# ARRC: Advanced Reasoning Robot Control—Knowledge-Driven Autonomous Manipulation Using Retrieval-Augmented Generation

Eugene Vorobiov, Ammar Jaleel Mahmood, Salim Rezvani, and Robin Chhabra

*Abstract*— We present ARRC (advanced reasoning robot control), a practical system that connects natural language instructions to safe, local robotic control by combining Retrieval-Augmented Generation (RAG) with RGB–D perception and guarded execution on an affordable robot arm. The system indexes curated robot knowledge (movement patterns, task templates, and safety heuristics) in a vector database, retrieves task-relevant context for each instruction, and conditions a large language model (LLM) to synthesize JSON-structured action plans. These plans are executed on a UFactory xArm 850 fitted with a Dynamixel-driven parallel gripper and an Intel RealSense D435 camera. Perception uses AprilTags detections fused with depth to produce object-centric metric poses; execution is enforced via a set of software safety gates (workspace bounds, speed/force caps, timeouts, and bounded retries). We describe the architecture, knowledge design, integration choices, and a reproducible evaluation protocol for tabletop scan/approach/pick–place tasks. Experimental results are reported to demonstrate efficacy of the proposed approach. Our design shows that RAG-based planning can substantially improve plan validity and adaptability while keeping perception and low-level control local to the robot.

## I. INTRODUCTION

Robotic manipulation is inherently dual in nature: it requires both high-level reasoning about objects and tasks, and precise low-level execution under physical constraints. Classical robotic systems are strong in motion planning and safety but brittle when adapting to new tasks or linguistic instructions. In contrast, vision-language-action (VLA) models such as RT-1 [1], PaLM-E [2], and RT-2 [3] excel in generalization and zero-shot task understanding, yet they frequently propose unsafe or infeasible plans due to limited embodiment grounding and lack of procedural rule enforcement.

Several works have sought to address these limitations. SayCan [4] integrates affordance estimators with language models to filter infeasible actions, while CLIPort [5] fuses semantic embeddings with transport operators for spatially grounded pick-and-place. VIMA [6] leverages multimodal prompts to generalize across manipulation tasks with few demonstrations. Although these approaches improve task grounding, they remain constrained by their training distributions and cannot readily incorporate new procedural or safety knowledge without retraining.

Retrieval-augmented generation (RAG), originally proposed for natural language processing (NLP) [7], enables

Authors are with the Department of Mechanical, Industrial and Mechatronics Engineering, Toronto Metropolitan University, Toronto, Canada
evorobiov@torontomu.ca
  ammar.j.mahmood@torontomu.ca
  salim.rezvani@torontomu.ca
  robin.chhabra@torontomu.ca

dynamic grounding of model outputs in external, updatable knowledge sources. Recently, robotics researchers have begun exploring retrieval-based compliance frameworks. Say-Comply [8], for instance, uses retrieval to enforce adherence to operational manuals in field robotics, demonstrating improved compliance and correctness. However, such approaches often stop at high-level compliance and do not integrate low-level metric grounding, perception-driven plan validation, or safety gating for execution.

In this paper, we bridge these gaps by introducing a RAG-enabled robotic manipulation pipeline–called ARRC (advanced reasoning robot control)–that unifies perception, retrieval, and safe plan execution. We deploy this system on a UFactory xArm 850 equipped with a RealSense D435 and a Dynamixel gripper, integrating retrieval of robot-centric safety heuristics and procedural templates at inference time. This design offers both adaptability and reliability, enabling the injection of new task knowledge or safety rules without retraining. Our contributions are as follows:

- We develop a hybrid VLA architecture integrating RAG for dynamic injection of procedural and safety knowledge.
- We demonstrate a real-world implementation with RGB-D perception, JSON-structured plan generation, and strict execution safety gates.
- We propose a reproducible ablation protocol to quantify contributions of retrieval, vision gating, and safety checks.

## II. RELATED WORK

### A. End-to-End and Vision-Language-Action Models

End-to-end VLA models have made significant strides in enabling robots to map visual observations and natural language instructions directly to motor actions. RT-1 ("Robotics Transformer") [1] trains a large-scale transformer on over 130,000 robot episodes across diverse tasks, demonstrating strong generalization to new objects and unseen task variations. RT-2 [3] and PaLM-E [2] extend this paradigm by integrating large multimodal models with robotic control, achieving impressive zero-shot capabilities. Other notable efforts include VIMA [6] and SayCan [4], which integrate language reasoning with affordance estimation or reward modeling. LM-Nav [9] further highlights the potential of language-conditioned navigation by combining pretrained vision-language-action models with planning modules for long-horizon mobility. Inner Monologue [10] explores reasoning and correction through explicit language feedback

during execution, enabling agents to adjust plans dynamically. Recent work on using GPT-4 for robotics [11] demonstrates how general-purpose LLMs can guide robotic task planning when provided with appropriate grounding and safety filtering. While these methods excel at leveraging large datasets to generalize across familiar tasks, they remain limited when new safety rules, procedural knowledge, or task constraints must be introduced post-training, often requiring expensive retraining or fine-tuning. Additionally, many end-to-end pipelines do not explicitly enforce low-level safety constraints, which can be critical in real-world robotic deployments.

### B. Affordance, Skill Library, and Hybrid Methods

Hybrid approaches attempt to combine learned policies with structured skill libraries or affordances. CLIPort [5] fuses visual-semantic embeddings with transport operators for spatially grounded pick-and-place, enabling more sample-efficient manipulation in tabletop settings. VIMA [6] leverages few-shot multimodal prompts to extend policy generalization across manipulation tasks. SayCan [4] couples high-level LLM plans with affordance and value estimators to filter infeasible actions, producing safer and more practical execution sequences. Neural task programming [12] and skill-chaining approaches [13] encode reusable motion primitives for modular policy composition. More recent methods, such as Huang *et al.* [14], show that LLMs can act as zero-shot planners for robotic tasks, though they require additional grounding for feasibility. LLM-BT [15] embeds structured behavior trees with large language models, improving task decomposition and policy reusability. While these approaches improve feasibility and generalization within known domains, they often rely on fixed libraries or pre-defined heuristics, limiting adaptability when new tasks, safety rules, or procedural knowledge need to be integrated dynamically.

### C. Retrieval-Based and Compliance-Grounded Planning

Retrieval-augmented generation (RAG) [7] in NLP improves grounding by combining a language model with an external knowledge corpus. Robotics researchers have begun adapting these ideas for compliant, knowledge-grounded planning. SayComply [8] retrieves compliance rules and manuals to constrain LLM outputs, improving task correctness and adherence to operational procedures. Similarly, behavior tree embeddings [16] capture structured task-level knowledge to support retrieval, task comparison, and indexing, but they do not integrate real-time perception or low-level control for safety-aware execution. Other related works include structured affordance retrieval [17], multi-step plan synthesis with semantic grounding [18], and Inner Monologue-style online corrections [10], which improve execution reliability but often do not incorporate validated real-robot execution or dynamic safety checks. Our approach extends these paradigms by combining retrieval-augmented LLM planning with local perception, validated execution, and low-level safety enforcement, allowing for flexible adaptation to new tasks, objects, and procedural constraints in real-world robotic environments.

### D. Gaps and Motivation for RAG-Driven Execution

From the survey, we identify several limitations in existing approaches. End-to-end VLA models exhibit strong generalization but are prohibitively expensive to update when introducing new tasks or safety constraints. Hybrid and affordance-based methods alleviate infeasibility, yet they remain restricted by fixed skills and affordance libraries. Retrieval-based compliance systems such as SayComply enforce high-level rules, but they rarely integrate metric perception, plan validation, or execution safety. Similarly, task-structure embeddings (e.g., BT embeddings) enable efficient knowledge indexing but often lack actionable primitives and safety gating within execution pipelines. To overcome these limitations, we present a new framework that unifies retrieval-augmented generation with robot-centric knowledge, metric perception, strict safety gating, and validated execution on real hardware—delivering both the flexibility of language-driven systems and the reliability required for safe robotic autonomy.

## III. System Overview

### A. Perception

AprilTags, combined with depth data from the Intel RealSense D435, provide marker-based detections that are fused to recover metric 3D poses in the robot frame. With sub-pixel accuracy, AprilTags enable robust and precise pose estimation, making them well-suited for manipulation tasks [19].

### B. Retrieval & Planning

We construct a curated robotics knowledge base comprising movement primitives, templates, safety heuristics, short demonstration transcripts, and parameterized affordances. This knowledge base is embedded using a SentenceTransformers model and indexed in ChromaDB (or a FAISS index). At inference time, the retrieval module selects the top-$k$ relevant context snippets, which are concatenated with the current observation summary and provided to an LLM (e.g., Gemini, PaLM-E-style models). The LLM then generates a structured JSON plan, enabling downstream execution [20].

### C. Execution

The JSON plan—represented as a sequence of named actions with bounded parameters—is first validated by a plan checker, synchronized with the latest object observations, and then executed through the XArm Python SDK. Execution is safeguarded by software safety gates, including workspace and joint limits, capped Cartesian speeds and accelerations, gripper torque/time gating, per-step timeouts, and bounded retries. Low-level controllers and perception remain local to the robot, while the LLM planner can be configured to run either locally (on-device) or via cloud APIs.

## IV. Proposed Method

In this section, we detail each system component and our design choices.

### A. Perception: AprilTags & Depth Fusion

The perception pipeline employs AprilTag detection (Tag-Standard41h12 family) to generate marker-based detections with unique object IDs. For each tag, we compute a robust depth estimate by taking the median of depth pixels within the detected tag region. This depth is back-projected through the calibrated camera intrinsics to obtain a 3D point in camera coordinates, which is subsequently transformed into the robot base frame via the extrinsic calibration $^{\text{base}}T_{\text{cam}}$. Each detected object is published as a compact observation message containing {tag_id, bbox, position_xyz, confidence}. A lightweight temporal filter smooths frame-to-frame jitter, reducing flicker and improving stability of the reported detections.

### B. Knowledge Base and Retrieval

The knowledge base is designed as a structured collection of short textual entries (50–400 tokens) that capture reusable domain knowledge for robotic manipulation. These entries span movement primitives (e.g., hover approaches, grasp offsets, retreat motions), parameterized task templates (such as SCAN AREA → APPROACH → GRASP → RETREAT), safety heuristics (safe hover heights, per-object velocity or force limits, and recovery strategies), as well as affordance notes derived from short demonstrations (e.g., "grasp handle from the side and lift with slow ascent").

Each document is embedded using a SentenceTransformers model from the Sentence-BERT family and stored in a vector index. In our implementation we employ ChromaDB, though FAISS or similar vector stores can serve as alternatives. At inference time, natural language instructions—optionally fused with compact visual summaries—are embedded and used to query the index. The retrieval step returns the top-$k$ semantically relevant documents, which are then passed forward to condition the planning module [21].

### C. RAG Prompting and JSON Plan Synthesis

The retrieved knowledge snippets are concatenated into a structured prompt that also incorporates a compact environment summary (object classes with recent positions), explicit system constraints (workspace limits, maximum speeds), and a JSON plan schema with few-shot exemplars. This prompt is submitted to a LLM, which in our prototype was accessed via a cloud API but can be replaced with any compatible model. The LLM generates a JSON-structured plan consisting of a goal description, optional reasoning metadata, and a sequence of action steps with bounded parameters (e.g., APPROACH_OBJECT {label:"bottle", hover_mm:30, timeout_sec:8}). By grounding generation in retrieval-augmented context, the system improves planning specificity and adaptability without retraining. The resulting plan is then parsed and validated by the plan checker before execution.

### D. Plan Representation, Validation, and Safety

Plans follow a restricted JSON schema:

```
{
  "goal": "place bottle on tray",
  "steps": [
    {"action":"SCAN_AREA",
     "params":{...}},
    {"action":"APPROACH_OBJECT",
     "params":{"label":"bottle",
               "hover_mm":40}},
    {"action":"MOVE_TO_POSE",
     "params":{"xyz_mm":[...],
               "rpy_deg":[...]}}
  ]
}
```

Before execution, each action step is validated against system constraints, including parameter bounds such as speeds and positions, and overall plan length. Steps that require current environmental information are synchronized with the latest perception data, and high-risk operations may optionally include human-in-the-loop confirmation. During execution, the executor enforces runtime safety through mechanisms such as per-step timeouts, gripper aborts triggered by load or duration limits, and emergency retreat in response to repeated failures.

### E. Gripper, Actuation, and Communication

We employ a parallel-jaw gripper actuated by Dynamixel motors (Protocol 2.0), providing precise and repeatable grasping for manipulation tasks. The software supports two configurations: a generic MX-series Dynamixel setup that allows direct control over position, speed, and torque, and an optional simplified driver tailored for the XL330-M288 model. Gripper parameters—including servo IDs, communication baudrate, open and closed positions, speed and torque limits, and load thresholds—are defined in config/gripper/gripper_config.yaml, allowing flexible adjustment for different tasks and objects. The controller enforces safe operational bounds on all commands and integrates load- and time-based gating to automatically detect grasp success or failure. This ensures that grasping actions remain robust under varying object geometries, weights, and environmental conditions, while preventing damage to the hardware.

### F. Executor API and Safety Gates

The executor exposes a compact set of high-level actions that the JSON planner can invoke, including SCAN_AREA, APPROACH_OBJECT, MOVE_TO_POSE, OPEN_GRIPPER, CLOSE_GRIPPER, RETREAT_Z, etc. Each high-level action is mapped to a sequence of atomic XArm SDK commands with pre-validated parameters. Runtime safety is enforced through multiple mechanisms: workspace boundaries and maximum single-move distances are respected, Cartesian and joint velocities and accelerations are capped, and gripper operations are monitored via load- and time-based thresholds

with emergency open procedures triggered on anomalies. Each step also supports bounded retries, and overall plan execution is subject to a configurable timeout.

Critically, the executor operates as the final stage of our end-to-end pipeline, translating LLM-generated JSON plans—conditioned on retrieval-augmented knowledge and real-time perception—into safe, executable robot motions. All safety parameters, including limits, thresholds, and retry counts, are fully configurable, allowing the system to adapt to different manipulators, end-effectors, and task environments. This integration ensures that high-level reasoning and planning performed by the LLM results in robust, reliable, and safe manipulation in real-world scenarios.

## V. ALGORITHMIC PROCESSES

This section details the core algorithmic components that enable the system's reasoning and execution capabilities. All the steps are described in the Figure.1.

### A. RAG-Based Planning Algorithm

The planning process follows a structured retrieval-augmented generation pipeline.

1) **Query Processing**: a natural-language instruction is embedded with SentenceTransformers to produce a query vector $q \in \mathbb{R}^d$.
2) **Knowledge Retrieval**: The query is compared against the knowledge-base embeddings $\{k_i\}_{i=1}^N$ using cosine similarity:

$$\text{sim}(q, k_i) = \frac{q^\top k_i}{\|q\|\|k_i\|} \qquad (1)$$

The top-$m$ most similar knowledge entries are retrieved and stored in the vector $\mathcal{R}$.
3) **Context Assembly**: Retrieved knowledge entries are concatenated with a symbolic summary of the current environment state $s_t$ (e.g., object classes and poses) to form the full planning context $C = [\mathcal{R}^\top, s_t^\top]^\top$.
4) **Plan Generation**: The context $C$ is passed to an LLM, which produces a JSON-structured plan $\pi = \{g, A\}$, where $g$ is the high-level goal and $A = \{a_1, \dots, a_l\}$ is a sequence of actions (action queue). Each action is represented as a tuple of an action label and bounded parameters (e.g., APPROACH_OBJECT {label: "bottle", hover_mm: 30, timeout_sec: 8}).

### B. Hierarchical Scanning Algorithm

The manipulator system employs a two-phase scanning routine for robust object detection. In the first phase, a **horizontal scan** is performed across the workspace at a fixed height, where the perception module inspects each sampled position for potential targets. If no objects are detected, the system transitions to a fallback **arc scan**, which uses three predetermined joint-space configurations (*LEFT*, *CENTER*, and *RIGHT*). These configurations are hardcoded to ensure safe, repeatable coverage of the workspace and avoid kinematic singularities, providing a deterministic alternative to Cartesian waypoint-based exploration.

### C. Coordinate Transformation Pipeline

Object coordinates obtained from AprilTags are first expressed in the eye-in-hand camera frame and then transformed into the robot base frame using calibrated transformations. These base-frame coordinates are subsequently used for inverse kinematics computations during manipulation. The transformed coordinates are stored in memory and retrieved on demand, depending on the specific objects referenced in the execution plan generated by the LLM.

We obtain the AprilTag position in the camera frame as $P_a^{\text{cam}} \in \mathbb{R}^3$, whenever the tag is in the field of view of the camera. The corresponding object position in the manipulator base frame $P_a^{\text{base}} \in \mathbb{R}^3$ is then computed via calibrated homogeneous transformations.

$$\begin{bmatrix} P_a^{\text{base}} \\ 1 \end{bmatrix} = H_{\text{TCP}}^{\text{base}}(\theta) H_{\text{cam}}^{\text{TCP}} \begin{bmatrix} P_a^{\text{cam}} \\ 1 \end{bmatrix}, \qquad (2)$$

where $H_{\text{TCP}}^{\text{base}}$ is the $4 \times 4$ forward kinematics transformation from the Tool Center Point (TCP) to the base frame, parameterized by the joint configuration $\theta$, and $H_{\text{cam}}^{\text{TCP}}$ is the fixed, calibrated extrinsic transformation from the eye-in-hand camera to the TCP.

Once the transformed position $P_a^{\text{base}} \in \mathbb{R}^3$ is obtained, it is validated against the robot's workspace bounds identified by the robot's reachable $\mathcal{W}$. Specifically, the position is checked to ensure it lies within the reachable volume of the manipulator, does not violate joint limits, and respects task-level safety constraints (e.g., avoiding collisions with the table or exceeding height limits). Only poses that pass these checks are forwarded to the inverse kinematics solver and subsequently used in the execution of the planned task.

### D. Safety-Constrained Execution

The execution module enforces a set of safety constraints to guarantee reliable operation during plan execution. First, all commanded target positions are validated against the robot's reachable workspace, ensuring all $P_a \in \mathcal{W}$. Motion commands are further regulated through velocity limits: Cartesian translations satisfy $\|\dot{P}\| \leq v_{\max}$, while joint velocities are constrained by $|\dot{\theta}_i| \leq \dot{\theta}_{i,\max}$. Each action step is assigned a maximum allowable duration $t_{\max}$; exceeding this limit triggers an abort and initiates a return-to-safe configuration. Gripper operations are protected through real-time load monitoring, with an emergency release executed if excessive force or unexpected obstruction is detected.

By integrating workspace validation, dynamic constraints, temporal bounds, and force-aware gripper control, the execution framework achieves robust, safe deployment of RAG-generated plans without sacrificing real-time responsiveness or hardware longevity.

## VI. EXPERIMENTAL SETUP

The experimental protocol is designed to provide a reproducible evaluation of the proposed manipulator system. We describe the hardware and software stack, the benchmark tasks and scenes, and the evaluation metrics used.

Fig. 1. High-level architecture. Perception produces object-centric observations; the RAG planner retrieves task knowledge and synthesizes a JSON plan; the executor validates and executes actions through the XArm SDK with safety gates.

## A. Hardware and Software

Experiments are conducted on a UFactory xArm 850 six-degree-of-freedom manipulator controlled over Ethernet through the official Python SDK. The manipulator is equipped with a parallel-jaw gripper driven by Dynamixel actuators, with both MX-series and an optional XL330 configuration supported. Perception is provided by an Intel RealSense D435 RGB–D camera, which is rigidly mounted and extrinsically calibrated with respect to the robot base frame, ensuring accurate recovery of 3D object poses from AprilTag detections. All computation runs on a host workstation running Ubuntu 22.04 and Python 3.10. The perception system is implemented in `vision_system/pose_recorder.py`, high-level planning is handled by `rag_system/true_rag_planner.py`, and low-level execution is carried out by `robot_controller/executor.py` interfacing with the robot through `robot_control/main.py`. ROS 2 is optionally used to publish vision messages, although the entire pipeline can run in a standalone lightweight mode. By default, LLM calls are executed via cloud APIs, but the system can be reconfigured to use a local LLM when available.

Table I summarizes the programmatically enforced con-

Fig. 2. A representative scene where the robot is performing manipulation.

straints and boundaries that regulate system resources and ensure stable performance.

| Parameter | Value |
|---|---|
| Max Cartesian speed | 150 mm/s |
| Max joint speed | 30 deg/s |
| Max Cartesian acceleration | 1000 mm/s$^2$ |
| Workspace $x$ | [150, 650] mm |
| Workspace $y$ | [-300, 300] mm |
| Workspace $z$ | [50, 500] mm |
| Gripper max speed | 300 (SDK units) |
| Gripper max force | 100 (SDK units) |
| Safety check interval | 0.1 s |
| Max single move distance | 400 mm |

*B. Tasks and Scenes*

We benchmark the system on three canonical tabletop manipulation tasks that reflect common challenges in domestic and industrial settings. The workspace is populated with a mix of everyday household objects, including bottles, boxes, and hand tools such as screwdrivers, arranged in moderate clutter to ensure that perception and motion planning must handle partial occlusions and close object proximities. All tasks are physically instantiated on a laboratory workbench, with object arrangements randomized across trials while remaining within a bounded workspace region to maintain consistency. Fig. 2 illustrates a representative cluttered scene where the robot performs a manipulation task..

The first task, **scan**, requires the manipulator to execute a workspace sweep using the mounted depth camera in order to detect and localize visible objects using AprilTags. The outcome is a set of object observations that are published for downstream planning. The second task, **approach**, evaluates the system's ability to condition on a class label such as "bottle" and to generate a collision-free plan that positions the end-effector at a stable hover pose above the identified

object. Finally, the **pick–place** task tests full manipulation capability: the robot must execute a grasp on the specified object, lift it from the tabletop, transport it to a designated drop location, and release it without premature slippage or collision.

To support reproducibility, scenes are generated from a limited set of objects and placement strategies, with identical layouts repeated across multiple trials for each task. This ensures that performance differences reflect system behavior rather than uncontrolled environmental variation, while still introducing sufficient diversity to challenge perception and planning modules. Fig. 3 demonstrates the adaptive reasoning pipeline that was used in the experiments to pickup and place a screwdriver.

Each condition is repeated for at least ten trials, with results reported as averages and standard deviations.

## VII. RESULTS

In the test case, an object was positioned to partially occlude a screwdriver at an angle. The system demonstrated advanced reasoning by: (1) initially performing a horizontal scan that failed to detect the target object due to occlusion, (2) automatically transitioning to an arc scan strategy to view the workspace from multiple angles, (3) successfully detecting the screwdriver from the arc scan position, and (4) executing a successful pick and place operation. This example showcases the system's ability to reason about occlusion and adapt its scanning strategy procedurally.

Representative failures we observed during prototyping include:

i) missed detections under partial occlusion (recovery: rescan + change point of view),

ii) low-confidence objects causing approach timeouts (recovery: fall back to scripted grasp),

iii) fragile objects requiring gentle close and tactile sensing (not yet instrumented).

Table. II provides a quantitative measure of the system's performance, which directly supports the superiority of the proposed ARRC framework.

**Plan Validity:** In 8 out of 10 trials, the plan was valid, which suggests the system consistently generates a correct and executable path. The two failures ("No" in trials 3 and 5) align with the "representative failures" mentioned, such as missed detections under partial occlusion, which could lead to a failed plan.

**Scan:** The "T" in all trials confirms that a target was successfully detected after a scan. This reinforces the example of the system's ability to adapt its scanning strategy to overcome occlusion.

**Approach (%):** This is a critical metric for a robotic manipulation task. The "Approach" value represents the accuracy of the robot to place the gripper in the vicinity of the object, which was quantified using the ground truth location of the objects. The high percentages across the board, with an average of 87.1%, indicate that the system is consistently successful in performing precise pick-and-place
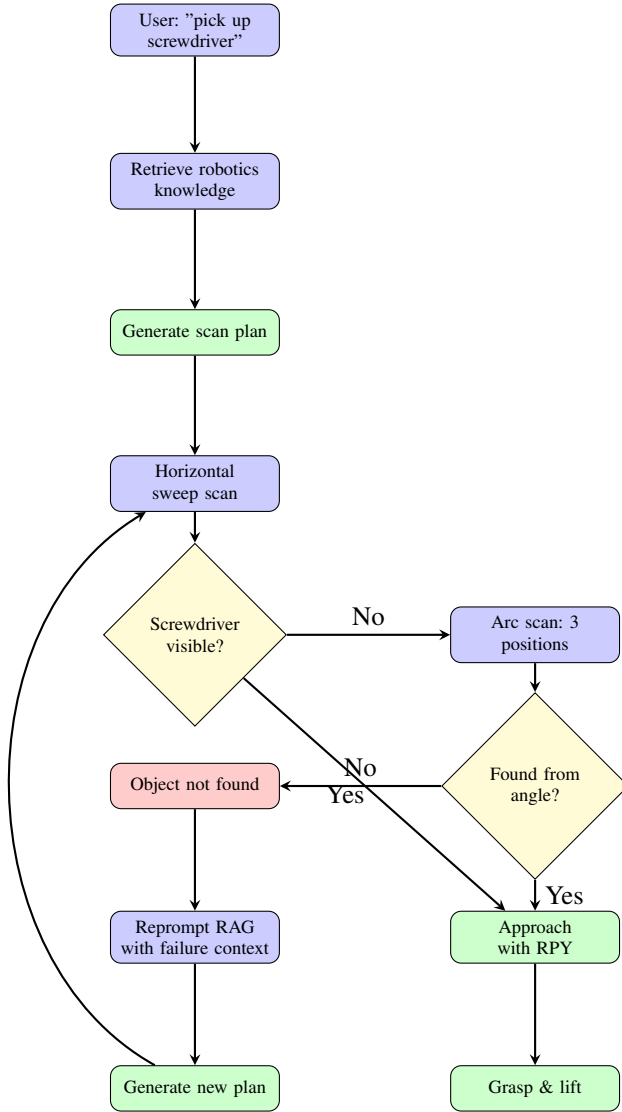
Fig. 3. Adaptive reasoning flow for "pick up the screwdriver" with failure recovery. The system demonstrates intelligent fallback from horizontal to arc scanning, and when objects remain undetected, it reprompts the RAG system with failure context to generate alternative strategies.

operations. This high accuracy is particularly noteworthy given the complexity of the tasks.

**P&P:** The "Yes" in 8 out of 10 trials confirms that the pick-and-place operation was successfully executed. The two "No" results correlate with the "Plan Validity" failures, as a failed plan would naturally prevent a successful P&P execution.

**Time Frame:** The time frame of the operations is relatively consistent, with an average of 1.225 seconds. This suggests that the system's reasoning and execution are not only accurate but also efficient.

The qualitative observations of system behavior are consistent with the quantitative results reported in Table II. For example, the screwdriver case under partial occlusion illustrates the system's ability to perform higher-level reasoning, ultimately yielding successful task execution. This

is supported by the high Approach and Pick–Place success rates recorded in the table.

Conversely, the representative failure modes, such as missed detections under occlusion, account for the two trials in which Plan Validity and Pick–Place were recorded as unsuccessful. The alignment between these qualitative insights and the quantitative measures provides a balanced assessment of both the strengths and current limitations of the proposed system.

## VIII. CONCLUSIONS

We have presented a practical, retrieval-augmented manipulation system that integrates local RGB–D perception, a vector-indexed robotics knowledge base, and LLM-driven plan synthesis to translate natural language instructions into validated, executable actions on a UFactory xArm 850. Experimental results demonstrate that retrieval-augmented generation improves plan specificity by providing concise, robot-relevant context, enabling more reliable and parameterized plans compared to LLM-only approaches. Furthermore, keeping perception and low-level control local proves effective for real-time execution and safety enforcement, while latency introduced by retrieval and LLM inference can be mitigated through capturing strategies.

Despite these strengths, the current system has some limitations. The knowledge base is curated and static, limiting adaptability to lifelong updates. Tasks requiring tactile feedback or precise torque control, such as dowel insertion or screw driving, remain out of scope. Additionally, the LLM planner may occasionally generate physically infeasible actions, highlighting the necessity for robust plan validation and symbolic feasibility checks. These limitations underscore important avenues for future research.

Future work will focus on scaling the system to more complex manipulation scenarios, including multi-arm coordination, large-scale benchmarking across diverse objects and lighting conditions, and integration of tactile sensing and force-aware motion primitives. On-device acceleration of retrieval and LLM inference is another promising direction to reduce latency and improve autonomy. Overall, our study demonstrates that combining retrieval-augmented language reasoning with local perception and safety-guarded execution provides a practical and reproducible pathway toward robust, adaptable robotic manipulation in real-world environments.

## REFERENCES

[1] A. Brohan *et al.*, "Rt-1: Robotics transformer for real-world control at scale," in *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
[2] D. Driess *et al.*, "Palm-e: An embodied multimodal language model," *arXiv preprint arXiv:2303.03378*, 2023.
[3] X. Chen *et al.*, "Rt-2: Vision-language-action models transfer web knowledge to robotic control," *arXiv preprint arXiv:2307.15818*, 2023.
[4] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," *arXiv preprint arXiv:2204.01691*, 2022.
[5] M. Shridhar *et al.*, "Cliport: What and where pathways for robotic manipulation," in *Conference on Robot Learning (CoRL)*, 2021.
[6] Y. Jiang *et al.*, "Vima: General robot manipulation with multimodal prompts," *arXiv preprint arXiv:2210.03094*, 2022.

TABLE II

RESULTS OF THE EXPERIMENTAL TRIALS, INCLUDING PLAN VALIDITY, APPROACH ACCURACY, AND PICK-AND-PLACE SUCCESS.

| # | Plan Validity | Scan | Approach Accuracy (%) | P&P | Time Frame |
|---|---|---|---|---|---|
| 1 | Yes | T | 86% | Yes | 1.23 |
| 2 | Yes | T | 86.80% | Yes | 1.21 |
| 3 | No | T | 80.30% | No | 1.22 |
| 4 | Yes | T | 88.20% | Yes | 1.24 |
| 5 | No | T | 79.30% | No | 1.21 |
| 6 | Yes | T | 88.90% | Yes | 1.23 |
| 7 | Yes | T | 96.10% | Yes | 1.23 |
| 8 | Yes | T | 89.20% | Yes | 1.23 |
| 9 | Yes | T | 88.20% | Yes | 1.22 |
| 10 | Yes | T | 88.80% | Yes | 1.23 |
| Success Rate. | 80% | 100% | 87.1% | 100% | 1.225 |

[7] P. Lewis *et al.*, "Retrieval-augmented generation for knowledge-intensive nlp tasks," in *NeurIPS*, 2020.

[8] M. F. Ginting, D.-K. Kim, S.-K. Kim, B. J. Krishna, M. J. Kochenderfer, S. Omidshafiei, and A.-a. Agha-mohammadi, "Saycomply: Grounding field robotic tasks in operational compliance through retrieval-based language models," in *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025, pp. 13 730–13 736.

[9] D. Shah, B. Osiński, S. Levine *et al.*, "Lm-nav: Robotic navigation with large pre-trained models of language, vision, and action," in *Conference on robot learning (CoRL)*. PMLR, 2023, pp. 492–504.

[10] W. Huang *et al.*, "Inner monologue: Embodied reasoning through planning with language models," in *6th Conference on Robot Learning (CoRL)*, vol. 205. PMLR, 2023, pp. 1769–1782.

[11] R. Mon-Williams, G. Li, R. Long, W. Du, and C. G. Lucas, "Embodied large language models enable robots to complete complex tasks in unpredictable environments," *Nature Machine Intelligence*, pp. 1–10, 2025.

[12] D. Xu, S. Nair, Y. Zhu, J. Gao, A. Garg, L. Fei-Fei, and S. Savarese, "Neural task programming: Learning to generalize across hierarchical tasks," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 3795–3802.

[13] T. Haarnoja, V. Pong, A. Zhou, M. Dalal, P. Abbeel, and S. Levine, "Composable deep reinforcement learning for robotic manipulation," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 6244–6251.

[14] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch, "Language models as zero-shot planners: Extracting actionable knowledge for embodied agents," in *International conference on machine learning*. PMLR, 2022, pp. 9118–9147.

[15] J. Ao, F. Wu, Y. Wu, A. Swiki, and S. Haddadin, "Llm-as-bt-planner: Leveraging llms for behavior tree generation in robot task planning," in *2025 IEEE International Conference on Robotics and Automation (ICRA)*, 2025, pp. 1233–1239.

[16] Y. Cao and C. G. Lee, "Behavior-tree embeddings for robot task-level knowledge," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 12 074–12 080.

[17] B. J. Gutierrez, N. McNeal, C. Washington, Y. Chen, L. Li, H. Sun, and Y. Su, "Thinking about gpt-3 in-context learning for biomedical ie? think again," *arXiv preprint arXiv:2203.08410*, 2022.

[18] M. Shridhar, L. Manuelli, and D. Fox, "Perceiver-actor: A multi-task transformer for robotic manipulation," in *Conference on Robot Learning (CoRL)*. PMLR, 2023, pp. 785–799.

[19] E. Olson, "Apriltag: A robust and flexible visual fiducial system," in *2011 IEEE International Conference on Robotics and Automation (ICRA)*, 2011, pp. 3400–3407.

[20] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, pp. 3982–3992, 2019.

[21] "Chroma: A vector database for building ai applications," 2023, accessed: Sep. 16, 2025. [Online]. Available: https://www.trychroma.com/