mobile Devices." IEEE Transactions on Mobile Computing.
12. Liu, W., et al. (2016). "SSD: Single Shot MultiBox Detector." European Conference on Computer Vision.
13. Pawluk, D., et al. (2020). "Accessibility Challenges in Assistive Technologies for the Blind." ACM Transactions on Accessible Computing.