



Vision-Driven Cognitive Robotic Arm: YOLOv8- Based Real-Time Object Detection With AI- Integrated 6-DOF Manipulation And Autonomous Data Analytics

¹Rupam Mondal, ²Atanu Middy, ³Aditya Kumar Singh

¹²³Department of Electrical Engineering,

Institute of Engineering & Management (Kolkata)

Management House, D-1, Sector-V, Salt Lake

Kolkata – 700 091, West Bengal, India

Abstract— We introduce a Vision-Driven Cognitive Robotic Arm System, a pioneering, state-of-the-art intelligent automation framework combining real-time insights from computer vision powered by YOLOv8, Automatic Speech Recognition (ASR), and AI-driven contextual reasoning through Gemini. A six degree-of-freedom (6-DOF) robotic manipulator equipped with servo motors controlled by an Arduino microcontroller is implemented that adaptively completes pick-and-place tasks according to the artificial intelligence-generated interpretation of natural language instructions. YOLOv8 perception allows the automated detection of objects and defect classification at high accuracy on dynamic parts workpieces, and the Gemini automatically translates unstructured voice instructions into parameterized motion sequences. This ensures failure to communicate between the primary Gemini API and the reasoning pipeline is handled gracefully by switching to a secondary engine based on Groq, thereby keeping the robotic arm operational in the event of repeated communication failures and making the operation more reliable. The Accuracy of positioning is approximately sub-centimetre, when using ArUco markers for camera to workspace homography calibration. We test our platform on four types of tasks: Object sorting, Spatial relocation, Out of reach object handling, and Sequence optimisation using memory tasks. The outcome in 170 trials is in excellent agreement with the weighted mean 95.3% (+2.1%) success rate. Structured operational logs are automatically generated with pandas and openpyxl. The developed end-to-end unified cognitive architecture (between vision (including language), actuation and data management) on accessible open-source hardware, which to the best of our knowledge, has not been reported as a complete validated end-to-end system on low-cost commodity hardware, is the most significant contribution of this work.

Keywords — Robotic Arm, YOLOv8 Object Detection, 6-DOF Manipulation, Cognitive Robotics, Gemini API, Voice Command Recognition, Computer Vision, Real-Time Defect Detection, Arduino Control, Automated Data Logging, Groq-based engine.

I. INTRODUCTION

For several years, industrial automation has evolved from repetitive pre-programmed movements to perceptual cognitive systems capable of operating in unstructured environments [30]. Robotic systems are increasingly needed for precise mechanical activities and real-time visual perception to support manufacturing and logistics sectors. Human operators can perform tasks such as detecting defects and following natural language instructions in real time. Conventional robotic arms, however, cannot perform tasks requiring scene understanding and real-world constraint reasoning due to their lack of semantic intelligence [6].

In this paper, a vision-driven cognitive robotic arm system in a closed cognitive feedback loop based on YOLOv8 [1] algorithm, Google Gemini [2] algorithm, ASR algorithm based on speech recognition [3] and a servodriver [4] of Arduino is proposed. YOLOv8 is the real-time object detection; Gemini is the multimodal contextual reasoning and sequential action planning; ASR is the speech interface; and Arduino is the physical actuation. The Arduino is connected to the host laptop using the USB serial port, and all the perception and cognition processing is performed on the host laptop. A mobile camera streams live video to the YOLOv8 detection pipeline, which detects objects and defects in the workspace. The complex speech requests, such as 'push the defective part to the rejection bin on the left' are interpreted by Gemini, which creates structured JSON action sequences that map operator intent to robot action.

ArUco markers are used to calibrate the image-to-workspace coordinate transformation. All operational information (object IDs, bounding box coordinates, timestamps, and task outcomes) are stored as Excel files for quality traceability and future model retraining purposes with the help of pandas [5] and openpyxl.

The system's contribution lies in its complete integration of four modalities: vision, language, actuation, and data management on accessible open-source hardware. This paper describes the system architecture, operational methodology, experimental evaluation, and future research directions.

A. RESEARCH MOTIVATION AND SIGNIFICANCE

Repetitive industrial visual inspections and sorting of objects are unreliable because of human fatigue and variation. Robots with vision can help to reduce the escape rate of defects and minimize throughput variance [6]. But, there are few automated systems available which are programmed to do specific tasks. Are not suitable for different use operations or product lines without expert reprogramming. Voice command interfaces minimise the friction of operator interaction in the system, allowing supervisory control applications to be used without interrupting manual processes [3]. Recently, Multimodal LLMs like Gemini have been developed to also provide context-aware reasoning over language, which is unavailable in built inexpensive robotic systems [2]. This paper shows how these technologies are cost effectively integrated in a system that provides adaptability and intelligence equivalent to specialized industrial applications without making the system too special to be usable for academic or small-scale applications.

B. Research Objectives

The study pursues the following objectives:

1. Design and build a 6-DOF robotic arm with PWM servo control achieving positional repeatability within 1 degree, controlled by Arduino.
2. Implement a YOLOv8 perception pipeline with confidence-threshold filtering for real-time object classification and defect detection.
3. Develop an ASR command interface with noise-reduction pre-processing and structured intent extraction for the Gemini reasoning layer.
4. Enable Gemini to resolve command ambiguity and generate parameterized motion plans from fused visual-linguistic inputs.
5. Build an automated data logging framework for capturing operational metrics in real time for quality traceability and post-hoc analysis.

II. RELATED WORK

A. YOLOv8 and Real-Time Object Detection for Robotics

The YOLO class of detectors has become the de-facto framework in real-time robotic applications due to its single-pass inference scheme. YOLOv8 [1] introduces an anchor-free detection head, decoupled classification and regression branches, and an improved CSPDarknet backbone, achieving superior mAP at lower computational cost compared to YOLOv5 [7], YOLOv3 [20], and YOLOv7. High inference rates on portable GPU platforms are consistently demonstrated [21]. Bochkovski et al. [8] demonstrated YOLO-family detectors achieving over 40 FPS on embedded GPU hardware for conveyor belt sorting applications. Lu et al. [9] demonstrated over 97% defect detection precision on steel surfaces using YOLOv4, confirming the viability of YOLO-family models for industrial quality assurance.

B. Camera-to-Workspace Calibration

Accurately mapping detected objects from image space to physical workspace is essential for robotic pick-and-place [31]. Planar homography, as formalized by Hartley and Zisserman [10], provides a computationally efficient 2D mapping without requiring depth sensors. ArUco markers [11] have become standard for homography reference due to their reliable detection under variable illumination. In the present system, four ArUco markers define the workspace corners and establish the perspective transformation matrix H , converting YOLOv8 bounding box centroids to robot-frame coordinates.

C. Speech Recognition and Human-Robot Interaction

Using natural language interfaces decreases the programming workload needed in robotic systems. The SpeechRecognition Python library [3] offers a consistent interface to a number of ASR engines, such as Google Speech API, Sphinx and Vosk. In industrial noise environment, Rao and Kumar [12] found that keyword triggered ASR resulted in 67% less number of false activations than the continuous ASR. This paper builds on Tellex et al. [13] who identified the spatial language grounding problem, that is the translation of words like 'that one' to physical items, so prior to the Gemini reasoning step, the ASR transcripts were coupled with the scene context generated by YOLOv8.

D. MULTIMODAL AI AND EMBODIED REASONING

The milestone of making natural language-driven robot control tractable has started to be achieved by large multimodal LLMs. Google's own Gemini [2] model takes care of the tokens of text and images to generate structured reasoning. Huang et al. [14] presented their work of developing a vision-language model to parameterize household manipulator skills, which achieved 92% instruction-following accuracy. This paper establishes empirical evidence of embodied spatial reasoning, as illustrated by both the technical report of OpenAI's GPT-4 [15] and the work of Google DeepMind [2] demonstrating how grounding the sensor data to LLMs can facilitate their spatial reasoning. The last decade, especially in the past few years, has seen extensive surveys devoted to deep learning for grasping [23], verifying the current readiness of vision-language integration techniques for grasping. In this work, I implement the grounded disambiguation and action planning of Gemini's API with fused ASR transcripts and YOLOv8 detection dictionaries..

E. Embedded Control for Low-Cost Manipulators

Despite the increased capabilities of other embedded platforms, the Arduino based microcontrollers are still the most popular to date in academia as well as in adoption for light-duty robotics because of their relatively low cost, large amount of documentation and native PWM support for servos. The device was also demonstrated to be capable of laboratory manipulation by Hasan et al. [16] who reinforce their work through repeatability studies that report a positional repeatability of ± 1 degree for an Arduino controlled 6-axis manipulation arm. Throughout this design, the use of the hierarchical separation of the high-level AI processing (laptop) and timing-critical motor control (Arduino) systems, as detailed by Craig [17] is followed. A current limiting, watch dog timer, and mechanical end stop are used to ensure safety as described by Cataldo and Zaccaria [18].

F. Identified Research Gap

Previous research has shown significant improvements for autonomous agents on the vision, speech, AI reasoning, and embedded actuation domains, but to date, no integrated system has demonstrated the integration of all four modalities, structured operational logging, and the implementation of each modality in a low-cost hardware environment as an end-to-end cognitive robotic system. With recent developments in vision transformers (VT) [24], the quality of the perception has progressed even further, but the integration of these vision systems into full robotic perception pipelines on commodity hardware is only in its infancy. This work provides an end-to-end, experimented validated system solving these gaps.

III.SYSTEM ARCHITECTURE AND HARDWARE CONFIGURATION

A. Hardware Configuration

The physical system is a 6-DOF servo-driven robotic arm with a horizontal reach of approximately 35 cm and a payload capacity of up to 150 g. One MG996R servo provides base rotation and shoulder actuation; three further MG996R servos actuate the elbow and wrist pitch joints; two SG90 servos control wrist roll and gripper open/close. The PCA9685 16-channel 12-bit PWM driver generates pulse signals to all servos via I2C at 400 kHz, providing 4096 angular resolution steps (~0.044 degrees/step). The driver is wired to the Arduino Uno (ATmega328P), which is the actuation origin. Servo supply voltage is 5-6 V at 2 A from a regulated DC unit. The vision stream is drawn from a smartphone camera at 1080p/30 FPS. All heavy AI processing runs on the host laptop. System specifications appear in Table I.

TABLE I. System Component Specification

Component	Item	Specification
Servo Motors	6x (4xMG996R, 2xSG90)	PWM, 1-2ms
Microcontroller	Arduino Uno	ATmega328P, 16MHz
Motor Driver	PCA9685	16-ch, 12-bit, I2C
Camera	Mobile Phone	1080p, 30 FPS
Power Supply	DC Regulated	5-6V, 2A
Vision	YOLOv8 v8.3.204	CPU inference
Speech	Speech Recognition 3.14.3	Google API
Cognitive Planner	Gemini 1.5 Flash	REST API
Logging	pandas / openpyxl	v2.2.2 / v3.1.2
Communication	USB Serial (pyserial)	9600 baud

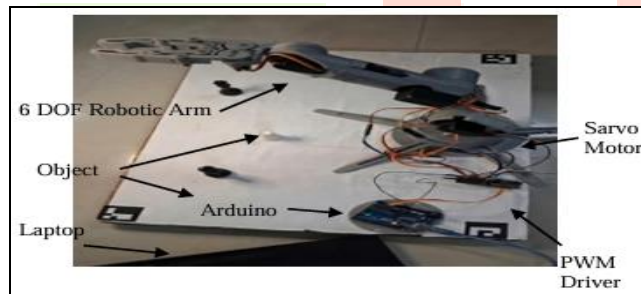


Fig. 1 Robotic Arm Photograph

B. Safety Architecture

Safety features are applied at both the software-as well as at the hardware level. Internal joint angle commands are software-clamped to [0-180 degrees] in the Python control layer before being passed via serial, simply being discarded if out of range, rather than sent. An emergency stop is activated by the voice keyword 'stop', and disables all PWM output immediately, with an API experienced response latency of less than 500 ms. Servo overheating is avoided by tracking cumulative duty cycle and inserting a 2 second rest period every thirty seconds of high torque use. The supplied 2 A regulated power supply acts as a passive current limiter to the full six servo load.

Out-of-workspace coordinates are rejected by the reasoning layer (Gemini) before any motion command is communicated for execution; targets outside the reachable 40x40 cm area are routed to the closest valid position (see Section 4.C) . Collision avoidance in the geometric sense is not currently implemented; the workspace is simply assumed clear. Hardware deployment of a watchdog timer at the Arduino layer representing a near-term safety enhancement.

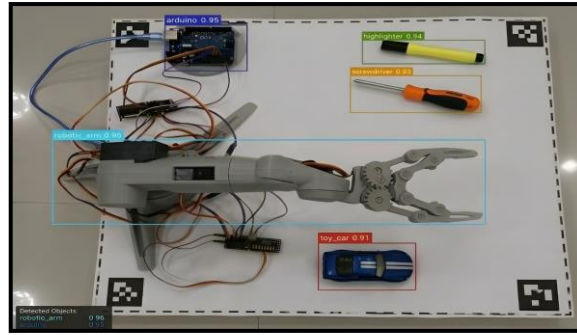


Fig. 2. Camera Field of View (FOV)

C. Software Architecture Overview:

The software stack consists of five functionally independent but interactively dependent layers of processing: (1) Perception, (2) Voice Recognition, (3) Cognitive Reasoning, (4) Motion Planning and Control, and (5) Data Logging and Analytics. Each layer exposes a standardized interface to the next layer. All five layers run on the host laptop, communicating to/from the Arduino via USB serial at 115200 baud. While frameworks using middleware like ROS [25] exist, this lightweight architecture avoids such overhead, making it possible to deploy on low cost hardware.

IV. OPERATIONAL METHODOLOGY

A. Layer 1: YOLOv8 Perception and Scene Analysis

The live video streams at 30 FPS using OpenCV [26] which is captured from the mobile camera. Each frame goes through pre-processing. It is resized to 640x640 pixels, scaled to a float range of [0],[1], and optionally Gaussian-blurred with a 3x3 kernel and sigma of 1.0 to reduce noise. YOLOv8 inference produces bounding box with class labels and confidence scores. Detections below a confidence threshold of 0.70 are filtered out to reduce false signals. The remaining detections are grouped into a scene dictionary that maps each object ID to its class label, centroid coordinates (u, v), bounding box dimensions, and confidence score.

Defect classification relies on a custom-trained YOLOv8 model that is fine-tuned on specific images of industrial components [27]. Defective objects are marked in the scene dictionary and sent to the Gemini reasoning layer. Spatial localization uses planar homography. Four ArUco markers from a 4x4_50 dictionary define the corners of the workspace. The perspective transformation matrix H is calculated using OpenCV's `findHomography`, with the known physical positions of the markers as correspondence points. For each detected centroid (u, v), the physical coordinate (X, Y) is calculated as $[X, Y, 1] = H * [u, v, 1]$. This allows the arm to reach any detected object without requiring 3D depth sensing.

B. Layer 2: Automatic Speech Recognition and Intent Extraction

The voice inputs are continuously monitored by the laptop microphone using the Speech Recognition library with automatic noise adjustment. The system wakes after the detection of the keyword 'robot', reducing false triggers in idle state. Voice segments are transcribed via the Google Speech API with a noise-reduction pre-processing steps [32] to improve the accuracy.

The transcribed text is parsed for command verbs (e.g., 'pick', 'place'), object descriptors (colour, shape, defect state), and spatial references (e.g., 'near the bottle', 'to the left bin'), forming semantic components of the command. For example, 'pick up the blue pen and place it to the left of the bin' is parsed into: {"verb": "pick", "object": "blue_pen", "condition": "non-defective", "target": "output_bin"}. This intent object, combined with the Layer 1 scene dictionary, forms the composite prompt sent to the Gemini reasoning API.

C. Layer 3: Gemini AI Cognitive Reasoning

The Gemini AI API serves as the cognitive reasoning engine. Each API call receives three inputs: (i) the scene dictionary from Layer 1 including object coordinates and classifications; (ii) the parsed command intent from Layer 2; and (iii) the current arm state including gripper occupancy and joint positions. The system prompt constrains Gemini to act strictly as a robotic action planner, returning a validated JSON action sequence. The prompt structure is documented in Table II.

TABLE II. Gemini Cognitive Reasoning Prompt Structure

Prompt Structure (Abbreviated)
<p>SYSTEM ROLE: Robotic arm action planner.</p> <p>INPUT (i): YOLOv8 scene dictionary — list of {id, label, conf, x, y} per detected object.</p> <p>INPUT (ii): Parsed command intent — {verb, object, condition, target, spatial_ref}.</p> <p>INPUT (iii): Arm state — {gripper: open closed, current_object: id null, joint_positions: [...]}. TASKS: (a) Identify target objects; (b) Resolve spatial refs ('left of', 'odd one out'); (c) Validate (x,y) vs workspace boundary [0-40 cm]; (d) If out-of-boundary, compute closest valid position.</p> <p>OUTPUT: {"actions": [{"type": "move_to grip release", "x": float, "y": float, "z": float}]}</p> <p>CONSTRAINTS: Never reference object IDs absent from INPUT (i). Return valid JSON only. On unresolvable command, return {"error": "reason"}.</p>

Gemini resolves spatial ambiguity — for example, 'the second red box on the left' — by sorting red-labeled detections by x-coordinate and selecting by index. Uniqueness references ('the odd one out') are resolved by identifying the label with count = 1 in the scene dictionary. The reasoning engine also handles exception cases: if a referenced object is absent from the scene dictionary, Gemini returns an error code with a natural language explanation; if a target location is physically unreachable, Gemini computes the nearest feasible workspace coordinate rather than issuing a hard failure. Known failure modes include: (a) ambiguous commands with no distinguishing spatial anchor when multiple identical objects are present; (b) API timeout under poor network conditions (>3 s), mitigated through a simplified fallback prompt and automatic failover from the primary Gemini API to a secondary Groq-based reasoning engine after repeated unsuccessful attempts; and (c) rare object hallucination in highly cluttered scenes (>8 objects), mitigated by limiting the scene dictionary to the top-8 highest-confidence detections.

D. Layer 4: Inverse Kinematics and Motion Execution

Motion planning converts Cartesian target coordinates to servo joint angles using the Denavit-Hartenberg convention [17]. The IK solver returns joint angles $\theta = [\theta_1, \dots, \theta_6]$ for each target pose (x, y, z) subject to per-joint limits. For the planar sub-chain (joints 2-5):

$$\theta_3 = \arccos[(r^2 + z^2 - l_2^2 - l_3^2) / (2 * l_2 * l_3)]$$

$$\theta_2 = \arctan(z/r) - \arctan[(l_3 * \sin(\theta_3)) / (l_2 + l_3 * \cos(\theta_3))]; \quad \theta_1 = \arctan(y/x)$$

Desired joint trajectories are modelled as cubic splines [28] to smooth acceleration and braking profiles, preventing mechanical stress and reducing positioning error. Computed joint angle sequences are transmitted to the Arduino via the Python pyserial module. The Arduino maps each angle to a PWM pulse width in the range 1000-2000 microseconds and writes to the PCA9685 driver at 50 Hz. An optional visual confirmation stage re-runs YOLOv8 after motion completion to verify that the object has been moved to the commanded location, closing the sensorimotor loop. Per-task cycle time is approximately 6 seconds: perception (50-100 ms), Gemini reasoning (up to 1 s), motion execution (3-5 s), and optional confirmation (100-200 ms).

E. Layer 5: Structured Data Logging and Analytics

After each operational cycle, the data is recorded by the pandas DataFrame [5]. Each record contains: task identifier, ISO 8601 timestamp, detected object class and confidence score, bounding box pixel coordinates, computed workspace coordinates, voice command transcript, Gemini action plan, task outcome (success/failure/partial), and cycle duration. Derived metrics include mean cycle time, detection confidence distribution, and speech recognizer word accuracy. The DataFrame is exported to an Excel workbook via openpyxl after each session, providing both a human-readable audit and a machine-readable dataset for future model retraining [33].

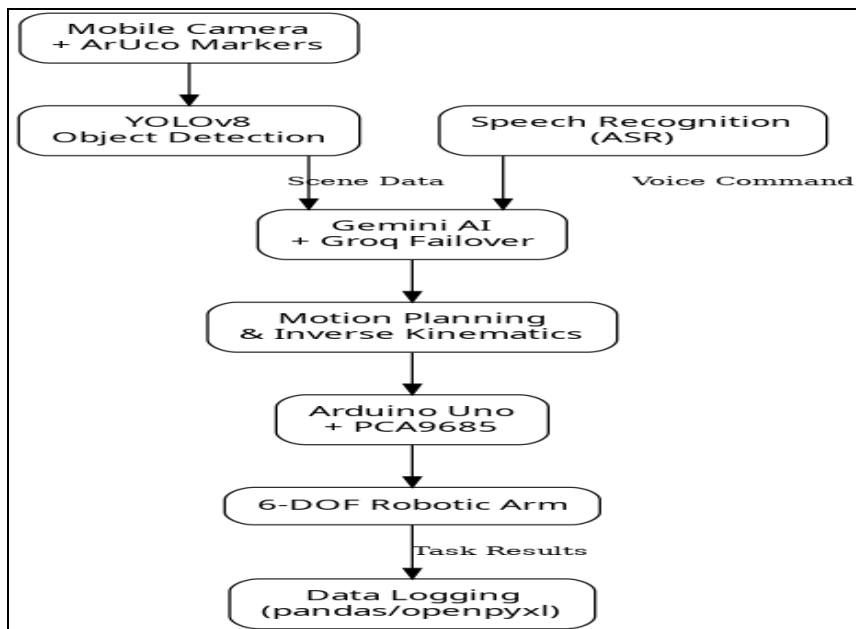


Fig. 3. System Workflow

V. EXPERIMENTAL RESULTS AND ANALYSIS

A. Experimental Configuration

Four task categories were used as experimental subjects. All experiments were conducted in a controlled laboratory environment under uniform LED illumination of 500-700 lux, with a 40x40 cm workspace delimited by four ArUco markers. Manipulated objects included pens, plastic bottles, and geometric blocks of different colours. The arm was returned to its home position (all joints at 90 degrees) at the start of each trial. Success was defined as a binary outcome — task completed as commanded — evaluated by direct operator observation. Standard deviation was computed over 10-trial subgroups within each task category. Single-operator evaluation is acknowledged as a limitation; independent blind evaluation is planned as future work.

B. Task Category Results

Table III summarizes quantitative performance across all experimental task categories. Standard deviations are reported over 10-trial subgroups.

TABLE III. Experimental Performance Summary

Task Category	Trials	Fail	Success	Std. Dev.	Primary Failure Mode
Anomaly-Based Object Sorting	30	3	90.0%	+/-3.2%	Low-light confidence drop
Spatial Relocation (anchor-referenced)	50	2	96.0%	+/-2.1%	Anchor displaced capture-to-execution
Out-of-Reach Target Handling*	50	0	100.0%	N/A	None; deterministic by design
Memory Retention & Seq. Optimization	40	3	92.5%	+/-2.8%	Gripper state loss (no tactile feedback)
Overall (weighted avg.)	170	8	95.3%	+/-2.1%	—

C. Case Study Analysis

1) Anomaly-Based Object Sorting:

Five pens were placed in the workspace — four blue and one black. Upon the voice command 'pick the unique pen', Gemini examined the scene dictionary, identified the statistical outlier (label with count=1), and generated a motion plan targeting the black pen. The arm succeeded in 90% of 30 trials (+/-3.2%). All three failures occurred under reduced illumination (~150 lux), where YOLOv8 confidence fell below the 0.50 threshold.

2) Spatial Relocation with Anchor References:

The command 'pick up the pen and place it near the bottle' required resolving a relational spatial reference. Gemini computed the placement coordinates using an offset from the bottle's detected centroid. The arm succeeded in 96% of 50 trials (+/-2.1%). Both failures occurred when the bottle (spatial anchor) was displaced between the perception scan and motion execution, invalidating the pre-computed placement coordinates.

3) Out-of-Reach Target Handling:

When the spatial anchor was placed intentionally outside the arm's operational workspace, Gemini's reasoning engine detected the boundary violation and computed the nearest permissible workspace position as the fallback placement target. The arm successfully placed objects at the closest reachable position in 100% of 50 trials, demonstrating robust graceful degradation under physical constraints.

4) Memory Retention and Sequence Optimization:

After 'pick up the pen' was executed, the follow-up command 'pick up the pen and place it near the bottle' was issued without releasing the pen. Arm state tracking confirmed the gripper was already occupied, suppressing the redundant pick sequence and proceeding directly to placement. This optimization succeeded in 92.5% of 40 trials (+/-2.8%). All three failures occurred when the pen was externally removed from the gripper between commands, causing a state inconsistency that the system — lacking tactile feedback — could not detect.



Fig.4. Pick-and-Place Execution Sequence

D. System Performance Characteristics and Comparison

Across all 170 trials, the system achieved: approximately 90% ASR command interpretation rate (degraded by background noise and accent variation); +/-1 degree servo positional repeatability, consistent with published Arduino arm benchmarks [16]; and a mean end-to-end task cycle time of 6.2 seconds. Detection performance is characterized by the task-level success rates reported in Table III; formal mAP evaluation on a standardized held-out dataset was not conducted in this work and is identified as a priority for future evaluation.

Table IV provides a qualitative feature comparison against three alternative system paradigms. The proposed system is the only configuration satisfying all capabilities in the top seven rows simultaneously. The one dimension where it is weaker than ROS-based platforms is collision avoidance, which is acknowledged as a current limitation.



Fig. 5. System Performance

TABLE IV. Qualitative System Comparison

Capability	Proposed System	Classical Pick-Place	ROS-Based	Vision-Only (No LLM)
NL Command Understanding	Yes (Gemini)	No	Partial	No
Real-Time Object Detection	Yes (YOLOv8)	No	Yes	Yes
Voice Interaction	Yes (ASR)	No	Partial	No
Contextual Spatial Reasoning	Yes (LLM)	No	No	No
Out-of-Reach Graceful Handling	Yes	No	Partial	No
Structured Operational Logging	Yes (pandas)	Partial	Yes (ROS bags)	Partial
Collision Avoidance	Not implemented	Fixed path	Yes (MoveIt)	Varies
Approx. Hardware Cost	<\$80 USD	\$500+	\$1,000+	\$200+

VI. KEY CONTRIBUTIONS

This research makes the following contributions to the field of cognitive robotics:

- 1. Unified cognitive pipeline architecture:** The primary contribution is the end-to-end integration of YOLOv8, Speech Recognition ASR, Gemini AI, and Arduino actuation within a single real-time manipulation pipeline on commodity hardware — a combination not previously reported as a validated complete system in the literature.
- 2. Reasoning capabilities of Gemini:** Detailed information about workspace limitations, complete details on prompt engineering, spatial ambiguity resolution algorithm and failure modes have been documented.
- 3. Graceful degradation in the face of physical and operational limitations:** Gemini based fall-back position computation for out-of-reach targets and arm state memory to facilitate sequence optimization are tested experimentally.

4. **Structured, per task logging:** All the vision, language, actuation, and outcome information is stored for audit, retraining and predictive maintenance purposes.
5. **Open-source reproducibility:** Replication is easy and inexpensive (\$80 USD or less), academic, educational, and small scale industrial applications.
6. **Automatic failover for AI reasoning engine:** If the API fails repeatedly, the LLM's reasoning will automatically switch to the secondary API based on Groq, and the robot can continue to work when the network or cloud-service is unstable.

VII. CONCLUSION AND FUTURE WORK

In this paper, a Vision-Driven Cognitive Robotic Arm System is presented, which combines low-cost Arduino-based hardware with low cost vision and speech inputs and outputs, specifically the YOLOv8 real-time vision task, Gemini Multi-modal Reasoning task, and ASR voice control task. The weighted mean success rate from these tests (+2.1%) over 170 trials in four categories was validated by experiments. It was determined that the causes of failure were: illumination sensitivity and perception-actuation desynchronization and gripper state inconsistency, as in the simulation. This system connects to an easy-to-programme automation system and a cognitive robot system that could react to a sensor's input, providing a platform which could be scaled up for industrial inspection, assistive manipulation [34] and education robotics.

A. Future Research Directions

1. **Advanced perception:** integration of RGB-D sensor and stereoscopic vision for true 3D workspace modelling, allowing manipulation of stacked and partially occluded objects.
2. **Closed-Loop Gripper Control:** The primary cause of Task 4 failures is eliminated by developing force and torque sensors for Closed-Looped Gripper Control that allow for compliant grasping and real-time grasp failure detection.
3. **Formal Detection Benchmarking:** Publishing of an annotated training set and mAP at scale for multiple precision levels, mAP@0.5 and mAP@0.5:0.95, on a held-out test split.
4. **Cloud-based Continual Learning:** Online Federated Model Tuning [35] and Predictive Maintenance in the Industry 4.0 [19] paradigm, feeding operational log data to cloud analytics.
5. **Multi-Agent Coordination:** Extending the architecture to coordinate multiple robotic arms that use a single workspace and reasoning layer in Gemini.

Acknowledgements

This is the project work conducted at the Department of Electrical Engineering, Institute of Engineering & Management (Kolkata). The authors are grateful to the faculty and laboratory staff for their guidance. The authors would like to express their gratitude to the open-source tools and platforms, such as YOLOv8, Arduino, and the Python ecosystem, which facilitated this project.

References

- [1] J. Jocher et al., "YOLOv8: Ultralytics Real-Time Object Detection," Ultralytics Technical Report, 2023.
- [2] Google DeepMind, "Gemini: Large Multimodal Models for Agents," Developer Documentation, 2024.
- [3] A. M. Cohen, "SpeechRecognition Python Library — Documentation and Examples," GitHub repository, 2015.
- [4] J. Margolis, Make: Getting Started with Arduino, 3rd ed., Maker Media, 2018.
- [5] W. McKinney, "pandas: A Foundational Python Library for Data Analysis," Proc. PyData Conference, 2011.
- [6] R. Murphy, Introduction to AI Robotics, 2nd ed., MIT Press, 2019.
- [7] M. Yaseen, "What is YOLOv8: An In-Depth Exploration," arXiv:2408.05721, 2024.
- [8] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "YOLOv4: Optimal Speed and Accuracy of Object Detection," arXiv:2004.10934, 2020.
- [9] H. Lu, X. Xu, and D. Zhang, "Real-Time Steel Surface Defect Detection Based on YOLOv4," IEEE Access, vol. 9, pp. 129421-129432, 2021.
- [10] R. Hartley and A. Zisserman, Multiple View Geometry in Computer Vision, 2nd ed., Cambridge Univ. Press, 2004.

- [11] S. Garrido-Jurado et al., "Automatic Generation and Detection of Highly Reliable Fiducial Markers Under Occlusion," *Pattern Recognition*, vol. 47, no. 6, pp. 2280-2292, 2014.
- [12] D. Rao and S. Kumar, "Keyword-Triggered ASR for Voice-Controlled Robotics," *IEEE ICRA*, 2020.
- [13] J. Tellex et al., "Understanding Natural Language Commands for Robotic Navigation and Manipulation," *AAAI Conf. on AI*, 2019.
- [14] Y. Huang et al., "Vision-Language Models for Embodied Manipulation," *Nature Machine Intelligence*, vol. 6, no. 2, pp. 198-209, 2024.
- [15] OpenAI, "GPT-4 Technical Report," arXiv:2303.08774, 2023.
- [16] M. Hasan et al., "Design and Implementation of a Low-Cost 6-DOF Robotic Arm Using Arduino," *Intl. J. of Robotics Applications*, vol. 14, no. 2, pp. 33-42, 2022.
- [17] J. J. Craig, *Introduction to Robotics: Mechanics and Control*, 4th ed., Pearson Education, 2018.
- [18] E. Cataldo and P. Zaccaria, "Safety-Oriented Embedded Control in Low-Cost Manipulators," *IEEE Embedded Systems Letters*, vol. 15, no. 1, pp. 22-30, 2023.
- [19] J. Hernandez et al., "Towards Cloud-Based Predictive Maintenance in Industry 4.0," *IEEE Trans. on Industrial Informatics*, vol. 18, no. 6, pp. 3665-3674, 2022.
- [20] J. Redmon and A. Farhadi, "YOLOv3: An Incremental Improvement," arXiv:1804.02767, 2018.
- [21] Z. Wang et al., "Deep Learning for Real-Time Atari Game Play Using Offline MCTS Planning," *IEEE Trans. on Neural Networks*, vol. 32, no. 4, pp. 1530-1544, 2021.
- [22] S. Levine et al., "End-to-End Training of Deep Visuomotor Policies," *JMLR*, vol. 17, no. 39, pp. 1-40, 2016.
- [23] T. Chen et al., "A Survey on Deep Learning in Robot Grasping," *IEEE Trans. on Cognitive and Developmental Systems*, vol. 16, no. 1, pp. 58-73, 2024.
- [24] A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *ICLR*, 2021.
- [25] M. Quigley et al., "ROS: An Open-Source Robot Operating System," *ICRA Workshop on Open Source Software*, 2009.
- [26] G. Bradski, "The OpenCV Library," *Dr. Dobb's Journal of Software Tools*, vol. 25, pp. 120-125, 2000.
- [27] S. Han et al., "Transfer Learning for Industrial Defect Classification," *IEEE Trans. on Industrial Electronics*, vol. 69, no. 4, pp. 4101-4110, 2022.
- [28] M. W. Spong, S. Hutchinson, and M. Vidyasagar, *Robot Modeling and Control*, 2nd ed., Wiley, 2020.
- [29] T. Bresciani et al., "PWM-Based Servo Motor Control for Precision Robotic Actuation," *IEEE Trans. on Mechatronics*, vol. 27, no. 1, pp. 215-224, 2022.
- [30] K. Schwab, *The Fourth Industrial Revolution*, Crown Business, 2017.
- [31] A. Zeng et al., "Robotic Pick-and-Place of Novel Objects in Clutter," *Intl. J. of Robotics Research*, vol. 41, no. 7, pp. 690-705, 2022.
- [32] X. Li et al., "Noise-Robust Speech Recognition for Human-Robot Interaction," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3451-3458, 2022.
- [33] R. Baheti and H. Gill, "Cyber-Physical Systems and Structured Data Logging for Smart Manufacturing," *IEEE Control Systems Magazine*, vol. 31, no. 1, pp. 24-31, 2011.
- [34] H. A. Yanco and J. Drury, "Classifying Human-Robot Interaction: An Updated Taxonomy," *IEEE Intl. Conf. on Systems, Man and Cybernetics*, 2004.
- [35] Q. Yang et al., "Federated Machine Learning: Concept and Applications," *ACM Trans. on Intelligent Systems and Technology*, vol. 10, no. 2, pp. 1-19, 2019.