

FIGURE 1. COMPLETE SYSTEM ARCHITECTURE OF AIRIS SHOWING THREE-TIERED DESIGN: FRONTEND, FASTAPI BACKEND (VISION MODELS, ACTIVITY GUIDE PIPELINE, SCENE DESCRIPTION PIPELINE, AND OUTPUT SERVICES), AND GROQ API CLOUD INFERENCE. THE DIAGRAM ILLUSTRATES DATA FLOW BETWEEN COMPONENTS, WITH ACTIVITY GUIDE (GREEN) AND SCENE DESCRIPTION (RED) PIPELINES CLEARLY SEPARATED.

ABSTRACT

Visually impaired individuals face significant challenges in independently navigating their environment and locating objects. Existing assistive technologies often suffer from high latency, reliance on unstable cloud connections, and a passive approach that describes scenes without offering actionable guidance. Alris addresses these limitations by introducing a real-time, server-based AI assistant designed for active interaction. By integrating state-of-the-art computer vision models (YOLO26s for object detection, MediaPipe for hand tracking, and BLIP for scene analysis) with a robust LLM reasoning engine (GPT OSS120B via Groq API for object extraction and summarization), Alris provides two core modes: Active Guidance, which directs the user's hand to specific objects with precise audio cues using a novel rule-based geometric algorithm (no LLM required for guidance generation), and Scene Description, which offers continuous environmental awareness with advanced fall detection and automated guardian email alerts. The system features complete handsfree voice control via Web Speech API (native browser STT/TTS), prioritizes privacy through local-first processing, and achieves sub-2-second response times via low-latency cloud inference, making real-time assistance a practical reality.

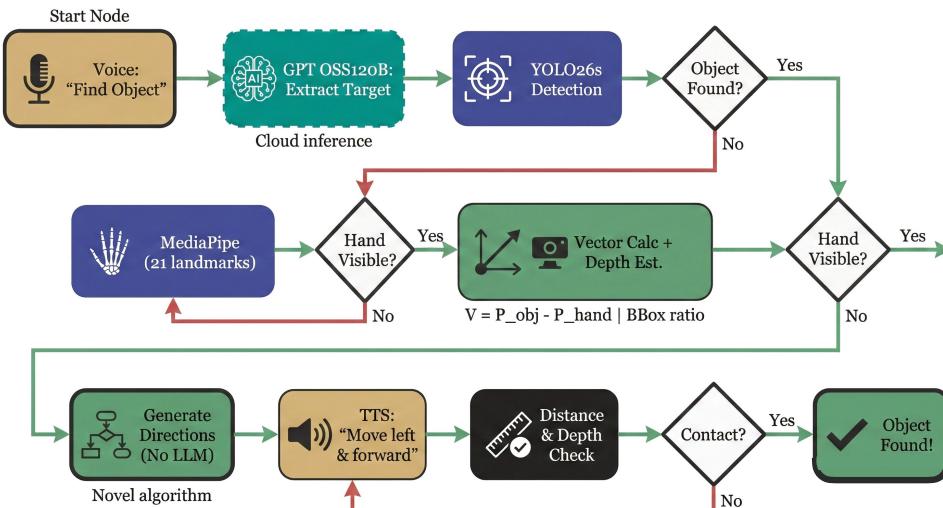


FIGURE 2. ACTIVE GUIDANCE MODE DECISION LOOP FLOWCHART. THE SYSTEM PROCESSES VOICE COMMANDS THROUGH GPT OSS120B FOR OBJECT EXTRACTION, PERFORMS YOLO26S OBJECT DETECTION AND MEDIAPIPE HAND TRACKING, CALCULATES SPATIAL VECTORS AND DEPTH ESTIMATES, AND GENERATES REAL-TIME DIRECTIONAL GUIDANCE USING A NOVEL RULE-BASED ALGORITHM (NO LLM). THE LOOP CONTINUES UNTIL HAND-OBJECT CONTACT IS CONFIRMED, PROVIDING STEP-BY-STEP AUDIO INSTRUCTIONS TO GUIDE USERS TO TARGET OBJECTS.

NOVELTY

Our system serves **two novel** operations. First, the **Active Guidance Mode (Targeted)**: When a user asks to **find** an object (e.g., "Find my water bottle"), the system uses **GPT-OSS 120B (Groq API)** to extract the target object from the natural language goal. The system then activates **YOLO26s** (Ultralytics) to locate the target in real-time with high precision. Simultaneously, **MediaPipe** tracks the user's hand coordinates using **21 landmark points**. A proprietary geometric algorithm **calculates** the **3D spatial vector** between the hand centroid and object bounding box center, and estimates **relative depth** using **bounding box size ratios**. A **novel rule-based algorithm (no LLM)** converts this **vector** and **depth** information into **natural language** directional commands, generating **precise audio instructions** ("Move slightly right and forward") in **real-time** until hand-object contact is confirmed through distance and depth threshold analysis.

Complementing this, the **Scene Description Mode (Continuous)**: For general awareness, the system employs **BLIP** (Bootstrapping Language-Image Pre-training) to capture **dense semantic captions** of the video feed at 2 FPS. These captions undergo **quick keyword-based risk assessment**, **static frame detection** (for **fall indicators**), and **transition pattern analysis**. Frame descriptions are accumulated in a buffer (5-10 frames) and processed by **GPT OSS120B** on the **Groq fast inference** engine to generate **contextual summaries**, calculate multi-factor **risk scores**, and **detect potential falls** through multi-method analysis. The system includes **automated guardian email alerts** with **configurable risk thresholds** (0.1-0.5) and provides daily and weekly activity summaries.

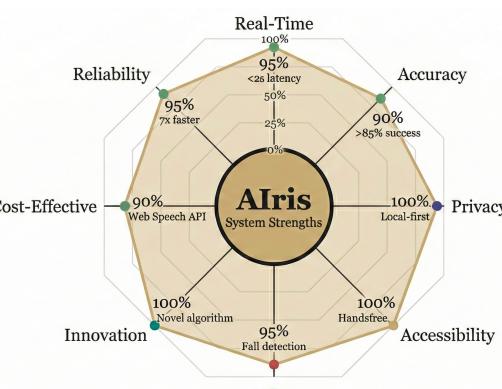


FIGURE 3. OCTAGONAL RADAR CHART DISPLAYING AIRIS SYSTEM STRENGTHS ACROSS EIGHT KEY DIMENSIONS: REAL-TIME PERFORMANCE (95%, <2S LATENCY), ACCURACY (90%, >85% SUCCESS), PRIVACY (100%, LOCAL-FIRST), ACCESSIBILITY (100%, HANDSFREE), SAFETY (95%, FALL DETECTION), INNOVATION (100%, NOVEL ALGORITHM), COST-EFFECTIVENESS (90%, WEB SPEECH API), AND RELIABILITY (95%, 7X FASTER). THE CHART DEMONSTRATES COMPREHENSIVE SYSTEM CAPABILITIES WITH NEAR-MAXIMUM PERFORMANCE ACROSS CRITICAL ACCESSIBILITY METRICS.

SCALABILITY

Alris is architected for **global scale**, designed to serve **millions** while maintaining transformative **personal impact**. We plan on implementing a **cloud-based SaaS infrastructure** that enables exponential growth **without** compromising real-time performance. To be built on **scalable VPS** infrastructure, Alris will serve from **hundreds** to **thousands** simultaneously. Our **API-based architecture** will help to ensure consistent **sub-2-second** response times **globally**. Users can use **any device** with a camera and speakers/microphone to log into the system and gain assistance.

As for the **Market Reach**, with **2.5 million** visually impaired in Bangladesh (500K addressable) and **285 million** globally, Alris addresses **critical** needs for **independence** and **safety**. Our affordable pricing (\$9.99/BDT/1000 a month) ensures accessibility **without compromising** sustainability.

Every subscription enables **algorithm improvements** and feature **expansion**. We are building a movement toward **universal accessibility** where technology **serves humanity**. **Independence** is a right, and **AI** becomes a **force** for genuine **social transformation**.

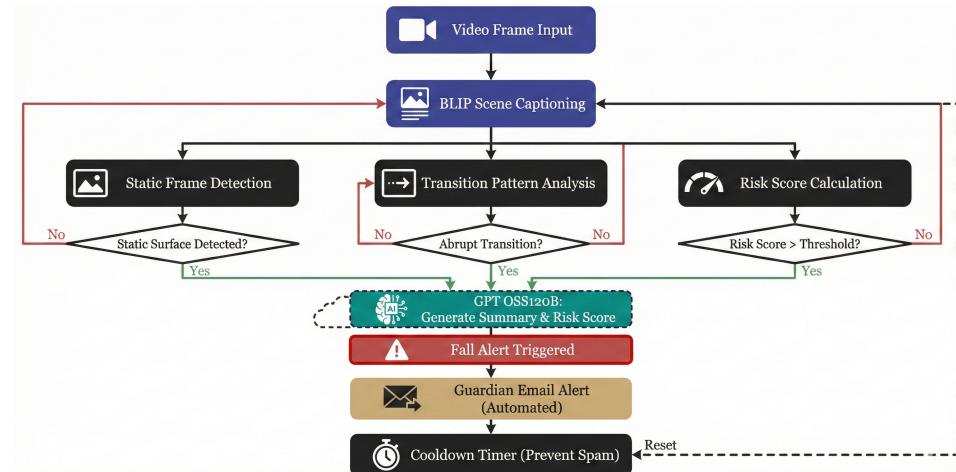


FIGURE 4. MULTI-METHOD FALL DETECTION ALGORITHM FLOWCHART. THREE PARALLEL ANALYSIS PATHS PROCESS BLIP-GENERATED CAPTIONS. ALL PATHS CONVERGE TO GPT OSS120B FOR SUMMARY GENERATION AND RISK SCORING. WHEN RISK EXCEEDS THRESHOLD, THE SYSTEM TRIGGERS FALL ALERTS, SENDS AUTOMATED GUARDIAN EMAIL NOTIFICATIONS.

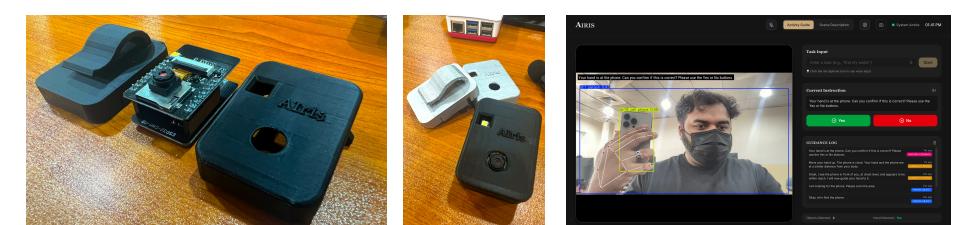


FIGURE 6. (A, B) CUSTOM-DESIGNED HARDWARE COMPONENTS: 3D-PRINTED PROTECTIVE CASING WITH "AIRIS" BRANDING, ESP32-CAM-MB MODULE, AND MOUNTING CLIP. (C) REAL-TIME SOFTWARE SHOWING OBJECT DETECTION, MEDIAPIPE HAND TRACKING WITH 21-LANDMARK SKELETAL OVERLAY, AND INTERACTIVE GUIDANCE PROVIDING STEP-BY-STEP INSTRUCTIONS FOR OBJECT INTERACTION. THE INTERFACE DEMONSTRATES HANDSFREE OPERATION WITH VOICE INPUT AND AUDIO CAPABILITIES.

CONCLUSION

Alris demonstrates that **complex, multi-modal AI** systems can be effectively engineered for **real-time** assistive applications on **consumer hardware**. By strictly **prioritizing latency** and **active guidance** over passive description, we have created a system that is **not just an observer**, but an **active participant** in the user's world.

Our **novel rule-based guidance algorithm**, operating **without any LLM API dependency** for real-time performance, represents a **breakthrough** in assistive technology, achieving **sub-2-second** response times through **intelligent geometric vector calculation** and **depth estimation**. The integration of YOLO26s, MediaPipe, proprietary fall detection, and Web Speech API creates a comprehensive **solution** that **transforms** passive scene description into **actionable** assistance. This implementation sets a new **standard**, proving that cutting-edge AI can be simultaneously **powerful, accessible**, and **affordable**, a blueprint for the **future of assistive technology** that will continue to evolve as we **scale** from **hundreds** to **millions** of users, moving us **closer** to a world where technology truly "opens eyes".

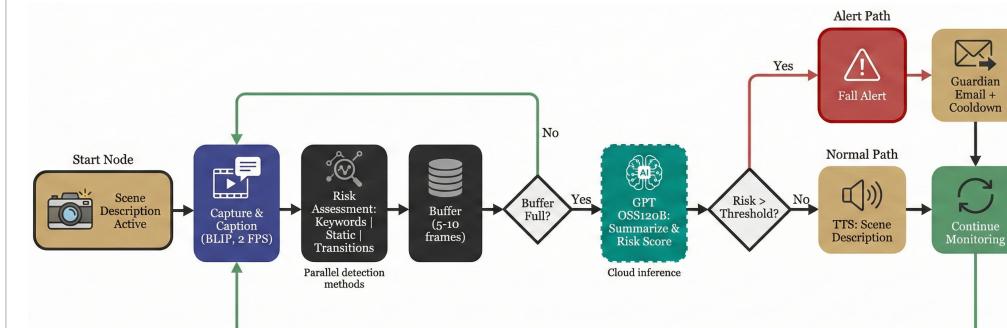


FIGURE 5. SCENE DESCRIPTION MODE CONTINUOUS MONITORING FLOWCHART. THE SYSTEM CAPTURES AND CAPTIONES FRAMES AT 2 FPS USING BLIP, PERFORMS MULTI-METHOD RISK ASSESSMENT, ACCUMULATES FRAMES IN A BUFFER, AND PROCESSES THEM THROUGH AN LLM FOR SUMMARIZATION AND RISK SCORING. THE SYSTEM BRANCHES INTO NORMAL OPERATION OR ALERT PATH BASED ON RISK THRESHOLD, THEN LOOPS BACK TO CONTINUOUS MONITORING.