Building Al-Powered Solution for Assisting Visually Impaired Individuals

Subtitle

: Leveraging Generative AI and Computer Vision for Enhanced Accessibility

Presented by

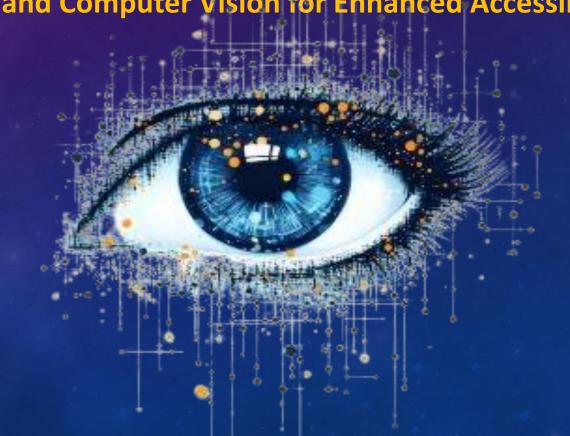
: Satve Pukhraj

Organization

: Innomatics Research Labs

Date

: 24/11/2024



Introduction

The goal of this project is to develop an Al-powered solution that assists visually impaired individuals in perceiving and interacting with their surroundings more effectively. Visually impaired people often face significant challenges in understanding their environment, reading visual content, and performing tasks that typically rely on sight. This solution aims to address these issues by leveraging advanced technologies such as generative AI, computer vision, and text-to-speech synthesis. By processing realtime images, the system can generate descriptive textual outputs that help users understand the content of their surroundings, convert text from images into audible speech, and detect objects or obstacles to ensure safe navigation. Furthermore, it provides personalized assistance for various daily tasks, such as recognizing items, reading labels, or offering context-specific guidance. The solution's primary objective is to enhance accessibility and independence for visually impaired individuals, allowing them to navigate their environment with greater confidence and ease.



Problem Statement

Challenges Faced by Visually Impaired Individuals:

- Navigating the Environment
- Reading Visual Content
- Interacting with Objects

- : Lack of real-time perception of objects and obstacles in the surroundings.
- : Difficulty in reading printed text, labels, or signs.
- : Trouble recognizing and interacting with items in their environment, especially in public spaces or at home.

❖ Impact:

- Visually impaired individuals often face increased dependency on others.
- Limited access to information that is presented visually in daily life.
- Safety concerns while navigating unfamiliar or cluttered spaces.

***** Key Needs:

- Real-time assistance for understanding their surroundings.
- Seamless access to visual content through text-to-speech and object recognition.
- Personalized guidance for tasks like reading labels, identifying items, and detecting obstacles.

Problem Analysis

Accessibility Issues:

- A large portion of daily information is conveyed visually (e.g., reading, navigation, safety signals).
- Visual content in public spaces (e.g., street signs, labels in stores) is often inaccessible.

Existing Solutions:

- Screen readers and OCR tools (like Tesseract) help, but they are limited in scope.
- Existing AI models provide scene analysis, but they don't integrate multiple aspects like real-time obstacle detection and personalized task guidance.

Need for an Integrated Solution:

- There is a need for a comprehensive solution that:
 - Provides real-time scene understanding through generative AI.
 - Allows for **OCR and text-to-speech conversion** to read and interpret text.
 - Detects objects and obstacles for safe navigation.
 - Offers personalized assistance for various daily tasks.

Solution Overview

- ❖ Objective: To develop an integrated Al-powered solution that:
 - Generates descriptive textual output for images, helping visually impaired users understand the scene.
 - Converts text to speech from images, providing accessibility to written content.
 - **Detects objects and obstacles** in real-time, ensuring safe navigation.
 - Assists in daily tasks by recognizing objects, and labels, and providing context-based help.

❖ How it Works:

- User Uploads Image: The user uploads an image of their surroundings or a document.
- **Generative AI & CV**: Al interprets the image and detects objects, generating text descriptions.
- OCR: Text is extracted from the image and converted to speech for accessibility.
- Audio Output: The system reads the descriptions and text aloud using a text-to-speech engine.

Features

1.Real-Time Scene Understanding:

- 1. Using generative AI to interpret and describe the content of an image in natural language.
- 2. Helps users understand the surroundings, including the layout of objects, people, and possible hazards.

2.OCR & Text-to-Speech Conversion:

- 1. OCR extracts text from images (documents, labels, etc.).
- 2. The extracted text is converted into speech for easy access.

3.Object and Obstacle Detection:

- 1. Real-time identification of objects and obstacles in the user's environment.
- 2. Alerts the user to avoid obstacles or informs them about the presence of objects (e.g., furniture, vehicles).

4.Personalized Assistance:

- 1. Context-aware guidance based on the specific task (e.g., reading a document, identifying items).
- 2. Helps with daily activities, like navigating a room, recognizing items on a table, or finding a specific book in a shelf.

1) Real-Time Scene Understanding

To generate descriptive textual output that helps visually impaired individuals understand the content of an uploaded image.

Process:

- **Generative Al Models**: The system uses advanced Al models, such as **CLIP** (Contrastive Language-Image Pre-training) or **GPT-4**-based models, which combine image understanding and natural language generation. These models analyze the visual content of an image and generate detailed descriptions in natural language.
- Scene Analysis: The system evaluates the objects, people, text, and other elements within the image. It recognizes relationships between elements, such as a person standing next to a bench or a dog walking on a sidewalk.
- **Descriptive Output**: The result is a verbal description that conveys the content of the scene, which can include details such as the environment (e.g., "a park with a walking path"), the people or animals present, and key activities taking place.

Example:

• A picture of a busy street might generate the following description: "The image shows a busy street with several pedestrians walking along the sidewalk. There is a yellow taxi on the left side of the road, and a streetlamp is visible on the corner."

2) Text-to-Speech Conversion for Visual Content

To extract text from uploaded images (using Optical Character Recognition or OCR) and convert it into audible speech, making written content accessible.

Process:

- OCR (Optical Character Recognition): The system employs Tesseract or other OCR tools to detect and extract any text present within the image. This includes printed text, handwritten text, labels, signs, or any other textual content.
- **Text Extraction**: After extracting the text from the image, it is cleaned up (removing noise, fixing misreads) to ensure high accuracy.
- Text-to-Speech (TTS): The extracted text is then passed through a TTS engine like pyttsx3 or Google Cloud Text-to-Speech, which converts the text into natural-sounding speech.
- Audio Output: The system outputs the speech through speakers or headphones, providing the user with a verbal description of the text content.

Example:

• A picture of a restaurant menu might generate the following output: "The menu includes a Caesar salad, grilled chicken sandwich, and spaghetti bolognese. Prices range from \$8 to \$15."

3) Object and Obstacle Detection for Safe Navigation

To identify objects or obstacles within an image and highlight them, helping users navigate safely and enhancing situational awareness.

Process:

- Object Detection Models: The system uses Convolutional Neural Networks (CNNs), such as YOLO (You Only Look Once) or Faster R-CNN, to detect and classify objects in real time. These models are trained on large datasets to recognize various objects, such as furniture, people, vehicles, and other common items.
- **Obstacle Detection**: For navigation, the system focuses on detecting obstacles that might be in the user's path, such as furniture, street signs, or other impediments. The system can highlight these obstacles and indicate their proximity to the user.
- Real-Time Feedback: Once obstacles are detected, the system provides the user with spoken feedback, alerting them about obstacles and guiding them to avoid or navigate around them.

Example:

 A picture of a room with furniture might generate the following description: "There is a chair directly ahead. A table is positioned to the left, and a lamp is on the right. There is an obstacle ahead, approximately three feet away."

4) Personalized Assistance for Daily Tasks

To provide task-specific guidance based on the uploaded image, such as recognizing items, reading labels, and offering context-based information.

Process:

- Image Analysis for Task Recognition: Depending on the context (e.g., reading labels, identifying items in a room), the system uses object detection, scene recognition, and task-specific algorithms to analyze the image and offer detailed feedback.
- Task-Specific Guidance: The system is trained to identify specific tasks based on the context of the image. For example, it might recognize a grocery list, identify the names of objects on a table, or assist in reading a prescription label.
- **Contextual Awareness**: The system can be personalized by the user, learning to recognize frequent objects, locations, and user preferences over time. This allows it to provide more accurate and useful assistance based on individual needs.
- **Spoken Feedback**: The user receives real-time verbal feedback, whether it's for identifying a product, reading text on a label, or providing context for a particular item in the environment.

Example:

• A user uploads an image of a kitchen counter with several items, and the system might say: "On the counter, there is a can of soup, a bottle of olive oil, and a cutting board with a knife. The soup can is labeled 'Tomato Soup,' and the olive oil is marked with an expiration date of September 2025."

Technologies Used

- Generative Al:
 - •Utilized to describe images and generate contextual insights about the scene.
 - •APIs like google. Generative ai to generate textual descriptions of scenes.
- Computer Vision & Deep Learning:
 - •PyTorch and TorchVision for real-time object detection and scene analysis.
- ❖ OCR:
 - •Tesseract to extract text from images and documents.
- **❖** Text-to-Speech:
 - pyttsx3 for converting extracted text into audio.
- Other Tools:
 - •Streamlit for building the web application interface.
 - •PIL for image processing.
 - OpenCV for real-time object detection and scene processing.
 - NumPy and Requests for handling data and APIs.

Architecture & Workflow

1.User Input: The user uploads an image through the Streamlit interface.

2.Scene Interpretation:

- 1. Generative AI models analyze the image and provide textual descriptions.
- 2. Computer vision models detect objects, people, and obstacles in the environment.

- **3.OCR Processing**: Text is extracted from the image (if applicable) using Tesseract.
- **4.Text-to-Speech**: The system converts the output text to speech using pyttsx3.
- **5. Audio Output**: The results are read aloud to the user, providing feedback.

Use Case Scenarios

Scenario 1: Navigating a New Environment:

 The user uploads an image of a street or public space, and the system describes the scene, including obstacles and possible hazards.

Scenario 2: Reading Labels and Documents:

 The user uploads an image of a document, and the system reads aloud the text present in the image.

Scenario 3: Identifying Objects:

The user uploads an image of a room or workspace, and the system identifies items like objects
on a table, books on a shelf, etc.

Challenges & Limitations

Accuracy of Scene Interpretation:

• The quality of object detection and scene understanding may vary depending on image quality, lighting conditions, and object complexity.

Real-Time Processing:

 Generating descriptions and processing text-to-speech in real-time can require significant computational resources.

OCR Limitations:

 Tesseract may struggle with non-standard fonts or distorted text, impacting the accuracy of text extraction.

Safety Concerns:

 Incorrect object detection or missed obstacles could pose safety risks, particularly in dynamic, realworld environments.

Future Scope

•Multimodal Inputs: Integrating voice commands and other sensory inputs (e.g., haptic feedback) for a more intuitive user experience.

•Improved Object Detection: Using more advanced AI models for better accuracy in recognizing objects and obstacles in various environments.

•Context-Aware Assistance: Personalizing the system further to provide specific guidance based on user preferences, needs, and habits.

•Accessibility Integration: Expanding to support integration with existing accessibility tools and platforms (e.g., screen readers, smart assistants).

Conclusion

In conclusion, this Al-powered solution offers a transformative approach to assisting visually impaired individuals by leveraging technologies like Generative AI, OCR, computer vision, and text-to-speech. It provides real-time scene understanding, text-to-speech conversion, object and obstacle detection, and personalized task assistance, all aimed at enhancing accessibility, safety, and independence. By improving navigation, reading capabilities, and daily task completion, the system empowers users to interact with their environment more confidently. As the technology evolves, it holds the potential to further enhance the lives of visually impaired individuals, promoting greater inclusion and independence.