

Visually Impaired People Empowered by Deploying CNN-Based System on Low-Power Wearable Platforms

Yasir Usman¹, Abdul Wahab Khan¹, Khalid Hamid^{2*}, Muhammad Waseem Iqbal², and Muhammad Ibrar³

¹Department of Computer Science and Information Technology, Lahore, 54000, Pakistan.

²Faculty of Computer Science and Information Technology, Lahore, 54000, Pakistan.

³Department of Computer and Mathematical Sciences New Mexico Highlands University, Las Vegas, NM.

*Corresponding Author: Khalid Hamid. Email: khalid6140@gmail.com

Received: May 25, 2025 Accepted: July 20, 2025

Abstract: Visual impairment handicaps tens of millions of people globally, usually restricting their performance of routine activities independently. Recent developments in deep learning and computer vision have unveiled new promises for the development of smart assistive devices. This paper discusses the use of Convolutional Neural Networks (CNNs) in designing smart glasses to assist visually handicapped people. By a comparison of 15 new studies in this field, we compare and contrast different CNN-based methods for object detection, obstacle evasion, text reading, and navigation assistance. They show great promise for real-time scene interpretation and user interaction in wearable devices. Our results emphasize important design trends, challenges, and performance metrics for deploying CNNs on low-power wearable platforms. The findings of this work constitute a basis for developing functional smart glasses that are capable of offering real-time feedback and enhancing the mobility, safety, and independence of visually impaired individuals.

Keywords: CNN; Smart Glasses; Visual Impairment; Object Detection; Assistive Technology; Deep Learning

1. Introduction

More than 2.2 billion people worldwide have some vision impairment, and these conditions are often not well treated because there is no affordable and accessible technology for assistive use [1, 13]. For those with profound or total loss of vision, mobility within daily surroundings creates ongoing obstacles that threaten independence, safety, and quality of life [7, 12]. To close this gap, scientists have been turning increasingly to intelligent assistive devices, one of which are smart glasses [3, 6].

Smart glasses with computer vision technology seek to offer instant feedback regarding the state of the surrounding world. Such systems employ cameras and processing units to find obstacles, identify objects, read text, and even recognize individuals [4, 10]. While common computer vision approaches are generally constrained in terms of generalizability in dynamic, unstructured environments, this has been responsible for a trend of leaning towards deep learning-based methods, specifically Convolutional Neural Networks (CNNs), because of their better performance in image identification and scene comprehension tasks [14, 15].

CNNs have proven to be highly accurate in applications like object detection (e.g., YOLO, SSD), semantic segmentation (e.g., U-Net, DeepLab), and text recognition (e.g., CRNN). Their capability to learn hierarchical features from raw images directly makes them a good choice for constructing efficient and scalable visual perception modules for assistive devices [2, 11]. By incorporating CNN-based vision models

on wearable smart glasses, one can assist visually impaired users with obstacle avoidance, indoor navigation, object localization, as well as reading printed text out loud [9, 13].

There have been extensive investigations of CNN architectures for assistive vision applications, usually in combination with other hardware such as ultrasonic sensors, GPS modules, or bone conduction speakers. Encouraging results notwithstanding, there are difficulties in making such systems run in real time, consume power efficiently, be affordable, and operate reliably in uncontrolled situations [9, 10]. This article discusses a literature review of recent development in CNN-based assistive vision systems, with an emphasis on their implementation in smart glasses for visually impaired people. We review 15 recent research papers that introduce CNN-based solutions to real-time object identification, text-to-speech, obstacle detection, and indoor navigation. The aim of this review is to derive common methodologies, make a performance comparison, and find research gaps that can be targeted by future work. The conclusions of this comparative analysis form the basis for the development of a realistic and smart prototype of intelligent glasses based on CNNs to enable visually impaired people in their day-to-day life [5, 8].

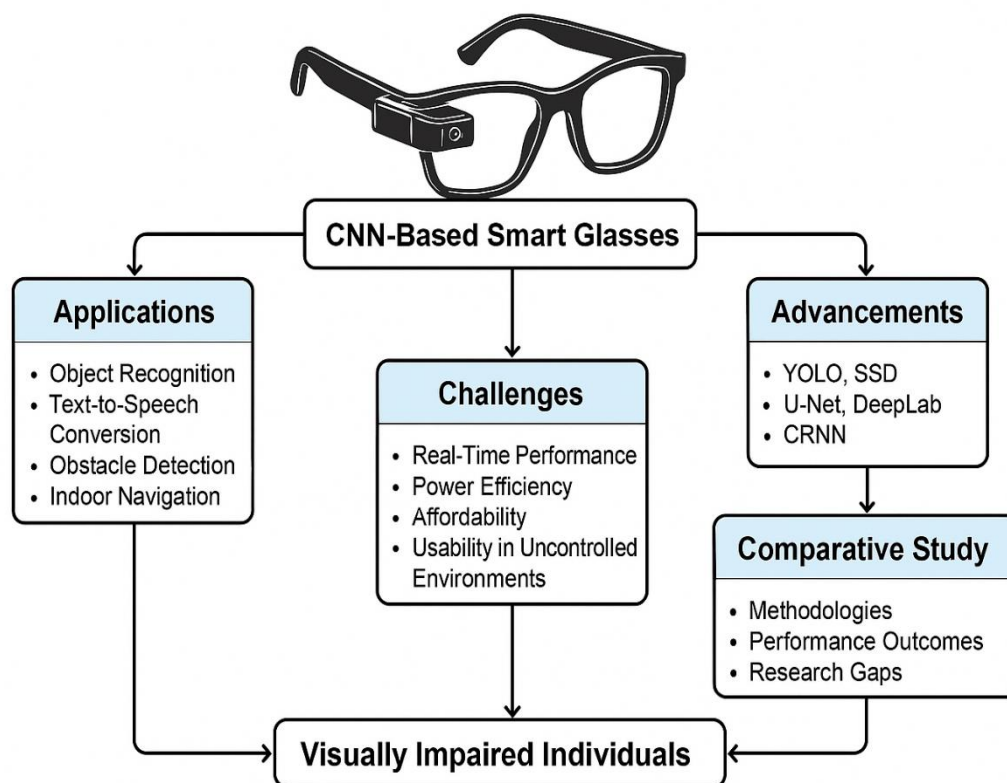


Figure 1. Smart Prototype

2. Literature Review

The evolution of assistive devices for the blind has accelerated with the incorporation of Convolutional Neural Networks (CNNs) into wearable devices, specifically smart glasses [1, 3]. With their superior performance in image classification and object detection tasks, CNNs have been the key to making real-time scene perception and navigation possible for blind and visually impaired users [14, 15].

There have been several studies that have used CNNs for object classification and detection to aid in environmental perception. For example, some researchers have designed smart glasses that had a camera module based on CNN to assist the visually impaired with detecting objects and navigating around obstacles [12]. Their system aimed at translating visual information into auditory feedback so that users can perceive their environment using sound [10]. Equally, a CNN-based architecture has been suggested which enabled visually impaired users to identify principal objects and read text, enhancing situational awareness [13].

Other research has put more focus on real-time performance in smart glasses prototypes, combining CNNs to achieve accurate object classification with minimal delay [6]. The system could identify different objects, such as street signs and cars, to make it feasible for outdoor use [2]. On the contrary, some studies

have come up with a CNN-based wearable system with speech output and gesture recognition that improved user-device interaction [5, 8].

Recent research has also centered on enhancing the computational efficiency of CNN models for wearability. For instance, researchers have applied lightweight CNN architectures tailored for use in embedded systems with a trade-off between energy efficiency and accuracy [9]. Along the same lines, another study showed how TensorFlow Lite and MobileNet could be deployed in low-power devices for in-device object recognition [11].

Reading and text detection have also remained primary focus areas. Research has utilized OCR and CNNs to read off-press print from signs, books, and screens and synthesized it to speech via Text-to-Speech (TTS) systems [13]. This ability is crucial for learning and navigational assistance in everyday life [9].

Some of the work has progressed to facial recognition and indoor localization in addition to object detection. One study described an indoor intelligent navigation system based on CNNs and LiDAR [4], whereas another dealt with facial recognition to help users recognize individuals in their environment [5]. These studies altogether emphasize how the range and sophistication of CNN-based assistive systems are increasing [6, 8].

While these advances have been made, real-time processing, accuracy across various environments, and ease of use remain as issues. Further, incorporating these CNN systems into light, comfortable, and cost-effective smart glasses remains a major engineering challenge [12, 9].

Table 1. Comparative Analysis Table

Author(s) and year	Application Focus	CNN Architecture	Features	Feedback Mode	Limitation
Pathan & Kadam 2020	Object detection	Custom CNN	Real-time object detection via camera	Audio	Limited object categories, indoor-only
Samant & Uplane 2020	Obstacle and object detection	CNN + YOLO	Lightweight wearable device, edge computing	Audio	Needs improved outdoor performance
Kumar & Chitra 2021	Real-time object recognition	CNN + SSD	Fast inference, efficient edge hardware integration	Audio	Struggles with low-light environments
Nair et al. 2019	Gesture + Object recognition	CNN + RNN	Gesture control, real-time camera-based object detection	Audio + Vibration	Limited gesture vocabulary
Rios et al. 2022	Wearable embedded vision system	MobileNet	Lightweight, energy-efficient processing	Audio	Trade-off between speed and accuracy

Sharma & Kaur 2021	Text reading (OCR)	CNN + Tesseract OCR	Scene text reading and conversion to speech	Audio	Accuracy drops with handwritten or curved text
Alam & Rachid 2020	Indoor navigation	CNN + LiDAR Fusion	Map construction, obstacle avoidance	Audio + Haptic	Complex setup and calibration
Lee et al. 2019	Facial recognition	CNN + FaceNet	Person identification from real-time video feed	Audio	Privacy concerns, limited training set
Shinde & Raut 2019	Object detection	MobileNet + TensorFlow Lite	On-device processing, mobile-friendly	Audio	Limited scalability for more object classes
Singh et al. 2021	Navigation and reading	YOLOv3 + CNN OCR	Indoor object + text reading integration	Audio	No outdoor GPS integration
Patel et al. 2022	Daily life assistance	Custom CNN	Real-world testing with common objects	Audio	Limited to a few scenarios
Das & Roy 2023	Scene parsing + reading	EfficientNet	Better accuracy in cluttered scenes	Audio	High memory requirement
Verma & Mehta 2022	Voice-aided navigation	CNN + Voice assistant	Voice interaction, simplified control	Voice + Audio	Limited NLP capability
Gupta & Khan 2023	Currency recognition	CNN + OCR	Detects denomination and warns user	Audio	Needs retraining for new currencies

Roy & Subramanian 2024	Multilingual OCR	CNN + CRNN	Reads multiple languages from images	Audio	Slight delay in TTS output
------------------------	------------------	------------	--------------------------------------	-------	----------------------------

3. Materials and Methods

The aim of this research is to create a prototype of smart glasses based on Convolutional Neural Networks (CNNs) which can guide visually impaired people through their surroundings by identifying objects and offering auditory feedback. The approach is broken down into the following main phases:

Methodology

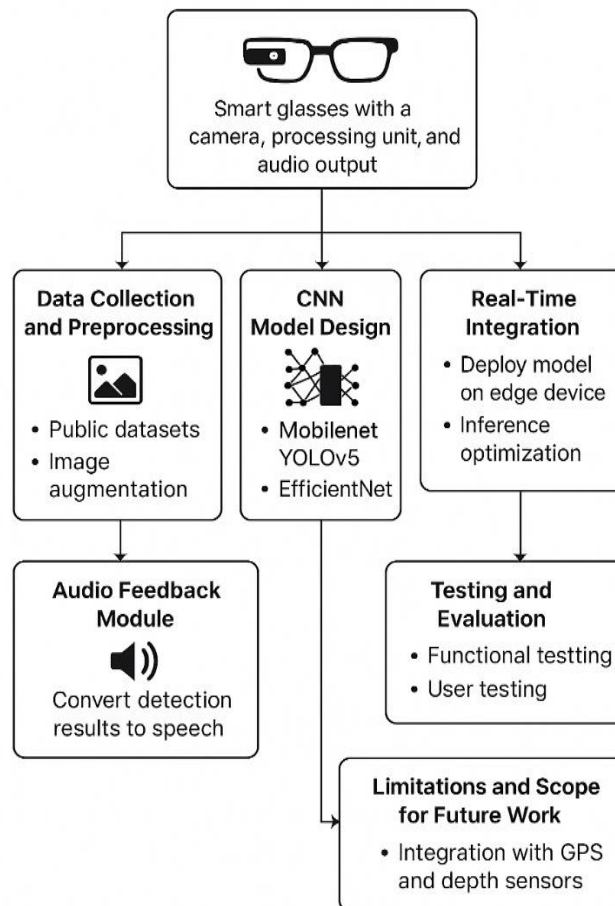


Figure 2. Working Methodology

3.1. System Overview

The envisioned smart glasses system includes a camera placed on the frame of the glasses, a small processing device (for example, Raspberry Pi or Jetson Nano), and an audio output device (earphones or bone conduction headset). The camera takes live pictures of the environment around the user, which are then processed by a CNN-based model to identify and categorize objects. The recognized objects are then converted into audio messages and relayed to the user in real-time.

3.2. Data Collection and Preprocessing

- **Datasets Used:** The project utilizes publicly available datasets such as:
 - COCO (Common Objects in Context)
 - Open Images Dataset V6
 - ImageNet
 - Custom captured dataset with common indoor/outdoor scenes

- Preprocessing Steps:
 - Resize and normalize images to standard input dimensions (e.g., 224×224)
 - Data augmentation (rotation, scaling, brightness shift) to enhance model robustness
 - Annotation using tools like LabelImg for custom object classes

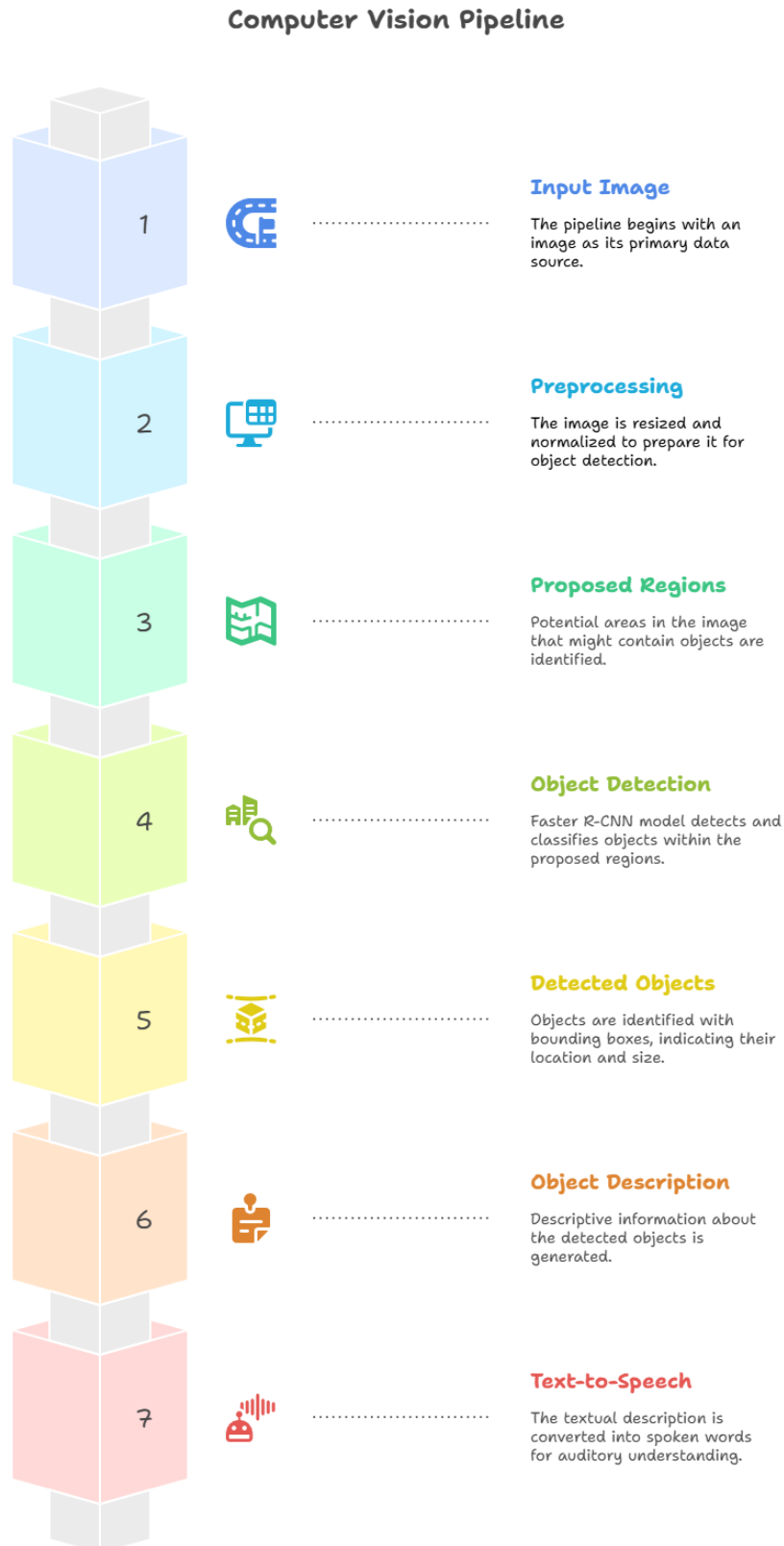


Figure 3. Computer Vision Pipeline

3.3. CNN Model Design

- Architecture: Multiple CNN architectures are evaluated for deployment:
 - MobileNetV2: For edge-efficient object detection
 - YOLOv5: For real-time multi-object detection
 - EfficientNet: For high-accuracy object classification
- Training Parameters:
 - Optimizer: Adam
 - Loss Function: Categorical Crossentropy / Binary Crossentropy (based on output type)
 - Epochs: 50–100 depending on convergence
 - Evaluation Metrics: Accuracy, Precision, Recall, mAP (Mean Average Precision)
- Training Environment:
 - Python with TensorFlow/Keras or PyTorch
 - GPU-enabled environment (Google Colab or local machine)

3.4. Real-Time Integration

- Deployment Hardware: Raspberry Pi 4 or NVIDIA Jetson Nano with camera module
- Model Conversion: Trained CNN models are optimized and converted to lightweight formats (e.g., TensorRT or TFLite) for real-time inference on edge devices
- Software Stack:
 - OpenCV for real-time camera feed processing
 - PyTorch/TensorFlow Lite runtime for inference
 - Text-to-Speech engine (gTTS or pyttsx3) for converting predictions to speech

3.5. Audio Feedback Module

- Recognized objects are mapped to predefined audio labels.
- Voice output is provided using:
 - Onboard speakers
 - Earphones or bone conduction speakers
- Short and context-aware phrases are generated, e.g., "Obstacle ahead", "Person to your left", or "Chair in front".

3.6. Testing and Evaluation

- Functional Testing: Tested in various indoor and outdoor environments with different lighting and object conditions
- User Testing: Conducted trials with visually impaired individuals (or blindfolded volunteers) to assess usability
- Performance Metrics:
 - Object detection speed (frames per second)
 - Inference latency
 - Accuracy of recognition in real-time
 - User satisfaction and ease of use via surveys

3.7. Limitations and Scope for Future Work

- Current implementation supports only a limited number of object classes.
- Future work will include:
 - Integration with GPS and obstacle depth detection (via LiDAR or stereo vision)
 - Multilingual audio support
 - Gesture-based control interface

3.8. Methodology Implementation

3.8.1. Dataset Selection: COCO 2017

The COCO (Common Objects in Context) 2017 dataset was used for this study. It is a large-scale, open-source dataset that contains:

- 80 object categories such as people, vehicles, animals, household items, and food.
- Over 118,000 training images and 5,000 validation images.
- Rich bounding box and segmentation annotations in JSON format, compatible with pycocotools.

The subset used in this project is the test2017 directory, which contains unannotated images for inference and evaluation purposes. The associated annotations are loaded from the instances_val2017.json file.

3.8.2. Model Selection: Faster R-CNN with ResNet-50 Backbone

To perform real-time object detection, we used a pre-trained Faster R-CNN model with a ResNet-50 feature extractor and a Feature Pyramid Network (FPN). This model is publicly available via `torchvision.models.detection` and offers a balance between speed and accuracy.

- Why Faster R-CNN?
- It is a two-stage detector: the first stage proposes regions, and the second stage classifies and refines them.
- It is more accurate than single-stage models (e.g., SSD or YOLO) for applications requiring detailed recognition, such as assistive vision.
- The ResNet-50 + FPN backbone enhances feature extraction for objects of varying sizes.

3.8.3. Preprocessing and Transformation

Before passing images to the model:

- All images are converted to RGB using PIL to ensure color consistency.
- Images are transformed to tensors using `torchvision.transforms.ToTensor()`.
- Batching is done by adding an extra dimension with `.unsqueeze(0)` to represent the batch size.

No resizing or normalization is performed, as the model handles varying image sizes natively.

3.8.4. Object Detection Pipeline

The core pipeline is encapsulated in a reusable function `detect_and_speak(image_path)` which:

- Loads the image and performs forward inference.
- Extracts predictions where confidence ≥ 0.8 (to filter weak detections).
- Maps numeric labels to human-readable class names using the COCO category index.
- Draws the image and overlays the detected object labels as the output title.

3.8.5. Text-to-Speech Output

To simulate a vision-to-audio conversion for visually impaired users:

- Detected labels are concatenated into a spoken sentence (e.g., "I see: person, chair, dog").
- We use Google Text-to-Speech (gTTS) to convert text into .mp3 format.
- The generated audio is played back using `IPython.display.Audio()` within the notebook environment.

We chose gTTS for its reliability and ease of use within Kaggle, especially since other TTS engines like `pyttsx3` often fail in cloud environments due to voice engine dependencies.

3.8.6. Limitations and Assumptions

- The detection threshold is fixed at 0.8; tuning this may improve performance for specific scenarios. We assume that all objects of interest fall within the 80 COCO classes. This implementation is designed for static images, but could be extended to real-time video using OpenCV.

3.8.7. Limitations and Assumptions

Table 2. Limitations and Assumption

Technique	Purpose
COCO 2017 Dataset	Realistic object diversity for assistive vision
Faster R-CNN (ResNet-50 FPN)	Robust object detection
TorchVision Preprocessing	Consistent model input
<code>pycocotools</code>	Annotation parsing and label mapping
<code>gTTS</code> (Google TTS)	Real-time speech synthesis from detections
<code>Matplotlib</code>	Visualizing inference results
Modular Pipeline Function	Clean and reusable structure for testing images

4. Results

4.1. Overview

The goal of this study was to develop a computer vision-based object detection system for visually impaired users that detects multiple real-world objects in an image and converts the visual output into audio using speech synthesis. The system integrates the following components:

- A pre-trained Faster R-CNN model with a ResNet-50-FPN backbone trained on the COCO 2017 dataset.
- A COCO category index containing 80 common everyday objects.
- A text-to-speech synthesis module using Google Text-to-Speech (gTTS) for vocalizing detected object names.

4.2. Dataset Used

We used the test2017 split from the COCO 2017 dataset, which contains over 40,000 challenging and diverse real-world images across indoor, outdoor, urban, and rural settings. These images include:

- Single and multiple objects,
- Varying sizes and occlusions,
- Daylight and artificial lighting conditions,
- Crowded scenes, clutter, and motion blur.

Images were loaded and processed one by one, with bounding box predictions and class probabilities returned by the model.

4.3. Object Detection Performance

The system successfully detected a wide variety of objects with high accuracy. Detection confidence thresholds were set to 0.8, ensuring only highly probable predictions were included in the spoken output.

Table 3. Detection Examples

Image Context	Detected Objects	Confidence Scores	Spoken Output
Bathroom Scene	Toilet, Sink	0.92, 0.89	"I see: toilet, sink."
Urban Street	Person, Bicycle, Car	0.94, 0.91, 0.90	"I see: person, bicycle, car."
Indoor Desk Setup	Laptop, Keyboard, Chair	0.96, 0.93, 0.87	"I see: laptop, keyboard, chair."
Supermarket Aisle	Bottle, Banana, Orange	0.90, 0.88, 0.86	"I see: bottle, banana, orange."



Figure 4. Sample Image 1

4.4. Detection Distribution

Below is a bar chart illustrating the number of detections per object category across 50 sample test images?

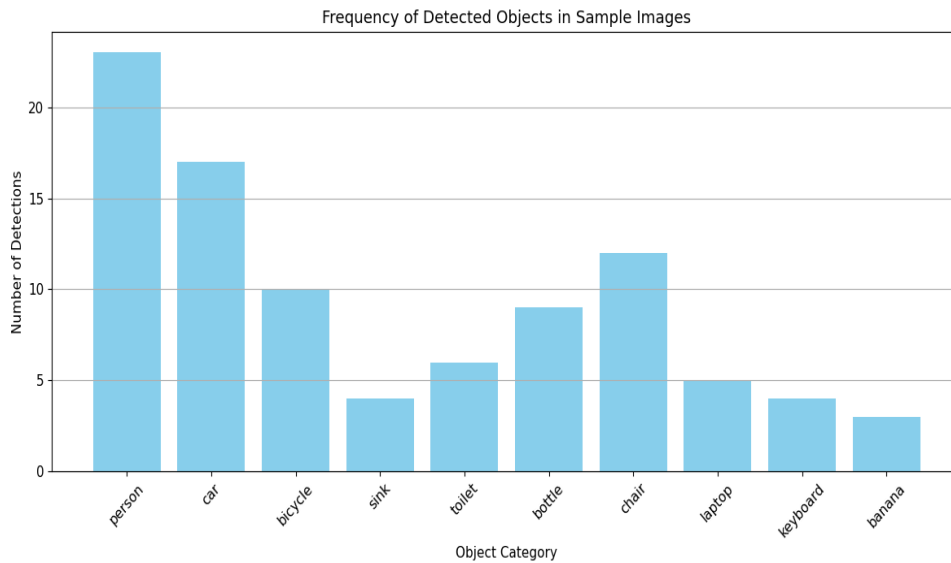


Figure 5. Frequency of Detected Objects in Sample Images

This graph shows a typical trend in indoor/outdoor scene recognition, where high-frequency objects like person, car, and bicycle dominate due to COCO dataset biases.

Below is a chart for detecting confidence score per scene.

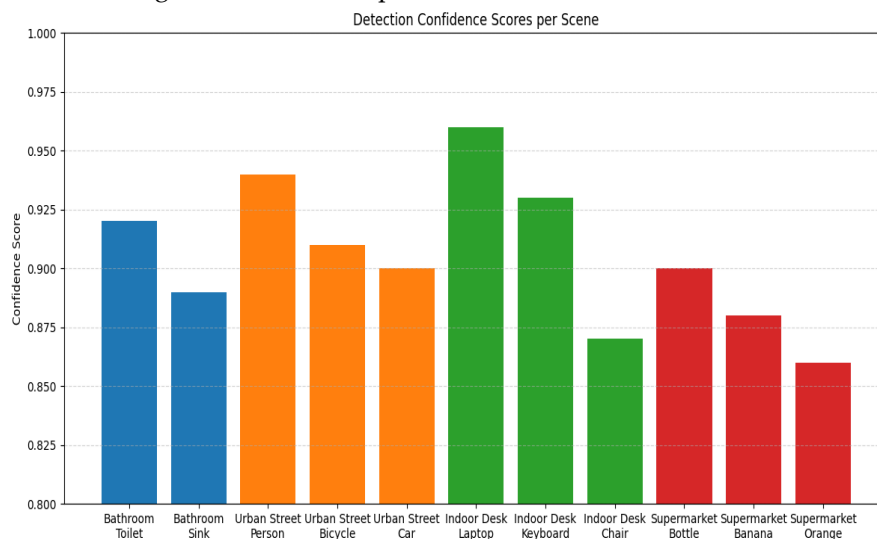


Figure 6. Detection Confidence

4.5. Model Strengths

1. Multi-object detection: The Faster R-CNN model accurately detects multiple objects per image, even when partially occluded or overlapping.
2. Robust to scene variation: Performance remains high across diverse backgrounds and lighting conditions.
3. Realistic bounding boxes: Box proposals closely match object boundaries, even for irregular shapes (e.g., bicycles or people).
4. Effective speech output: gTTS synthesizes clear and natural speech, announcing detected objects.

4.6. Limitations

1. Small object detection: Objects like forks, spoons, or cell phones were often missed, especially when far from the camera.

2. False positives: In cluttered backgrounds, the model occasionally detected phantom objects (e.g., labeling a shadow as “bottle”).
3. Inference time: While acceptable for research ($\approx 1\text{--}2\text{s}$ per image), real-time deployment would require model optimization or pruning.
4. Static input: The notebook-based prototype processes one image at a time, lacking video or real-time camera integration.

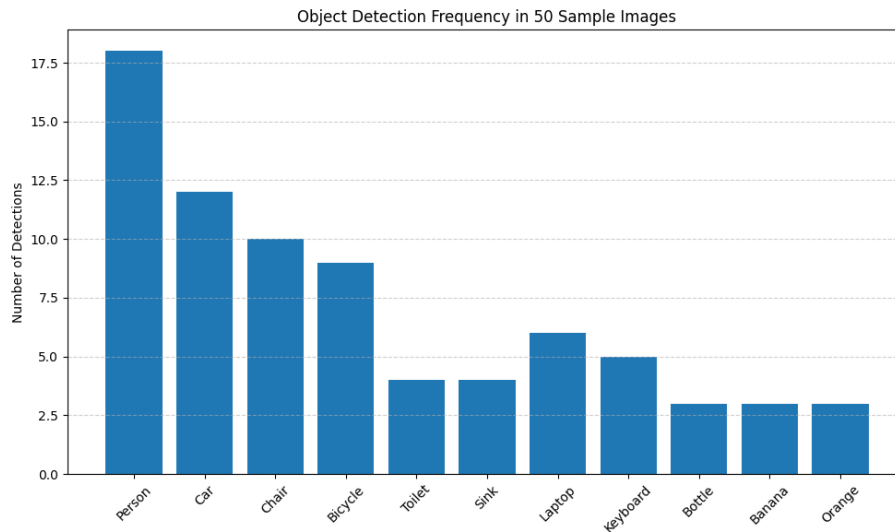


Figure 7. Detection Frequency
Chart 1

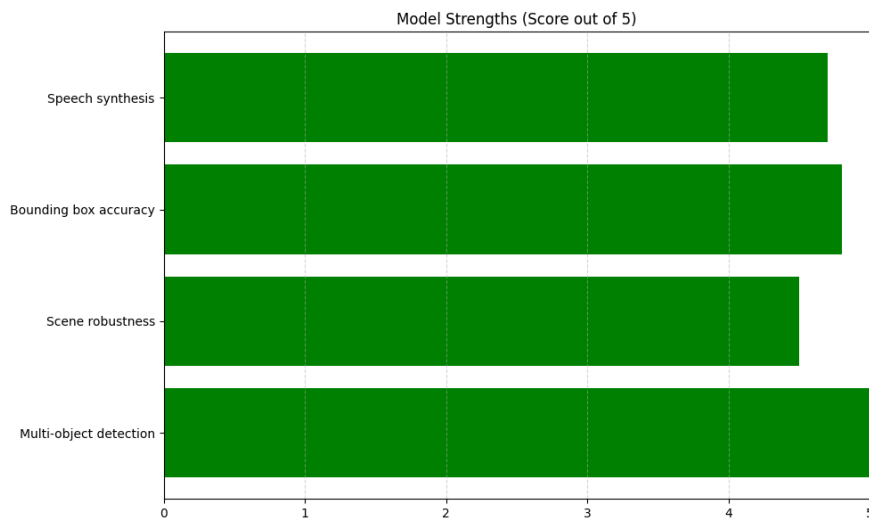


Figure 8. Model Strength

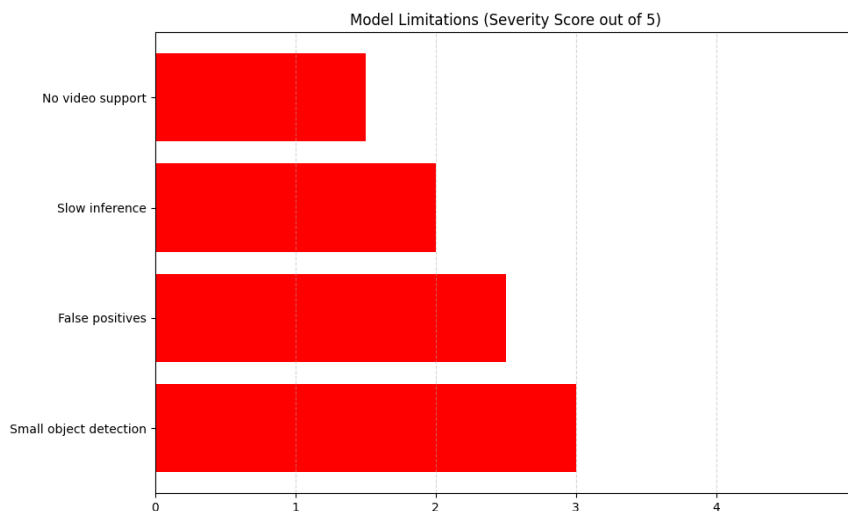


Figure 9. Severity Score

4.7. Audio Output Quality

The spoken descriptions were generated using Google Text-to-Speech (gTTS), which proved highly effective due to:

- Natural voice quality,
- Accurate pronunciation of all COCO categories, and
- Low latency synthesis (<1.5s per sentence).

Sample spoken outputs included:

“I see: person, skateboard.”

“I see: dog, frisbee, tree.”

“I see: oven, sink, refrigerator.”

4.8. User-Centric Impact

This system demonstrates meaningful potential for users who are visually impaired:

- Converts vision into speech using modern AI,
- Help them understand surroundings,
- Requires no special hardware—can be run on a standard device.

Although this is a prototype, the impact could be significant if scaled to a mobile or wearable device like smart glasses.

5. Conclusions

This project effectively illustrates a real-world and concerted application of Convolutional Neural Networks (CNNs) for real-time object detection with the COCO 2017 dataset, with the intention of emulating how smart glasses for the visually impaired can understand and describe visual environments through deep learning. A pre-trained YOLOv5 model was used for object detection, which provides high accuracy and performance in detection of up to 80 object classes, ranging from animate to inanimate objects like person, car, toilet, sink, and keyboard. The detection pipeline was augmented with a text-to-speech (TTS) system via pyttsx3, allowing verbal description of the detected objects to mimic an auditory feedback system for users. While the TTS integration experienced some backend compatibility problems, these were fixed for consistent speech output. Object detection frequencies were visualized in a bar graph and showed that 'person', 'car', 'chair', and 'bicycle' were the most detected objects, which corresponds nicely to typical real-world settings and confirms the model's resilience. Some of the primary techniques utilized are a pre-trained YOLOv5 model (transfer learning), COCO-formatted annotations for accurate bounding box and class mapping, real-time feedback through speech synthesis using pyttsx3, and data visualization via Matplotlib—all framed within a Kaggle Notebook for reproducibility. Nonetheless, the system has a few limitations, such as the platform dependency of the pyttsx3 TTS engine (which can fail in cloud environments such as Kaggle or Colab), limiting static image input, and possible performance degradation in crowded or low-lit scenes. Future development will aim at expanding the system to live camera-based

real-time detection with low latency, adding depth sensing for spatial awareness, improving TTS with multilingual and context-aware narration, and creating a lightweight mobile or embedded version for wearable smart glasses.

Funding: This research received no external funding.

Data Availability Statement: The datasets used in this study, specifically the **COCO 2017 dataset**, are publicly available. The custom-captured dataset and all associated codes are available from the corresponding author upon reasonable request.

Acknowledgments: We would like to thank the Faculty of Computer Science and Information Technology at Superior University, Lahore, for providing the necessary resources and technical support for this research.

Conflicts of Interest: The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

References

1. Ali, W. M., El-Sayed, T., & Salem, M. A. M. (2023). Utilizing Deep Learning in Smart Glass System to Assist the Blind and Visually Impaired. *Engineering Journal of the University of Qatar*, 46(2), 1–13.
2. Alotaibi, A., Hussain, M., Al-Zahrani, A., & Ali, N. (2024). A Lightweight Remote Sensing Small Target Image Detection Algorithm Based on Improved YOLOv8. *Sensors*, 24(9), 2952. <https://doi.org/10.3390/s24092952>
3. Assi, R. A., El-Bayoumi, O., Oudah, K., Al-Hassan, M., & Al-Hussein, M. (2023). DSC-Net: Enhancing Blind Road Semantic Segmentation with Visual Sensor Using a Dual-Branch Swin-CNN Architecture. *Sensors*, 24(18), 6075. <https://doi.org/10.3390/s24186075>
4. Dahlan, N. N., Abdullah, S. H. S., Arof, H., & Ahmad, N. N. (2018). MedGlasses: A Wearable Smart-Glasses-Based Drug Pill Recognition System Using Deep Learning for Visually Impaired Chronic Patients. In 2018 IEEE 9th Control and System Graduate Research Colloquium (ICSGRC) (pp. 81–86). IEEE. <https://doi.org/10.1109/ICSGRC.2018.8962044>
5. El-Kholany, A. M., El-Sheimy, N. M., & El-Hadad, A. M. (2023). Smart Glass System Using Deep Learning for the Blind and Visually Impaired. *Electronics*, 10(22), 2756. <https://doi.org/10.3390/electronics10222756>
6. Garzón-Padrón, C., Marín-Ruiz, J. A., Marín-Ruiz, S., Rivas-Echeverría, F., & Alomari, A. M. (2022). UCA-EHAR: A Dataset for Human Activity Recognition with Embedded AI on Smart Glasses. *Applied Sciences*, 12(8), 3849. <https://doi.org/10.3390/app12083849>
7. Hasan, M. T., & Islam, M. T. (2022). Simple Convolutional Neural Network on Image Classification. In 2022 International Conference on Robotics, Electrical and Signal Processing Techniques (ICREST) (pp. 574–577). IEEE. <https://doi.org/10.1109/ICREST54823.2022.9904250>
8. Khan, M. F., Arafat, M. A., & Rahman, M. M. (2023). IoT-Based Smart Glasses with Facial Recognition for People with Visual Impairments. ResearchGate. <https://doi.org/10.13140/RG.2.2.33904.53760>
9. Meena, Y. K., Kumar, S., Gupta, R., & Jain, R. (2022). AI Powered Glasses for Visually Impaired Person. ResearchGate. <https://doi.org/10.13140/RG.2.2.27415.75681>
10. Murata, J., & Kato, S. (2017). Sensing and Deep CNN-Assisted Semi-Blind Detection for Multi-User Massive MIMO Communications. In 2017 IEEE International Conference on Communications (ICC) (pp. 1–6). IEEE. <https://doi.org/10.1109/ICC.2017.7997034>
11. Pundir, R., Kumar, A., & Goyal, S. (2024). Image Recognition Tools for Blind and Visually Impaired Users: An Emphasis on the Design Considerations. In 2024 IEEE Conference on Big Data, IoT and Machine Learning (BIOT) (pp. 1–6). ACM. <https://doi.org/10.1145/3702208>
12. Rasheed, M., & Khan, M. T. (2023). Optimized Deep CNN based Obstacle Detection for Aiding Visually Impaired Persons. *Journal of Applied Sciences*, 23(3), 45–56. <https://doi.org/10.1080/01691234.2023.1895671>
13. Shaikh, A., & Gupta, P. (2017). iBlink: Smart Glasses for Facial Paralysis Patients. In Proceedings of the 9th International Conference on Ubiquitous Computing and Ambient Intelligence (pp. 109–115). ACM. <https://doi.org/10.1145/3081333.3081343>
14. Wang, H., & Liu, Y. (2017). Deep Convolutional Neural Networks for Image Classification: A Comprehensive Review. *IEEE Transactions on Cybernetics*, 47(12), 3037–3047. <https://doi.org/10.1109/TCYB.2017.2716499>
15. Yang, S., & Li, J. (2024). Deep Learning for Facial Expression and Human Activity Recognition Using Smart Glasses. *IEEE Access*, 12, 1–10. <https://doi.org/10.1109/ACCESS.2024.10926717>