

AssistedVision

Real-Time Navigation Aid for Visually Impaired Individuals
with Adaptive Turn Detection and Geometric Gap Analysis

Lakshminarayan Shrivs , Mahammad Afandiyev, Mengfan Zhang
Sudharsan Thiagarajan, Kumar Vaibhav

MSML640: Computer Vision

December 8, 2025

Abstract

Independent mobility is fundamental for quality of life, yet millions of visually impaired individuals face daily navigation challenges. This report presents **AssistedVision**, a comprehensive real-time navigation assistance system that addresses critical gaps in existing assistive technology through three novel contributions: (1) **Adaptive Turn Detection** achieving 90% accuracy in validating user compliance with directional instructions, reducing contradictory guidance from 67% to 10%; (2) **Geometric Gap Quantification** calculating real-world gap widths with 11.2cm mean absolute error using camera field-of-view geometry and monocular depth estimation; (3) **Dual-Method Wall Detection** combining PathFinder floor segmentation with depth-based analysis, achieving 92% recall and reducing false negatives by 60%. The system operates at 28 FPS on consumer hardware (\$300-700), demonstrating 76% collision reduction and 60% navigation speed improvement over baseline methods in user studies with 4.5/5 satisfaction rating. All six bonus tasks completed, with public dataset release and comprehensive ethical analysis.

Keywords: Assistive technology, Computer vision, Navigation systems, Depth estimation, Object detection, Behavioral tracking

GitHub Repository: <https://github.com/LNSHRIVAS/AssistedVision>

Dataset: https://drive.google.com/drive/folders/1TjCOPhU7Z5v3R_NyNBkg33-9IfXt7kSQ

Contents

1	Introduction	3
1.1	Motivation and Background	3
1.2	Problem Statement	3
1.3	Novel Contributions	3
2	Related Work	4
2.1	Assistive Navigation Systems	4
2.2	Object Detection and Depth Estimation	4
2.3	Research Positioning	4
3	System Overview	5
3.1	System Architecture	5
3.2	Application Context	6
3.3	Application Context	6
4	Method	6
4.1	Adaptive Turn Detection Algorithm	6
4.2	Geometric Gap Quantification	7
4.3	Dual-Method Wall Detection	7
5	Experiments and Results	8
5.1	Experimental Setup	8
5.2	Dataset Creation and Collection	8
5.3	Turn Detection Performance	9
5.4	Gap Width Estimation Accuracy	10
5.5	Wall Detection Evaluation	10
5.6	Real-Time Performance Analysis	10
5.7	User Evaluation Study	11
6	Bonus Tasks Summary	11
7	Limitations and Future Work	12
7.1	Current Limitations	12
7.2	Future Directions	12
8	Conclusion	13

1 Introduction

1.1 Motivation and Background

Independent mobility is a fundamental requirement for quality of life, yet millions of visually impaired individuals face daily challenges navigating complex environments. According to the World Health Organization (WHO), approximately 2.2 billion people worldwide have a vision impairment, with 36 million being blind and 217 million having moderate to severe vision impairment [1].

Traditional mobility aids such as white canes and guide dogs, while valuable, provide limited information about the surrounding environment and cannot anticipate dynamic obstacles or provide optimal path planning. Recent advances in computer vision, deep learning, and mobile computing have created unprecedented opportunities for developing intelligent navigation assistance systems.

However, most existing solutions focus solely on obstacle detection without considering the **behavioral aspects** of human navigation—such as whether users actually follow directional instructions, or whether gaps between obstacles are wide enough to pass through comfortably. Furthermore, many systems suffer from practical deployment challenges including high computational requirements, poor real-time performance, and inadequate handling of network latency in distributed architectures.

1.2 Problem Statement

Current assistive navigation systems face several critical limitations:

1. **Lack of User Compliance Verification:** Systems issue directional commands without confirming that users have completed the instructed turns, leading to contradictory instructions when network lag causes processing delays.
2. **Binary Obstacle Assessment:** Existing systems typically classify spaces as either “safe” or “unsafe” without quantifying navigable gaps between objects, forcing users to make spatial judgments they cannot visually verify.
3. **Single-Modality Detection Failures:** Reliance on a single detection method (e.g., only floor segmentation or only object detection) can result in false negatives, particularly for plain walls or uniform surfaces.
4. **Insufficient Spatial Information:** Audio feedback often lacks precise quantitative information about gap widths, distances, and spatial configurations that would enable informed navigation decisions.

1.3 Novel Contributions

This report presents AssistedVision, a comprehensive navigation assistance system with the following novel contributions:

1. **Adaptive Turn Detection System (90% Accuracy):** A behavioral tracking mechanism that monitors obstacle position history across 15 frames, validates turn completion using stability criteria, and prevents premature directional updates during active navigation maneuvers.
2. **Geometric Gap Quantification (11.2cm MAE):** A novel algorithm that calculates real-world gap widths between obstacles using camera field-of-view geometry, depth estimation, and the law of cosines for angled configurations.

3. **Dual-Method Wall Detection (92% Recall):** A redundant sensing approach combining PathFinder floor segmentation with depth-based center region analysis, reducing false negatives by 60% compared to single-method approaches.
4. **Clock-Position Audio Guidance:** An intuitive 12-position directional system optimized for audio-only navigation, with gap width announcements and direction consistency maintenance.
5. **Real-Time Performance on Consumer Hardware (28 FPS):** Optimization strategies achieving real-time processing on CPU-only systems using smartphone IP webcam streaming, at 1/17th the cost of LiDAR-based systems.

2 Related Work

2.1 Assistive Navigation Systems

The development of electronic travel aids (ETAs) for visually impaired individuals has been an active research area for over five decades. Early systems like the Sonic Pathfinder [2] (1974) and Mowat Sensor [3] (1997) used ultrasonic sensors to detect obstacles, but provided limited range and spatial information.

Vision-Based Systems: The advent of computer vision enabled more sophisticated approaches. Hub et al. [4] developed a stereo vision-based system for real-time indoor navigation.

Deep Learning Approaches: Recent works leverage deep learning for enhanced perception. Bai et al. [5] developed a CNN-based system for indoor scene understanding using AlexNet, for real-time obstacle classification.

Research Gap: While these systems detect obstacles effectively, most lack validation of user compliance with instructions or quantification of navigable gaps. The recent PathFinder system [6] addresses offline navigation with monocular depth-based pathfinding but does not validate turn completion or provide quantitative gap metrics.

2.2 Object Detection and Depth Estimation

The YOLO (You Only Look Once) family revolutionized real-time object detection [7, 8, 9]. YOLOv8 achieves state-of-the-art performance with 37.3 mAP on COCO at 280 FPS on GPU hardware.

For depth estimation, MiDaS [10, 11] introduced robust monocular depth estimation trained on 10 diverse datasets through multi-objective optimization. MiDaS v3.0 uses Vision Transformers (DPT-Large) and achieves 95.1% accuracy on KITTI with zero-shot transfer to unseen datasets.

2.3 Research Positioning

Our work differs from existing systems through:

- Behavioral validation via turn detection (absent in all prior systems)
- Quantitative gap analysis with geometric calculations
- Redundant wall detection combining multiple modalities
- Latency-aware design with adaptive feedback timing
- Real-time performance on consumer hardware at 1/17th the cost of LiDAR systems

3 System Overview

3.1 System Architecture

AssistedVision employs a distributed architecture optimized for real-time performance on consumer hardware. The system comprises four principal layers:

- **Perception Layer:**

- Object Detection: YOLOv8n (3.2M parameters, 37.3 mAP on COCO, conf=0.20)
- Depth Estimation: MiDaS v3.0 (DPT-Large, 95.1% accuracy)
- Floor Segmentation: PathFinder algorithm

- **Processing Layer:**

- Multi-Object Tracking: ByteTrack-inspired with Kalman filtering
- Turn Detection Module: 15-frame history buffer, multi-criteria validation
- Gap Quantification: Geometric calculator using camera FOV (60°) and depth
- Dual-Method Wall Detection: PathFinder + depth-based center analysis

- **Decision Layer:**

- Priority-based guidance logic (**Gaps, Walls, High-risk obstacles**)
- Clock-position encoding (12 directional positions)
- Audio cooldown management (8-second minimum)

- **Feedback Layer:**

- Text-to-Speech: PowerShell SAPI subprocess-based
- Audio Messages: Gap width + clock position format

AssistedVision System

Smartphone: Samsung Galaxy (IP Webcam, 640×480@30FPS)

Laptop: Intel i5-1135G7, 16GB RAM, CPU-only

Network: WiFi 802.11ac (50-150ms latency)

Cost: \$300-700

Perception Layer:

- **YOLOv8n:** 42ms, 37.3 mAP, Object detection
- **MiDaS v3.0:** 18ms (every 10 frames), Depth estimation
- **PathFinder:** Real-time floor segmentation

Processing Layer:

- **Turn Detection:** 15-frame history, 90% accuracy
- **Gap Quantification:** Geometric calculation, 11.2cm MAE
- **Wall Detection:** Dual-method (PathFinder + depth), 92% recall
- **Multi-Object Tracking:** ByteTrack-inspired, 1ms

Decision Layer:

- Priority: **Gaps, Walls, High-Risk Obstacles**
- Clock-position guidance (12 positions, 30° each)
- Audio cooldown: 8.0s (compensates network lag)

Feedback Layer:

- TTS: PowerShell SAPI (Volume 100, Rate 1-3)

- Message format: "Gap 75 cm at 11 o'clock"
- Performance:** 28.3 FPS — 650-750ms latency — Real-time processing

Hardware Configuration: Android smartphone (Samsung Galaxy) running IP Webcam app, laptop (Intel i5-1135G7, 16GB RAM, CPU-only), WiFi 802.11ac (50-150ms latency). Total system cost: \$300-700.

3.2 Application Context

The system serves as a secondary sensory channel to augment traditional aids (white canes, guide dogs), particularly for dynamic obstacles or overhead hazards. Tested in three environments: indoor corridors ($15m \times 2m$), office spaces ($8m \times 6m$), and outdoor walkways.

Quantitative Impact:

Hardware Configuration: Android smartphone (Samsung Galaxy) running IP Webcam app, laptop (Intel i5-1135G7, 16GB RAM, CPU-only), WiFi 802.11ac (50-150ms latency). Total system cost: \$300-700.

3.3 Application Context

The system serves as a secondary sensory channel to augment traditional aids (white canes, guide dogs), particularly for dynamic obstacles or overhead hazards. Tested in three environments: indoor corridors ($15m \times 2m$), office spaces ($8m \times 6m$), and outdoor walkways.

Quantitative Impact:

bbstacle Avoidance: 94% success rate vs. 78% baseline

- Collision Rate: 0.04 per trial vs. 0.17 for white cane baseline
- Navigation Speed: 0.8 m/s (60% faster than white cane-only at 0.5 m/s)
- User Satisfaction: 4.5/5 average rating (n=5 participants)

4 Method

4.1 Adaptive Turn Detection Algorithm

Motivation: Network latency (200-800ms) causes instructions to lag behind user movements. Without turn validation, users receive contradictory directions before completing instructed maneuvers, reducing compliance by 47%.

Algorithm Overview: The system maintains a 15-frame obstacle position history and validates turn completion using four criteria:

1. **Stability Check:** Recent positions vary by ≤ 1 clock position
2. **Movement Check:** Obstacle shifted ≥ 2 clock positions
3. **Centering Check:** Obstacle now at 11-1 o'clock (front-centered)
4. **Risk Reduction:** Obstacle risk decreased below 40%

Turn confirmation requires: (Stability AND Movement) OR (Centering AND Risk Reduction)

Parameters:

- 15-frame history buffer (0.5-second window)

- 2 clock-position threshold (30° angular change)
- 5-second timeout prevents indefinite waiting

Performance: Achieved 90% accuracy (45/50 trials), reduced contradictory instructions from 67% to 10%.

4.2 Geometric Gap Quantification

Problem: Calculate real-world width of passable gaps between obstacles for traversability assessment.

Case 1: Perpendicular Configuration (depth difference $<20\%$):

When objects are at similar depths:

$$\text{gap} = 2 \cdot d_{avg} \cdot \tan\left(\frac{\theta_{separation}}{2}\right) \quad (1)$$

where d_{avg} is average depth and $\theta_{separation} = \frac{\text{pixel_separation_FOV}}{\text{frame_width}}$.

Case 2: Angled Configuration (depth difference $\geq 20\%$):

Using law of cosines:

$$\text{gap}_{raw} = \sqrt{d_1^2 + d_2^2 - 2d_1d_2 \cos(\theta_{separation})} \quad (2)$$

with projection correction: $\text{gap} = \text{gap}_{raw} \cdot \frac{\min(d_1, d_2)}{\max(d_1, d_2)}$

Confidence Scoring:

$$\text{confidence} = 0.4 \cdot \frac{w}{2.0} + 0.4 \cdot \left(1 - \frac{d_{avg}}{4.0}\right) + 0.2 \cdot \text{perp_bonus} \quad (3)$$

Performance: 11.2cm MAE overall, 77% within 10cm of ground truth.

4.3 Dual-Method Wall Detection

Challenge: PathFinder floor segmentation alone achieved only 40% recall, missing plain painted walls and uniform surfaces.

Method 1: PathFinder Floor Segmentation

- Color/texture-based floor region detection
- Identifies where floor ends (walls, drop-offs)
- Strengths: Effective on textured floors (tiles, carpets)
- Weaknesses: Fails on uniform colors (white on white)

Method 2: Depth-Based Center Region Analysis

- Sample center region (1/3 of frame)
- Compare median disparity to 95th percentile
- If center disparity $> 70\%$ of max disparity \rightarrow wall detected
- 60-frame (2 sec) threshold prevents false positives

Combined Decision Logic:

$$\text{wall_detected} = \text{PathFinder OR (no_objects AND high_center_disparity)} \quad (4)$$

Performance: 92% recall vs. 40% (PathFinder-only) or 75% (depth-only), F1=90%.

5 Experiments and Results

5.1 Experimental Setup

Hardware: Samsung Galaxy Android 11 with IP Webcam app, Intel i5-1135G7 laptop (16GB RAM, CPU-only), WiFi 802.11ac.

Software: Python 3.13.5, PyTorch 2.0.1, OpenCV 4.8.0, YOLOv8n (ultralytics), MiDaS v3.0.

Test Environments:

1. Indoor Corridor: 15m×2m, fluorescent lighting, plain walls
2. Office Space: 8m×6m, complex layout, varied lighting
3. Outdoor Walkway: Natural lighting, varying surfaces, dynamic obstacles

Evaluation Data:

- 50 navigation trials for turn detection (10 per participant, 5 participants)
- 30 gap measurements with ground truth (tape measure, 1cm precision)
- 40 wall approach scenarios (20 textured, 20 plain surfaces)
- 5 user evaluation sessions (10 trials each)

5.2 Dataset Creation and Collection

The quality of a machine learning model depends heavily on the dataset used for training and evaluation. For this project, we independently designed, captured, organized, and validated an original dataset to support robust model development for assisted vision and object risk assessment tasks.

Recording Setup: All dataset videos were recorded using an Android smartphone (Samsung Galaxy) to ensure accessibility, portability, and realistic conditions. The camera was held at a consistent height of approximately 4.5 feet to simulate natural human perspective during navigation tasks. Each video was captured with an average duration of 10 seconds, providing sufficient temporal information for downstream processing while keeping file sizes manageable.

Environmental Diversity: To build a dataset that reflects real-world variability, we recorded videos across multiple locations and lighting conditions:

- **Morning:** Strong daylight with natural brightness
- **Evening/Noon:** Soft ambient light with varying brightness
- **Night:** Low-light recordings for challenging visibility conditions

Capturing under diverse environmental conditions increases model generalizability and robustness to domain shift.

Sensor-Aware Video Collection: Beyond standard video recording, we created a separate set of sensor-aware videos using the web tool av-record-video.netlify.app. This tool simultaneously captured:

- Accelerometer data
- Gyroscope readings
- Device orientation
- Recording timestamps

These videos provide deeper insight into camera behavior and motion dynamics, valuable for understanding input variability and calibrating model behavior.

Dataset Organization: All videos were systematically stored in Google Drive with six major folders reflecting different recording contexts:

1. **Morning** – Natural daylight recordings
2. **Evening/Noon** – Afternoon/evening with varying brightness
3. **Night** – Low-light testing scenarios

4. **Newaddition** – Expanded samples addressing coverage gaps
5. **AV-RecordedVideos** – Videos with additional sensor metadata
6. **Test** – Model demo outputs with processed videos and `outputlog.jsonl`

Code Improvements and Validation: During initial attempts, the project code required modifications. We created a dedicated branch (`AssistedvisionModified`) where we:

- Fixed critical execution issues
- Updated video paths and Python version dependencies in `run.sh`
- Improved polygon selection logic in `main.py` for accurate object risk assessment
- Ensured compatibility between dataset structure and processing pipeline

After modifications, we conducted multiple validation runs on the dataset. Each video was processed using the updated pipeline, with results (bounding polygons, risk assessments, visual overlays) saved as output demo videos. Both the dataset and demo outputs were uploaded to Google Drive for full reproducibility.

Final Dataset Characteristics:

- Recording device: Android smartphone
- Average video length: 10 seconds
- Camera height: 4.5 feet
- Lighting variations: Morning, Noon/Evening, Night
- Environmental diversity: Multiple indoor and outdoor locations
- Sensor-augmented videos: Available in AV-RecordedVideos folder
- Total organizational folders: 6
- Processing: All videos processed with logs and demo outputs

This comprehensive dataset now serves as a strong foundation for the project's machine learning, computer vision, and risk assessment tasks, with systematic organization enabling easy access, sharing, and version control through Google Drive.

5.3 Turn Detection Performance

Table 1: Turn Detection Results (50 Trials)

Metric	Value	Notes
True Positives	45/50	90% correctly validated
False Positives	3/50	6% premature detection
False Negatives	2/50	4% missed turns
Overall Accuracy	90.0%	45 correct / 50 total
Avg Detection Time	0.83±0.21s	Turn start to confirmation
Timeout Occurrences	1/50	2% exceeded 5-second limit
Contradictory Instructions:		
Without turn detection	34/50 (67%)	Baseline
With turn detection	5/50 (10%)	57 pp reduction

5.4 Gap Width Estimation Accuracy

Table 2: Gap Width Estimation Results by Configuration

Configuration	N	MAE (cm)	Max Error	<10cm	<20cm
Perpendicular ($\Delta d < 20\%$)	18	8.3	18cm	87%	100%
Angled ($\Delta d \geq 20\%$)	12	14.7	28cm	67%	92%
Overall	30	11.2	28cm	77%	97%

Depth-Stratified Analysis:

- 0.5-2.0m: MAE = 7.1cm (most accurate)
- 2.0-3.0m: MAE = 11.8cm (moderate)
- 3.0-4.0m: MAE = 16.4cm (higher uncertainty)

5.5 Wall Detection Evaluation

Table 3: Wall Detection Performance Comparison (40 Scenarios)

Method	Recall	Precision	F1-Score	Detected/Total
PathFinder only	40%	94%	56%	16/40 (1 FP)
Depth-based only	75%	81%	78%	30/40 (7 FP)
Dual method (Combined)	92%	88%	90%	37/40 (5 FP)

Surface Type Breakdown:

- Textured walls: Dual method 95% vs. PathFinder 85% vs. Depth 70%
- Plain walls: Dual method 90% vs. PathFinder 10% vs. Depth 80%
- False negative reduction: 60% (8% miss rate vs. 60% single-method)

5.6 Real-Time Performance Analysis

Table 4: System Latency Breakdown and FPS Performance

Component	Latency (ms)	% of Total
Camera capture	33	—
Network transmission (WiFi)	50-150	—
YOLOv8n inference	42	65%
MiDaS depth (amortized)	18	28%
Gap calculation	3	5%
Turn detection	1	1%
Risk scoring	2	3%
Total processing/frame	65ms	100%
TTS audio synthesis	500	—
End-to-end latency	650-750ms	—
Average FPS	28.3	Min: 24, Max: 32

Optimization Strategies:

- Frame skipping: Process every 2nd frame ($2\times$ speedup)
- Asynchronous depth: MiDaS every 10 frames ($10\times$ reduction)
- CPU-only: No GPU required, eliminates CUDA overhead
- Model selection: YOLOv8n (3.2M params) vs. YOLOv8x (68.2M) trades 8% mAP for $6\times$ speed

5.7 User Evaluation Study

Participants: N=5 (3 male, 2 female, ages 22-28), all sighted, blindfolded to simulate visual impairment. 10 trials each (50 total trials).

Table 5: Subjective User Ratings (5-point Likert scale)

Question	Mean \pm SD	Distribution
Clarity of audio instructions	4.6 ± 0.5	[4,4,5,5,5]
Confidence in navigation	4.2 ± 0.7	[3,4,5,5,4]
Comfort with gap guidance	4.4 ± 0.6	[4,4,5,5,4]
Usefulness of turn validation	4.8 ± 0.4	[4,5,5,5,5]
Overall satisfaction	4.5 ± 0.5	[4,4,5,5,5]

Qualitative Feedback Highlights:

- “Gap width information helps me decide confidently whether to go through or around.”
- “System doesn’t contradict itself anymore [compared to baseline].”
- “Clock positions are intuitive, I can immediately understand which direction.”
- “Knowing exact gap size (70cm vs. 90cm) makes a huge difference.”

Table 6: Objective Performance Metrics Comparison

Metric	AssistedVision	Audio-Only	White Cane
Collision rate (per trial)	0.04	0.17	0.22
Navigation speed (m/s)	0.8 ± 0.15	0.6	0.5 ± 0.12
Gap traversal success	94% (44/47)	N/A	65%
Time for 15m corridor (s)	22.3 ± 4.1	28.5	30.0
User satisfaction (1-5)	4.5 ± 0.5	3.8 ± 0.6	3.2 ± 0.7
Collision reduction	76% vs. audio-only, 82% vs. white cane		

6 Bonus Tasks Summary

The following bonus feature was implemented to extend the system’s capabilities:

1. **Android Mobile Mode: (2 pts):** Implemented a complete mobile deployment system enabling the use of Android smartphones as the primary camera and audio output device. The system architecture includes: (1) IP Webcam integration for video streaming from the phone camera to the laptop running the CV pipeline, (2) WebSocket-based audio streaming server (mobile_server.py) for delivering real-time spoken guidance directly to the phone, (3) Gyroscope data transmission from phone to laptop for potential future orientation-aware features,

- (4) Browser-based mobile companion interface (mobile.html) with motion/audio permissions. This bonus feature extends the system beyond desktop-only deployment, making it accessible for real-world mobile use cases where users can leverage their existing smartphone hardware. The implementation is fully documented in the README with setup instructions and is tested on Android devices.
2. **Data-in-the-Wild (1 pt):** Tested system in 3 diverse uncontrolled environments (indoor, office, outdoor). Performance degradation analysis: FPS 28-30 (indoor) to 24-28 (outdoor), turn accuracy 92% to 85%, gap MAE 9.8cm to 14.1cm. Domain shift observations documented.
 3. **User Study (2 pts):** Conducted with 5 participants, 50 trials total. Demonstrated 76% collision reduction, 60% speed improvement, 4.5/5 satisfaction. Qualitative feedback revealed usability insights and design recommendations.
 4. **Ethical Considerations (1 pt):** Analyzed 6 ethical issues including data privacy, user autonomy, accessibility equity, and bias in training data. Proposed mitigation strategies: local processing, adjustable risk tolerance, subsidized devices, diverse dataset fine-tuning.
 5. **Data Collection (2 pts):** Collected original dataset: 50 navigation trials (45 min video, 81,000 frames), 30 gap measurements, 40 wall scenarios. Applied augmentation (brightness, noise, blur). Public release on Google Drive with CC-BY-4.0 license.
 6. **Cross-Modal Integration (2 pts):** Integrated depth modality (MiDaS) with vision (YOLO + PathFinder). Demonstrated 130 pp improvement in wall detection recall (40% → 92%), enabled quantitative gap analysis (11.2cm MAE), improved risk scoring by 18%.
 7. **Mobile Optimization (2 pts):** Achieved 28 FPS on CPU-only consumer hardware through: YOLOv8n model selection (6× speedup), frame skipping (2×), asynchronous depth (10× reduction), IP webcam streaming. Cost: \$300-700 vs. \$5000 LiDAR systems.

Total Bonus Points: 10 completed (5 maximum applied per syllabus)

7 Limitations and Future Work

7.1 Current Limitations

- **Glass Detection:** Neither PathFinder nor depth-based method detects transparent surfaces (3/40 failures, 7.5% FN rate)
- **Monocular Depth Ambiguity:** Scale uncertainty causes 11.2cm MAE; stereo cameras could reduce to 2cm
- **Network Dependency:** Requires WiFi/cellular; offline on-device deployment needed
- **Turn Detection Robustness:** 6% false positive rate when users hesitate mid-turn
- **Battery Life:** Continuous streaming drains smartphone in 3 hours
- **Small Sample Size:** User study n=5 should be 20-30 for statistical significance

7.2 Future Directions

- **Short-term (6-12 months):** Gyroscope integration (90%→95% turn accuracy), vibration feedback for critical warnings, FOV calibration (11.2cm→8-9cm MAE)
- **Medium-term (1-2 years):** On-device smartphone deployment (TensorFlow Lite, 15 FPS target), longitudinal studies with visually impaired participants (n=20-30, 4-6 weeks), outdoor navigation extension (GPS, curb detection, crosswalk recognition)

- **Long-term (3-5 years):** Reinforcement learning-based personalized path planning and socially aware navigation in crowded environments, leveraging multi-modal sensor fusion (ultrasonic for detecting glass, thermal imaging for low-light conditions, and LiDAR for high-precision mapping), with the goal of commercializing this solution for large enterprises such as Google or Meta.

8 Conclusion

This report presented **AssistedVision**, a real-time intelligent navigation system for visually impaired individuals that addresses critical gaps in existing assistive technology through three primary innovations:

1. **Adaptive Turn Detection (90% accuracy):** First system to validate user compliance with directional instructions, reducing contradictory guidance from 67% to 10% through 15-frame obstacle position tracking and multi-criteria validation. Compensates for network latency inherent in distributed architectures.
2. **Geometric Gap Quantification (11.2cm MAE):** Novel algorithm calculating real-world gap widths using camera FOV geometry and monocular depth estimation. Distinguishes perpendicular and angled configurations, achieving 77% accuracy within 10cm. Provides quantitative spatial information enabling informed traversability decisions.
3. **Dual-Method Wall Detection (92% recall):** Redundant sensing combining PathFinder floor segmentation with depth-based center analysis, reducing false negatives by 60% compared to single-method approaches. Prioritizes safety through sensor fusion.

Practical Impact: The system demonstrates that contemporary deep learning (YOLOv8n, MiDaS) can be unified into an auditory assistive framework achieving real-world performance at accessible cost (\$300-700 vs. \$5000+ for LiDAR systems). User studies showed 76% collision reduction, 60% speed improvement, and 4.5/5 satisfaction rating.

Comprehensive Evaluation: Validated through 50 turn trials, 30 gap measurements, 40 wall scenarios, and user studies with 5 participants across 3 diverse environments. All six bonus tasks completed including data-in-the-wild testing, ethical analysis, and public dataset release.

Significance: AssistedVision is the first assistive navigation system to combine behavioral validation (turn detection), quantitative spatial reasoning (gap widths in cm), and redundant safety mechanisms (dual-method wall detection) in a real-time framework optimized for consumer hardware. The clock-position guidance system with 8-second audio cooldown specifically addresses latency challenges in network-distributed architectures—a practical consideration absent in prior research.

Future Vision: Short-term enhancements include gyroscope integration and vibration feedback. Medium-term directions focus on smartphone-only deployment and longitudinal studies with visually impaired participants. Long-term vision encompasses reinforcement learning for personalized path planning and multi-modal sensor fusion. With appropriate empirical refinement and user-centered design, systems like AssistedVision could transform independent mobility for millions of visually impaired individuals worldwide.

Individual Reflections

This section contains individual reflections from each team member on their contributions, learnings, and experiences throughout the project.

Team Member 1: Lakshminarayan

My role throughout this project was comprehensive and multifaceted. I led the overall project coordination and management from start to finish, ensuring alignment across all team activities. I conducted extensive research to identify the most effective approaches for implementing the AssistedVision system, focusing on how to deliver clear and actionable instructions to visually impaired users. I explored various methodologies and technical frameworks to determine the optimal implementation strategy.

A significant part of my contribution involved investigating how to integrate Large Language Models (LLMs) to provide proactive, context-aware instructions. I researched potential applications of LLMs for generating dynamic guidance based on real-time environmental analysis, exploring how natural language processing could enhance the system's ability to communicate spatial information effectively. This exploration opened new directions for making assistive navigation more intuitive and responsive to user needs.

Throughout the project lifecycle, I coordinated work across multiple platforms and maintained continuous communication with the team to ensure project milestones were met. My research-driven approach helped shape the project's technical direction, from initial concept through final implementation. This reinforced the importance of considering human factors in computer vision systems. Pure algorithmic accuracy is insufficient—systems must account for human response time, behavioral patterns, and real-world constraints like network lag. I gained hands-on experience with Kalman filtering for object tracking and learned how to balance precision with real-time performance requirements.

Broader Impact Reflection: Working on assistive technology provided perspective on how computer vision can meaningfully improve quality of life. The 76% collision reduction and 60% speed improvement metrics represent tangible safety and independence gains for visually impaired users. This experience motivated me to prioritize accessibility considerations in future projects.

Team Member 2: Mahammad

Primary Contributions: Focused on researching and exploring how object detection and risk calculation mechanisms would function within the AssistedVision framework. Investigated YOLOv8-based detection approaches and developed understanding of risk assessment methodologies for obstacle navigation. Contributed to the integration of the final codebase by merging the complete system implementation, which included 23 files encompassing detection algorithms (`detection.py`), risk calculation modules (`risk.py`, `prob_risk.py`), path finding logic (`path_finder.py`), depth estimation (`depth.py`), tracking systems (`tracker.py`), and comprehensive documentation. This merge brought together over 3700 lines of code representing the culmination of the team's development efforts.

Technical Challenges: The most significant challenge was achieving accurate gap width estimation with monocular depth data. MiDaS provides relative depth, requiring careful calibration with camera FOV geometry. The law of cosines approach for angled configurations initially produced 20-30cm errors, which I reduced to 11.2cm MAE through projection correction and confidence scoring.

Key Learnings: I deepened my understanding of camera geometry, coordinate transformations, and the fundamental limitations of monocular depth estimation. The project taught me to design robust systems that gracefully handle uncertainty—for example, the confidence scoring mechanism that accounts for depth, angle, and gap width when providing traversability assessments.

Broader Impact Reflection: Quantifying navigable spaces is crucial for enabling independent mobility. The ability to tell a user "this gap is 75cm wide" versus simply "obstacle ahead" empowers informed decision-making. This project highlighted how precise, quantitative

information can restore agency to users who lack visual feedback.

Team Member 3: Mengfan Zhang

Primary Contributions: Researched and implemented masking techniques to isolate threat objects in the visual detection pipeline. Collaborated with Mohammad to develop risk logic algorithms and safety zone definitions for obstacle threat assessment. Explored various masking approaches to prioritize dangerous objects while filtering out non-threatening elements in the scene. Contributed to defining critical, high-alert, and caution zones based on distance thresholds and user safety standards. Worked on integrating the masking and risk assessment modules with the overall YOLOv8 detection system to enable real-time threat prioritization.

Technical Challenges: Developing effective masking techniques while maintaining real-time performance posed significant challenges. Balancing masking precision with computational efficiency required exploring various approaches to isolate threat objects without introducing latency. Defining appropriate risk zone boundaries that account for user walking speed, reaction time, and safety standards required integrating insights from proxemics research and robotics safety standards. Coordinating with Mohammad to ensure the masking and risk assessment modules integrated seamlessly with the detection pipeline demanded careful attention to data flow and timing.

Key Learnings: This project deepened my understanding of how threat prioritization in visual systems requires both spatial awareness and intelligent filtering. I learned that effective masking isn't just about isolating objects—it's about understanding context and determining which elements truly require user attention. The experience taught me how safety-critical systems must balance multiple factors: detection accuracy, computational speed, and human response capabilities. Collaborating on risk logic development highlighted the importance of evidence-based design, drawing from established research in human factors and robot safety standards.

Broader Impact Reflection: Working on threat detection and risk assessment for assistive navigation reinforced that technical accuracy alone is insufficient for safety-critical systems. The system must account for human cognitive load—providing too many warnings reduces effectiveness, while missing critical threats endangers users. This project demonstrated how assistive technology must be designed with deep understanding of user limitations and capabilities, not just technical metrics. The safety zones we defined translate directly to user independence and confidence in real-world navigation.

Team Member 4: Sudharsan Thiagarajan

Primary Contributions: Designed and implemented core system modules: Adaptive Turn Detection, Geometric Gap Quantification, and Dual-Method Wall Detection. Integrated YOLOv8n, MiDaS, and PathFinder into a unified pipeline optimized for CPU-only hardware. Built intuitive clock-position audio guidance and optimized performance to 28 FPS with less than 750 ms latency. Led experimental validation across indoor, office, and outdoor environments, ensuring robustness under varied lighting and obstacle conditions. Contributed to ethical analysis and proposed mitigation strategies for accessibility equity, bias in training data, and user autonomy.

Technical Challenges: Tackled network latency by validating user turns before issuing new instructions, reducing contradictory guidance. Reduced monocular depth ambiguity with geometric corrections and confidence scoring to achieve reliable gap estimation. Combined modalities to detect plain walls and uniform surfaces that single-method approaches missed. Balanced accuracy vs. efficiency trade-offs to achieve real-time performance on consumer hardware without GPU acceleration. Addressed user hesitation mid-turn, refining thresholds to minimize false positives in behavioral validation.

Key Learnings: Navigation systems must validate user behavior, not just obstacles, to ensure safe and intuitive guidance. Lightweight models and multi-modal fusion can deliver robust accuracy while remaining deployable on affordable hardware. Usability hinges on intuitive guidance formats (clock positions, gap widths) rather than raw technical metrics. Iterative testing with users revealed that trust and confidence in the system are as important as technical precision. Ethical considerations (privacy, accessibility, bias) must be integrated into design from the start, not treated as add-ons.

Broader Impact Reflection: Demonstrated that affordable, CPU-only solutions can rival costly LiDAR systems, broadening accessibility for visually impaired individuals. Enhanced autonomy and confidence by bridging perception with behavioral validation, reducing collisions by 76 percent and improving navigation speed by 60 percent. Planning to extend this work into a publishable paper and evolve the prototype into a marketable product with real-world deployment to establishment like Google or Meta. Long-term vision includes on-device smartphone deployment, multi-sensor fusion (gyroscope, ultrasonic, thermal), and reinforcement learning for personalized path planning. Beyond technical impact, the project reflects a commitment to inclusive design, ensuring that assistive technology empowers users rather than constrains them.

Team Member 5: Kumar Vaibhav

Primary Contributions: Led the comprehensive dataset creation effort for the AssistedVision project, recording videos across diverse environmental conditions and locations. Captured dataset videos using an Android smartphone at a consistent 4.5-foot camera height to simulate natural human perspective. Recorded videos across multiple lighting conditions (Morning, Evening/Noon, Night) and various indoor and outdoor locations to ensure environmental diversity. Created sensor-aware videos using the av-record-video.netlify.app tool to simultaneously capture accelerometer, gyroscope, device orientation, and timestamp data alongside video footage. Organized the complete dataset into a systematic 6-folder Google Drive structure (Morning, EveningNoon, Night, Newaddition, AV-RecordedVideos, Test) for easy access and version control. Validated all recordings through the processing pipeline to ensure compatibility and data quality.

Technical Challenges: Balancing environmental diversity with practical recording constraints required careful planning and execution. Ensuring consistent camera height while capturing varied lighting conditions and obstacles across different locations demanded meticulous attention to detail. Integrating sensor metadata collection (accelerometer, gyroscope, orientation) with video recording required learning and utilizing specialized web tools. Debugging the processing pipeline to handle sensor metadata integration involved modifying `run.sh` and `main.py` files to ensure full compatibility between the dataset structure and the processing workflow.

Key Learnings: This project reinforced that dataset quality fundamentally determines model performance in machine learning systems. I learned that systematic documentation and organized folder structures are as critical as the algorithms themselves for reproducible research. The experience highlighted the often-overlooked but essential work of data engineering in ML projects. I gained appreciation for how proper version control through cloud storage (Google Drive) enables seamless collaboration and ensures data accessibility for all team members. The importance of capturing diverse environmental conditions to improve model generalizability and reduce domain shift became evident through this work.

Broader Impact Reflection: Creating datasets for assistive technology made me conscious of representation—datasets must reflect the diverse real-world conditions users actually encounter, not idealized laboratory scenarios. The inclusion of sensor data (accelerometer, gyroscope, orientation) demonstrated how multimodal information can significantly improve system robustness, a principle applicable across many real-world vision applications. Understanding

that the dataset serves as the foundation for a system designed to enhance independence and safety for visually impaired users added meaningful purpose to the meticulous data collection work.

Acknowledgments

We thank the participants in our user study for their valuable feedback and time. We acknowledge the open-source communities behind PyTorch, YOLOv8, and MiDaS for their foundational contributions.

References

- [1] World Health Organization. **World Report on Vision**. WHO, 2019.
- [2] L. Kay. “A Sonar Aid to Enhance Spatial Perception of the Blind: Engineering Design and Evaluation.” **Radio and Electronic Engineer**, vol. 44, no. 11, pp. 605–627, 1974.
- [3] I. Ulrich and J. Borenstein. “The GuideCane - A Computerized Travel Aid for the Active Guidance of Blind Pedestrians.” **IEEE International Conference on Robotics and Automation**, pp. 1283–1288, 1997.
- [4] A. Hub, J. Diepstraten, and T. Ertl. “Design and Development of an Indoor Navigation and Object Identification System for the Blind.” **ACM SIGACCESS Conference on Assistive Technologies**, pp. 147–152, 2004.
- [5] J. Bai, Z. Liu, Y. Lin, Y. Li, S. Lian, and D. Liu. “Wearable Travel Aid for Environment Perception and Navigation of Visually Impaired People.” **Electronics**, vol. 6, no. 3, p. 59, 2017.
- [6] D. Das et al. “PathFinder: Monocular Depth-Based Pathfinding for Visually Impaired Navigation.” **arXiv preprint arXiv:2504.20976**, 2025.
- [7] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. “You Only Look Once: Unified, Real-Time Object Detection.” **IEEE Conference on Computer Vision and Pattern Recognition**, pp. 779–788, 2016.
- [8] J. Redmon and A. Farhadi. “YOLOv3: An Incremental Improvement.” **arXiv preprint arXiv:1804.02767**, 2018.
- [9] G. Jocher et al. “Ultralytics YOLOv8.” GitHub repository, 2023. <https://github.com/ultralytics/ultralytics>
- [10] R. Ranftl, K. Lasinger, D. Hafner, K. Schindler, and V. Koltun. “Towards Robust Monocular Depth Estimation: Mixing Datasets for Zero-shot Cross-dataset Transfer.” **IEEE Transactions on Pattern Analysis and Machine Intelligence**, vol. 44, no. 3, pp. 1623–1637, 2020.
- [11] R. Ranftl, A. Bochkovskiy, and V. Koltun. “Vision Transformers for Dense Prediction.” **IEEE International Conference on Computer Vision**, pp. 12159–12168, 2021.