

AHA: A VISION-LANGUAGE-MODEL FOR DETECTING AND REASONING OVER FAILURES IN ROBOTIC MANIPULATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Robotic manipulation in open-world settings requires not only task execution but also the ability to detect and learn from failures. While recent advances in vision-language models (VLMs) and large language models (LLMs) have improved robots’ spatial reasoning and problem-solving abilities, they still struggle with failure recognition, limiting their real-world applicability. We introduce AHA, an open-source VLM designed to detect and reason about failures in robotic manipulation using natural language. By framing failure detection as a free-form reasoning task, AHA identifies failures and provides detailed, adaptable explanations across different robots, tasks, and environments. We fine-tuned AHA using FailGen, a scalable framework that generates the first large-scale dataset of robotic failure trajectories, the AHA dataset. FailGen achieves this by procedurally perturbing successful demonstrations from simulation. Despite being trained solely on the AHA dataset, AHA generalizes effectively to real-world failure datasets, robotic systems, and unseen tasks. It surpasses the second-best model (GPT-4o in-context learning) by 10.3% and exceeds the average performance of six compared models including five state-of-the-art VLMs by 35.3% across multiple metrics and datasets. We integrate AHA into three manipulation frameworks that utilize LLMs/VLMs for reinforcement learning, task and motion planning, and zero-shot trajectory generation. AHA’s failure feedback enhances these policies’ performances by refining dense reward functions, optimizing task planning, and improving sub-task verification, boosting task success rates by an average of 21.4% across all three tasks compared to GPT-4 models. Anonymous project page: aha-iclr.github.io.

1 INTRODUCTION

In recent years, foundation models have made remarkable progress across various domains, demonstrating their ability to handle open-world tasks (Driess et al., 2023; Alayrac et al., 2022; Achiam et al., 2023; Zhang et al., 2023). These models, including large language models (LLMs) and vision-language models (VLMs), have shown proficiency in interpreting and executing human language instructions (Ouyang et al., 2022), producing accurate predictions and achieving strong task performance. However, despite these advancements, key challenges remain—particularly with hallucinations, where models generate responses that deviate from truth. Unlike humans, who can intuitively detect and adjust for such errors, these models often lack the mechanisms for recognizing their own mistakes (Lin et al., 2021; Chen et al., 2021; Heyman, 2008).

Learning from failure is a fundamental aspect of human intelligence. Whether it’s a child learning to skate or perfecting a swing, the ability reason over failures is essential for improvement (Young, 2009; Gopnik, 2020; Heyman, 2008). The concept of improvement through failures is widely applied in training foundation models and is exemplified by techniques such as Reinforcement Learning with Human Feedback (RLHF) (Ouyang et al., 2022; Christiano et al., 2017), where human oversight and feedback steers models toward desired outcomes. This feedback loop plays a critical role in aligning generative models with real-world objectives. However, a crucial question persists: How can we enable these models to autonomously detect and reason about their own failures, particularly in robotics, where interactions and environments are stochastic and unpredictable?

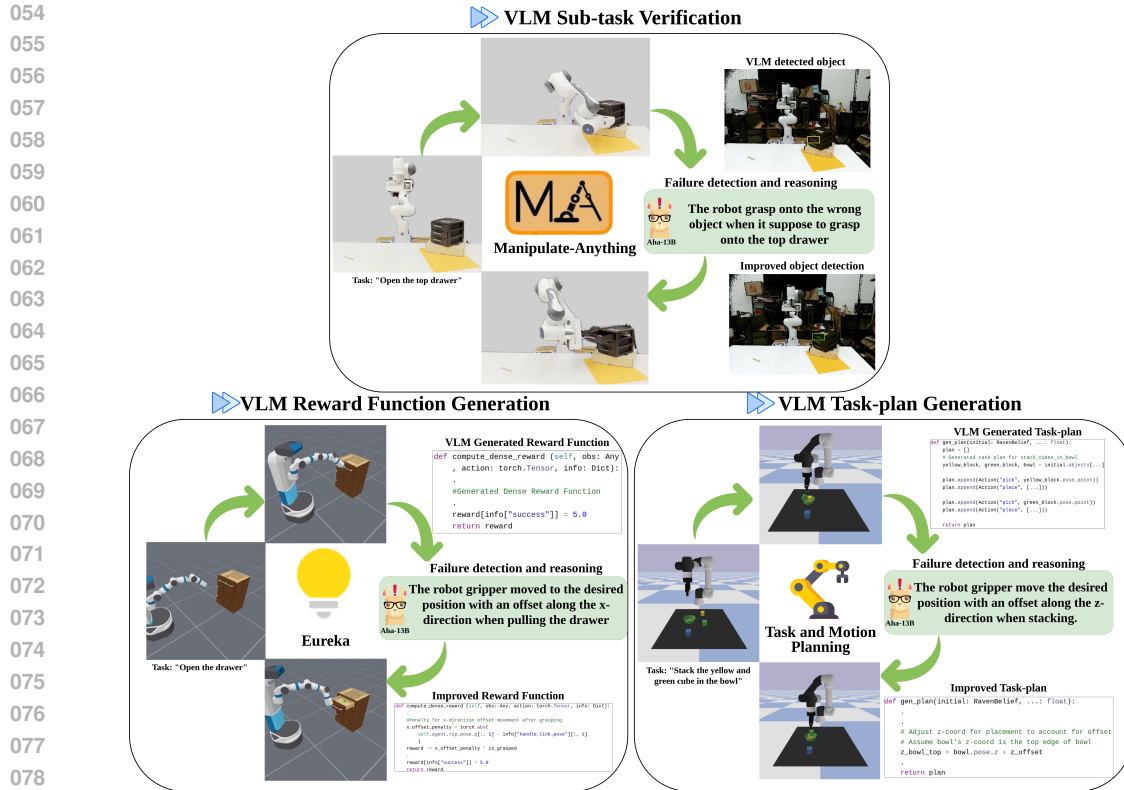


Figure 1: AHA is a Vision-Language Model designed to detect and reason about failures in robotic manipulation. As an instruction-tuned VLM, it can enhance task performance in robotic applications that utilize VLMs for reward generation, task planning, or sub-task verification. By incorporating AHA into the reasoning pipeline, these applications can achieve accelerated and improved performance.

This need is particularly pressing in robotics, where foundation models such as VLMs and LLMs are increasingly used to address open-world tasks. Recent advancements have enabled these models to tackle spatial reasoning, object recognition, and multimodal problem-solving—skills vital for robotic manipulation (Reid et al., 2024; OpenAI, 2024; Yuan et al., 2024; Chen et al., 2024; Wang et al., 2023b). VLMs and LLMs are already being integrated to automate reward generation for reinforcement learning (Ma et al., 2023; 2024), develop task plans for motion planning (Curtis et al., 2024), and even generate zero-shot robot trajectories (Huang et al., 2023; 2024a; Duan et al., 2024; Huang et al., 2024b). While these models excel at task execution, they often face challenges in detecting and reasoning over failures—skills that are crucial for navigating dynamic and complex environments. For example, if a robot drops an object mid-task, a human observer would immediately recognize the error and take corrective action. How can we empower robots with similar capabilities, allowing them not only to perform tasks but also to detect and learn from their mistakes?

To learn from their mistakes, robots must first detect and understand why they failed. We introduce AHA, an open-source VLM that uses natural language to detect and reason about failures in robotic manipulation. Unlike prior work that treats failure reasoning as a binary detection problem, we frame it as a free-form reasoning task, offering deeper insights into failure mode reasoning. Our model not only identifies failures but also generates detailed explanations. This approach enables AHA to adapt to various robots, camera viewpoints, tasks, and environments in both simulated and real-world scenarios. It can also be integrated into downstream robotic applications leveraging VLMs and LLMs, shown in Figure 1. We make the following three major contributions:

1. We introduce FailGen, a data generation pipeline for the procedural generation of failure demonstration data for robotic manipulation tasks across simulators. To instruction-tune AHA, we developed FailGen, the first automated data generation pipeline that procedurally creates the AHA dataset—a large-scale collection of robotic manipulation failures with over 49K+ image-query

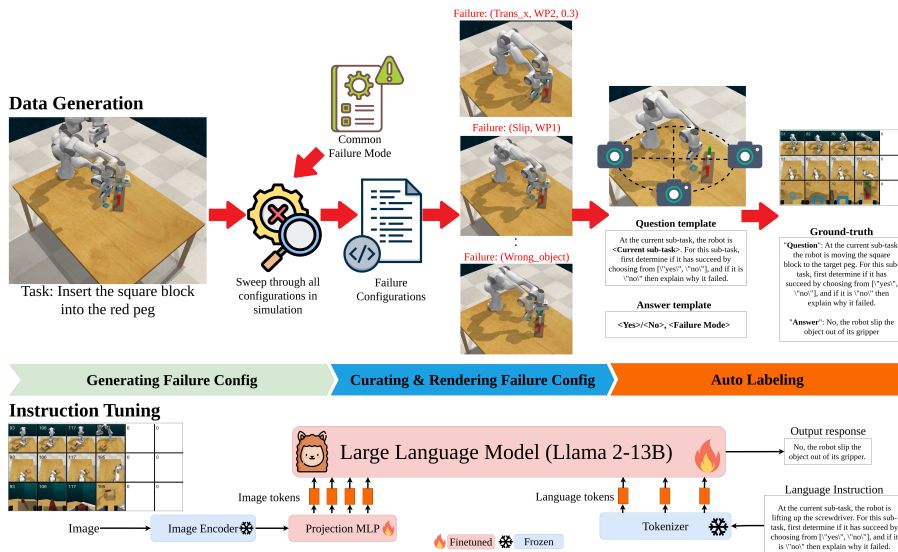


Figure 2: **Overview of AHA Pipeline.** (Top) The data generation for AHA is accomplished by taking a normal task trajectory in simulation and procedurally perturbing all keyframes using our taxonomy of failure modes. Through `FailGen`, we systematically alter keyframes to synthesize failure demonstrations conditioned on the original tasks. Simultaneously, we generate corresponding query and answer prompts for each task and failure mode, which are used for instruction-tuning. (Bottom) The instruction-tuning pipeline follows the same fine-tuning procedure as LLaVA-v1.5 Liu et al. (2023a), where we fine-tune only the LLM base model—in this case, LLaMA-2-13B and the projection linear layers, while freezing the image encoder and tokenizer.

pairs across 79 diverse simulated tasks. Despite being fine-tuned only on the AHA dataset, AHA demonstrates strong generalization to real-world failure datasets, different robotic systems, and unseen tasks, as evaluated on three separate datasets not included in the fine-tuning. `FailGen` is also flexible data generation pipeline integrates seamlessly with various simulators, enabling scalable procedural generation of failure demonstrations.

2. We demonstrate that AHA excels in failure reasoning, generalizing across different embodiments, unseen environments, and novel tasks, outperforming both open-source and proprietary VLMs. Upon fine-tuning AHA, we benchmarked it against six state-of-the-art VLMs, both open-source and proprietary, evaluating performance across four metrics on three diverse evaluation datasets, each featuring different embodiments, tasks, and environments out-of-distribution from the training data. AHA outperformed GPT-4o model by more than 20.0% on average across datasets and metrics, and by over 43.0% compared to LLaVA-v1.5-13B (Liu et al., 2023a), the base model from which AHA is derived. This demonstrates AHA’s exceptional ability to detect and reason about failures in robotic manipulation across embodiment and domains.

3. We show that AHA enhances downstream robotic applications by providing failure reasoning feedback. We demonstrate that AHA can be seamlessly integrated into robotic applications that utilize VLMs and LLMs. By providing failure feedback, AHA improves reward functions through Eureka reflection, enhances task and motion planning, and verifies sub-task success in zero-shot robotic manipulation. Across three downstream tasks, our approach achieved an average success rate 21.4% higher than GPT-4 models, highlighting AHA’s effectiveness in delivering accurate natural language failure feedback to improve task performance through error correction.

2 RELATED WORK

AHA enables language reasoning for failure detection in robotic manipulation, enhancing downstream robotics applications. To provide context, we review progress in: 1) failure detection in robotic manipulation, 2) data generation in robotics, and 3) foundation models for robotic manipulation.

Failure Detection in Robotic Manipulation. Failure detection and reasoning have long been studied in the Human-Robot Interaction (HRI) community (Ye et al., 2019; Khanna et al., 2023) and in works leveraging Task and Motion Planning (TAMP) (Garrett et al., 2020). With the recent widespread adoption of LLMs and VLMs in robot manipulation systems—either for generating reward functions or synthesizing robot trajectories (Ma et al., 2023; 2024) in a zero-shot manner—the importance of detecting task failures has regained prominence (Huang et al., 2023; Duan et al., 2024; Skreta et al., 2024; Ha et al., 2023). Most modern approaches focus on using off-the-shelf VLMs or LLMs as success detectors (Ma et al., 2022; Ha et al., 2023; Wang et al., 2023a; Duan et al., 2024), and some employ instruction-tuning of VLMs to detect failures (Du et al., 2023). Furthermore, hallucinations often occur in LLMs and VLMs. Methods that leverage these models for failure detection can mitigate this issue by detecting uncertainty in VLMs, as demonstrated in this work Zheng et al. (2024). However, these methods are often limited to binary success detection and does not provide language explanations for why failures occur. Our framework introduces failure reasoning in a new formulation, generating language-based explanations of failures to aid robotics systems that leverage VLMs and LLMs in downstream tasks. Additionally, we investigated whether AHA suffers from hallucinations by analyzing the prediction probabilities of sentence tokens. We found that AHA exhibits fewer hallucinations compared to other VLMs (see supplementary material).

Data Generation in Robotics There have been many methods in robotic manipulation that automate data generation of task demonstrations at scale (Mandlekar et al., 2023; Hoque et al., 2024), whether for training behavior cloning policies, instruction-tuning VLMs (Yuan et al., 2024), or curating benchmarks for evaluating robotic policies in simulation (Xie et al., 2024; Pumacay et al., 2024). A well-known example is MimicGen (Mandlekar et al., 2023), which automates task demonstration generation via trajectory adaptation by leveraging known object poses. Additionally, works like RoboPoint use simulation to generate general-purpose representations for robotic applications, specifically for fine-tuning VLMs. Similarly, systems like The Colosseum Pumacay et al. (2024) automate data generation for curating benchmarks in robotic manipulation. Our approach aligns closely with RoboPoint, as we also leverage simulation to generate data for instruction-tuning VLMs. However, unlike RoboPoint, we focus on synthesizing robotic actions in simulation rather than generating representations like bounding boxes or points.

Foundation Models for Robotic Manipulation. In recent years, leveraging foundation models for robotic manipulation has gained significant attention due to the effectiveness of LLMs/VLMs in interpreting open-world semantics and their ability to generalize across tasks (Duan et al., 2022; Hu et al., 2023; Firoozi et al., 2023; Urain et al., 2024). Two main approaches have emerged: the first uses VLMs and LLMs in a promptable manner, where visual prompts guide low-level action generation based on visual inputs (Liu et al., 2024a; Huang et al., 2024a;b). The second focuses on instruction-tuning VLMs for domain-specific tasks (Li et al., 2024). For example, RoboPoint (Yuan et al., 2024) is tuned for spatial affordance prediction, and Octopi (Yu et al., 2024) for physical reasoning using tactile images. These models generalize beyond their training data and integrate seamlessly into manipulation pipelines. Our approach follows this second path, developing a scalable method for generating instruction-tuning data in simulation and fine-tuning VLMs specialized in detecting and reasoning about robotic manipulation failures, with applications that extend beyond manipulation tasks to other robotic domains.



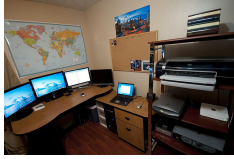
3 THE AHA DATASET

We leveraged `FailGen` to procedurally generate the AHA dataset from RL Bench tasks (James et al., 2020) and used it for the instruction-tuning of AHA. In this section, we begin by categorizing common failure modes in robotics manipulation and defining a taxonomy of failures in Section 3.1. Next, we explain how this taxonomy is used with `FailGen` to automate the data generation for the AHA dataset in simulation in Section 3.2.

3.1 FAILURE MODES IN ROBOTIC MANIPULATION

To curate an instruction-tuning dataset of failure trajectories for robotic manipulation tasks, we began by systematically identifying prevalent failure modes. Our approach involved a review of existing datasets, including DROID (Khazatsky et al., 2024) and Open-X Embodiment (Padalkar et al., 2023), as well as an analysis of policy rollouts from behavior cloning models. We examined

Table 1: **AHA datasets for instruction-tuning.** We combined the AHA dataset, our large-scale robotic manipulation failure dataset, with VQA and object detection data. By incorporating this diverse data mix into the fine-tuning process, AHA is able to reason about failures in robotic manipulation across different domains, embodiments, and tasks.

Source	The AHA dataset (Train)	VQA (Liu et al., 2023a)	LVIS (Gupta et al., 2019)
			
Quantity	49K	665K	100K
Query	For the given sub-tasks, first determine it has succeed by choosing from ["yes", "no"] and then explain the reason why the current sub-tasks has failed.	What is the cat doing in the image?	Find all instances of drawer.
Answer	No, The robot gripper rotated with an incorrect roll angle	The cat is sticking its head into a vase or container, possibly drinking water or investigating the interior of the item.	[(0.41, 0.68, 0.03, 0.05), (0.42, 0.73, 0.04, 0.08), ...]

failures occurring in both teleoperated and autonomous policies. Building upon prior works, such as REFLECT (Liu et al., 2023d), we formalized a taxonomy encompassing seven distinct failure modes commonly observed in robotic manipulation: incomplete grasp, inadequate grip retention, misaligned keyframe, incorrect rotation, missing rotation, wrong action sequence, and wrong target object.

Incomplete Grasp (No_Grasp) Failure: `No_Grasp` is an object-centric failure that occurs when the gripper reaches the desired grasp pose but fails to close before proceeding to the next keyframe.

Inadequate Grip Retention (Slip) Failure: `Slip` is an object-centric failure that happens after the object has been successfully grasped. As the gripper moves the object to the next task-specific keyframe, the grip loosens, causing the object to slip from the gripper.

Misaligned keyframe (Translation) Failure: This action-centric failure occurs when the gripper moves toward a task keyframe, but a translation offset along the X, Y, or Z axis causes the task to fail with respect to a fixed reference coordinate system.

Incorrect Rotation (Rotation) Failure: `Rotation` occurs when the gripper successfully reaches the correct position but rotates to an incorrect angle in roll, pitch, or yaw relative to a fixed reference point. Although it attempts the required rotation, the misalignment due to inaccurate rotation results in task failure.

Missing Rotation (No_Rotation) Failure: `No_Rotation` occurs when the gripper reaches the correct position but completely fails to perform the necessary rotation in roll, pitch, or yaw. The absence of any rotation when it is required leads to misalignment and ultimately causes the task to fail.

Wrong Action Sequence (Wrong_action) Failure: `Wrong_action` is an action-centric failure that occurs when the robot executes actions out of order, performing an action keyframe before the correct one. For example, in the task `put_cube_in_drawer`, the robot moves the cube toward the drawer before opening it, leading to task failure.

Wrong Target Object (Wrong_object) Failure: `Wrong_object` is an object-centric failure that occurs when the robot acts on the wrong target object, not matching the language instruction. For example, in the task `pick_the_red_cup`, the gripper picks up the green cup, causing failure.

3.2 IMPLEMENTATION OF THE AHA DATASET

The AHA dataset is generated with RL Bench James et al. (2020), utilizing its keyframe-based formulation to dynamically induce failure modes during task execution. RL Bench natively provides keyframes for task demonstrations, which enables flexibility in object manipulation (handling tasks

with varying objects) and the sequence of actions (altering the execution order of keyframes). Building on this foundation, we leverage `FailGen`, our custom environment wrapper around `RLBench` that allows for task-specific trajectory modifications through keyframes perturbations, object substitutions, and reordering of keyframe sequences. `FailGen` systematically generates failure trajectories aligned with the taxonomy defined in Section 3.1, yielding a curated dataset of 49k failure-question pairs.

To generate the AHA dataset, we systematically sweep through all keyframes in each `RLBench` task, considering all potential configurations of the seven failure modes that could result in overall task failure. By leveraging the success condition checker in the simulation, we procedurally generate YAML-based configuration files by sweeping through each failure mode across all keyframes. These files provide details on potential failure modes, parameters (such as distance, task sequence, gripper retention strength, etc.), and corresponding keyframes that `FailGen` should perturb to induce failure. Additionally, we incorporate language templates to describe what the robot is doing between consecutive keyframes. Using these descriptions along with the failure modes, we can systematically curate question-answer pairs for each corresponding failure mode.

For specific failure modes, `No_Grasp` is implemented by omitting gripper open/close commands at the relevant keyframes, effectively disabling gripper control. `Slip` introduces a timed release of the gripper shortly after activation. `Translation` and `Rotation` perturb the position and orientation of a keyframe, respectively, while `No_Rotation` constrains the keyframe’s rotational axis. `Wrong_Action` reorders keyframe activations to simulate incorrect sequencing, and `Wrong_Object` reassigns the keyframes intended for one object to another, maintaining the relative pose to mimic improper object manipulation. Using this pipeline, we also successfully generated a failure dataset from `ManiSkill` (Tao et al., 2024) and adapted `RoboFail` (Liu et al., 2023d) for the evaluation of AHA. This further demonstrates the generalizability and versatility of `FailGen` in generating failure cases across different simulation environments.

4 METHOD

This section outlines the failure reasoning problem formulation (Sec.4.1) used to fine-tune and evaluate AHA. Next, we discuss the curated data mix used for co-finetuning AHA (Sec.4.2). Finally, we detail the instruction fine-tuning pipeline and the model architecture selection for AHA (Sec.4.3).

4.1 FAILURE REASONING FORMULATION

We extend prior work (Skreta et al., 2024; Duan et al., 2024) by introducing a two-step framework for robot failure analysis that combines sub-task success detection and failure reasoning. Sub-task success is evaluated as a binary classification problem (*Yes/No*), while failure reasoning is performed using vision-language models (VLMs) to generate natural language explanations for the causes of failure. This approach allows for both precise failure detection and interpretability in robot manipulation tasks. Manipulation tasks are represented as trajectories consisting of a sequence of sub-tasks $\{S_0, S_1, \dots, S_T\}$, where each sub-task S_t is defined by two consecutive keyframes (K_t, K_{t+1}) . Each sub-task corresponds to an atomic manipulation action, such as “grasping a cube” in a stacking task. For each sub-task, the input to the VLM includes a query prompt and a structured image representation. The query prompt is generated using a template specific to the sub-task and describes the task context and success condition.

The image input is represented as a matrix $\mathbf{I} \in \mathbb{R}^{n \times T \times H \times W \times C}$, where rows correspond to camera viewpoints $\{V_0, V_1, \dots, V_{n-1}\}$ and columns correspond to temporal keyframes $\{K_0, K_1, \dots, K_T\}$. To capture the spatiotemporal progression of the task, frames are arranged in temporal order, and missing keyframes are replaced with white patches. We include several camera viewpoints to mitigate occlusions and ensure a comprehensive spatial context. This combined representation enables the VLM to reason over the robot’s trajectory and diagnose failure causes effectively, as demonstrated in Table 1.

4.2 SYNTHETIC DATA FOR INSTRUCTION-TUNING

To facilitate the instruction-tuning of AHA, we needed to systematically generate failure demonstration data. To achieve this, we developed `FailGen`, an environment wrapper that can be easily applied to

any robot manipulation simulator. `FailGen` systematically perturbs successful robot trajectories for manipulation tasks, transforming them into failure trajectories with various modes of failure as depicted in Figure 2 (Top image). Using `FailGen`, we curated the AHA dataset (Train) dataset by alternating across 79 different tasks in the RL Bench simulator, resulting in 49k failure image-text pairs. Furthermore, following proper instruction-tuning protocols for VLMs (Liu et al., 2023a) and building on prior works (Brohan et al., 2023; Yuan et al., 2024), co-finetuning is crucial to the success of instruction fine-tuning of VLMs. Therefore, in addition to the AHA dataset, we co-finetuned AHA with general visual question-answering (VQA) datasets sourced from internet data, which helps models retain pre-trained knowledge. Specifically, we included the VQA dataset (Liu et al., 2023a), containing 665k conversation pairs, and the LVIS dataset (Gupta et al., 2019), which comprises 100k instances with predicted bounding box centers and dimensions, as summarized in Table 1.

4.3 INSTRUCTION FINE-TUNING

We followed the instruction-tuning pipeline outlined by (Liu et al., 2023b). As depicted in Fig. 2, our model architecture includes an image encoder, a linear projector, a language tokenizer, and a transformer-based language model. The image encoder processes images into tokens, projected by a 2-layer linear projector into the same space as the language tokens. These multimodal tokens are then concatenated and passed through the language transformer. All components are initialized with pre-trained weights. During fine-tuning, only the projector and transformer weights are updated, while the vision encoder and tokenizer remain frozen. The model operates autoregressively, predicting response tokens and a special token marking the boundary between instruction and response.

4.4 IMPACT ON DOWNSTREAM TASKS

AHA integrates failure reasoning to address limitations in downstream robotics methods, improving reward synthesis, decision-making, and feedback efficiency. In reinforcement learning (RL), AHA refines reward synthesis by analyzing rollouts to provide failure explanations, enabling iterative adjustments to dense reward functions and improving sample efficiency, as demonstrated in approaches such as Eureka (Ma et al., 2023). In task and motion planning (TAMP) systems like PRO-C3S (Curtis et al., 2024), AHA enhances feedback loops by interpreting visualizations of failed plans, generating failure explanations, and informing language-model-based plan refinement. This process improves robustness in long-horizon tasks by addressing semantic errors overlooked by finite failure checks. In open-ended frameworks like Manipulate-Anything (Duan et al., 2024), AHA improves subtask verification by analyzing sequential frames for task progression errors, reducing failure propagation in zero-shot data generation. These integrations enable systematic reasoning improvements across RL, TAMP, and data generation, directly enhancing task success and robustness.

5 EXPERIMENTAL RESULTS

In this section, we evaluate AHA’s detection and reasoning performance against six state-of-the-art VLMs, including both open-source and proprietary models, some utilizing in-context learning. The evaluation spans three diverse datasets, covering out-of-domain tasks, various simulation environments, and cross-embodiment scenarios. We then assess AHA’s ability to retain general world knowledge after fine-tuning on domain-specific data. Finally, we explore its potential to improve downstream robotic manipulation applications.

5.1 EXPERIMENTAL SETUP

To quantitatively evaluate AHA’s detection and reasoning capabilities for failures in robotic manipulation, we curated two failure datasets and adapted an existing failure dataset for benchmarking. To ensure a fair comparison of free-form language reasoning, we also employed four different evaluation metrics to measure semantic similarity between sentences.

Benchmarks. We curated three datasets to evaluate AHA’s reasoning and failure detection capabilities, benchmarking against other state-of-the-art VLMs. The first dataset, AHA dataset (Test), includes 11k image-question pairs from 10 RL Bench tasks, generated similarly to the fine-tuning data via

Table 2: **Quantitative Evaluation on Failure Detection and Reasoning.** AHA-13B was evaluated and benchmarked against three open and three proprietary VLMs and one visual prompting baseline across three evaluation datasets. AHA-13B outperformed all other VLMs on every evaluation dataset and nearly every evaluation metric, with the exception of the AHA (Test) dataset, where GPT-4o exceeded by less than 3%.

Models	AHA dataset (Test set)				ManiSkill-Fail				REFLECT			
	ROUGE _L ↑	Cos Sim ↑	BinSucc(%) ↑	Fuzzy Match ↑	ROUGE _L ↑	Cos Sim ↑	BinSucc(%) ↑	Fuzzy Match ↑	ROUGE _L ↑	Cos Sim ↑	BinSucc(%) ↑	Fuzzy Match ↑
LLaVA-v1.5-13B (Liu et al., 2023a)	0.061	0.208	0.080	0.648	0.000	0.208	0.022	0.270	0.000	0.000	0.000	0.404
LLaVA-NeXT-34B (Liu et al., 2024b)	0.013	0.231	0.017	0.626	0.001	0.195	0.007	0.277	0.018	0.188	0.017	0.351
Owens-VL (Bai et al., 2023)	0.000	0.161	0.000	0.426	0.037	0.301	0.116	0.034	0.000	0.159	0.000	0.050
Gemini-1.5 Flash (Reid et al., 2024)	0.120	0.231	0.371	0.566	0.003	0.121	0.014	0.032	0.000	0.042	0.000	0.393
GPT-4o	0.251	0.308	0.500	0.784	0.142	0.335	0.688	0.453	0.114	0.318	0.554	0.438
GPT-4o-ICL (5-shot)	0.226	0.380	0.611	0.776	0.341	0.429	0.971	0.630	0.236	0.429	0.571	0.418
AHA-7B	0.434	0.574	0.691	0.695	0.609	0.680	1.000	0.532	0.204	0.394	0.625	0.439
AHA-13B (Ours)	0.446	0.583	0.702	0.768	0.600	0.681	1.000	0.633	0.280	0.471	0.643	0.465

Table 3: **Quantitative Evaluation on Standard VQA Benchmarks.** AHA-13B performs on par with LLaVA-13B Liu et al. (2023a), the VLM from which AHA adapts its fine-tuning strategy.

	MMBench (Liu et al., 2023c)	ScienceQA (Lu et al., 2022)	TextVQA (Singh et al., 2019)	POPE (Li et al., 2023)	VizWiz (Gurari et al., 2018)
LLaVA-13B (LLama-2) (Liu et al., 2023a)	67.70	73.21	67.40	88.00	53.01
AHA-13B (LLama-2)	65.20	71.94	65.20	85.74	53.45

FailGen (Section 3.2) but without overlapping with the finetuning dataset. It evaluates AHA’s ability to generalize to novel, out-of-domain tasks. The second dataset, ManiSkill-Fail, comprises 130 image-question pairs across four tasks in ManiSkill (Tao et al., 2024), generated using FailGen wrapper on the ManiSkill simulator. This dataset assesses AHA’s performance in a different simulator and under changing viewpoints. Lastly, we adapted a failure benchmark from the RoboFail dataset (Liu et al., 2023d), which features real-world robot failures in seven UR5 robot tasks, allowing for evaluation across simulation, real-world trajectories, and different embodiments.

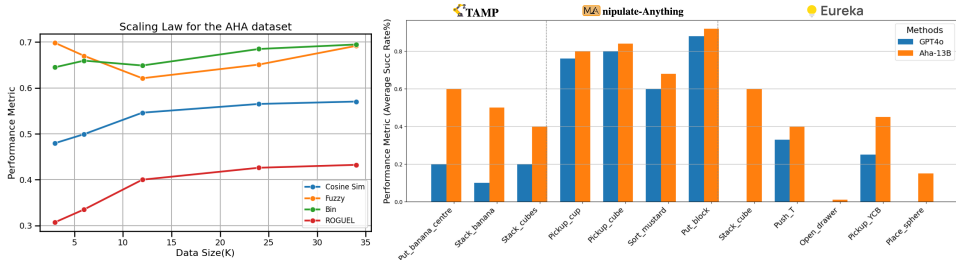


Figure 3: (Left) **Scaling law with the AHA dataset.** Scaling of effect of model performance with varying domain specific fine-tuning data. (Right) **Downstream Robotic Application Performance.** AHA-13B outperforms GPT-4o in reasoning about failures within these robotic applications, leading to improved performance of the downstream tasks.

Evaluation Metrics. To fairly evaluate success detection and language reasoning across all datasets and baselines, we employ four metrics. First, the **ROUGE-L score** measures the quality of generated text by focusing on the longest common subsequence between candidate and reference texts. Second, we use **Cosine Similarity** to assess similarity between texts or embeddings, avoiding the "curse of dimensionality". Third, **LLM Fuzzy Matching** utilizes an external language model—specifically, Anthropic’s unseen model, claude-3-sonnet—to evaluate semantic similarity in a teacher-student prompting format (Zhou et al., 2023). Lastly, we calculate a **Binary success rate** by comparing the model’s predictions directly against the ground truth for success detection.

5.2 QUANTITATIVE EXPERIMENTAL RESULTS

We contextualize the performance of AHA by conducting a systematic evaluation of failure reasoning and detection across these three datasets, general VQA datasets, and performed ablation studies.

AHA generalizes across embodiments, unseen environments, and novel tasks. To ensure fairness and eliminate bias in the detection and reasoning capabilities of AHA, we evaluated it on three different datasets that were never seen during fine-tuning, each designed to test a specific form of

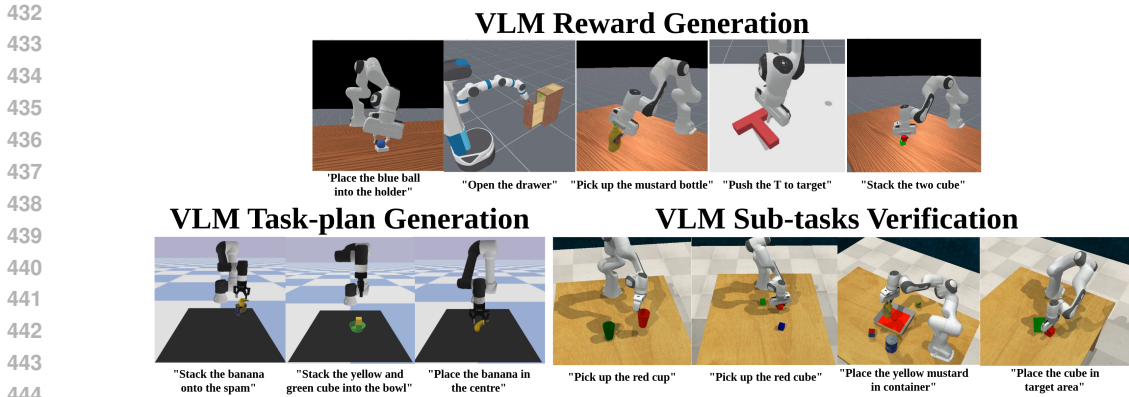


Figure 4: **Downstream Robotic Application.** We demonstrated that AHA can be integrated into existing LLM/VLM-assisted robotic applications to provide failure reasoning and feedback, helping to accelerate and improve task success rates in these systems.

generalization. First, on the AHA dataset (test) dataset, AHA demonstrated its ability to **generalize reasoning across tasks and new behaviors within the same domain, outperforming the second-best performing VLM, GPT-4o ICL**, by an average margin of 12.6% difference across all evaluation metrics. Second, we assessed AHA-13B on a dataset generated by the `Failgen` wrapper in a **different simulation domain**, ManiSkill, showing that our model outperforms GPT-4o-ICL by an average of 13.4% difference across all metrics as depicted in Table 2. Lastly, to demonstrate **generalization to real-world robots and different embodiments**, we evaluated AHA-13B on RoboFail (Liu et al., 2023d), where it outperforms GPT-4o-ICL by 4.9% difference.

AHA retains common sense knowledge. We evaluated AHA-13B’s performance on various VQA benchmarks and present the results in Table 3 . AHA-13B **performs comparably to LLaVA-v1.5-13B (LLama-2) (Liu et al., 2023a)** , with only a 1.5% margin difference as depicted in Table 3. Notably, LLaVA-v1.5-13B is a VLM trained on the same pre-trained weights as AHA-13B but fine-tuned on VQA data. This indicates that AHA-13B is capable of functioning as a general purpose VLM, in addition to excelling at failure reasoning.

AHA’s performance scales with data size. We evaluated Aha’s performance using a range of AHA data for instruction fine-tuning, spanning [3k, 6k, 12k, 34k, 48k, 60k], and co-trained individual checkpoints corresponding to these data sizes as shown in Figure 3 (Left). The model was then assessed on the ManiSkill-Fail dataset across four evaluation metrics. An average quadratic fit of 0.0022 across all four metrics demonstrates a **scaling effect with fine-tuning on our procedurally generated data pipeline**. This suggests that further scaling can improved model performance.

5.3 DOWNSTREAM ROBOTICS TASKS

We demonstrate that AHA’s failure detection and reasoning capabilities are useful across a wide spectrum of downstream robotics applications. This includes automatic reward generation for reinforcement learning applications (Ma et al., 2023), automatic task plan generation for task and motion planning applications (Curtis et al., 2024), and as an improved verification step for automatic data generation systems (Duan et al., 2024). Find videos, improved reward functions, task plans, and downstream application examples on the project page: supplementary materials or [aha-iclr.github.io](https://github.com/aha-iclr).

AHA enables efficient reward synthesis for reinforcement learning. To evaluate this downstream task, we adapted Eureka’s (Ma et al., 2023) implementation to the ManiSkill simulator, which offers more state-based manipulation tasks. We strictly followed the Eureka reward function generation and reflection pipeline, modifying it by incorporating perception failure feedback via either AHA-13B or GPT-4o (acting as a baseline) to enhance the original LLM reflection mechanism. Instead of only including a textual summary of reward quality based on policy training statistics for automated reward editing, we further incorporated explanations of policy failures based on evaluation rollouts. We evaluated our approach on five reinforcement learning tasks from ManiSkill, ranging from tabletop to mobile manipulation. To systematically assess the reasoning capabilities of different VLMs under

486 budget constraints, we sampled one reward function initially and allowed for iterations over two
487 sessions of GPT API calls. Each policy was trained using PPO over task-specific training steps and
488 evaluated across 1,000 test steps. During policy rollouts, we employed either AHA-13B or GPT-4o
489 for reward reflection to improve the reward function. Comparing the evaluated policy success rates
490 using different failure feedback VLMs, we observed that AHA-13B provided intuitive, human-level
491 failure reasoning that aided in modifying and improving generated dense reward functions. This
492 resulted in success across all five tasks within the budget constraints, and our approach **outperformed**
493 **GPT-4o by a significant margin of 22.34% in task success rate** shown in Figure 3 (Right).

494 **AHA refines task-plan generation for TAMP.** To demonstrate AHA’s utility within a planning
495 system, we incorporated our approach into PRoC3S (Curtis et al., 2024). The PRoC3S system solves
496 tasks specified in natural language by prompting an LLM for a Language-Model Program (LMP) that
497 generates plans, and then testing a large number of these plans within a simulator before executing
498 valid plans on a robot. If no valid plan can be found within a certain number of samples (100 in our
499 experiments), the LLM is re-prompted for a new LMP given failure information provided by the
500 environment. Importantly, as is typical of TAMP methods, the original approach checks for a finite
501 set of failures (inverse kinematics, collisions, etc.) from the environment, and returns any sampled
502 plan that does not fail in any of these ways. We incorporated a VLM into this pipeline in two ways:
503 (1) we prompt the VLM with visualizations of failed plan executions within the simulator, ask it to
504 return an explanation for the failure, and feed this back to PRoC3S’ LLM during the LMP feedback
505 stage, (2) after PRoC3S returns a valid plan, we provide a visualization of this to the VLM and ask
506 it to return whether this plan truly achieves the natural language goal, with replanning triggered
507 if not. We compared GPT-4o and AHA-13B as the VLM-based failure reasoning modules within
508 this implementation of PRoC3S across three tasks (shown in Figure 4). Each task was evaluated
509 over 10 trials, with a maximum of 100 sampling steps and three feedback cycles provided by either
510 GPT-4o or AHA-13B. The success rate for each task was recorded. As shown in Figure Figure 3
511 (Right), utilizing AHA-13B for **failure reasoning significantly improved the task success rate and**
outperforming GPT-4o by a substantial margin of 36.7%.

512 **AHA improves task verification for zero-shot robot data generation.** To demonstrate
513 AHA’s utility in zero-shot robot demonstration generation, we integrated our approach into the
514 Manipulate-Anything framework. This open-ended system employs various Vision-Language
515 Models (VLMs) to generate diverse robot trajectories and perform a wide range of manipula-
516 tion tasks without being constrained by predefined actions or scenarios. A critical component
517 of Manipulate-Anything is its sub-task verification module, which analyzes past and current
518 frames to decide whether a sub-task has been achieved before proceeding or re-iterating over the
519 previous sub-task. We replaced the original VLM (GPT-4V) in the sub-task verification module with
520 AHA-13B and evaluated performance across four RLbench tasks (Figure 4), conducting 25 episodes
521 for each task. Our results show that **substituting the sub-task verification module’s VLM with**
522 **AHA improved reasoning accuracy and overall task success by an average of 5%.**

524 6 CONCLUSION

526 **Limitations.** AHA currently outputs language reasoning that is closely aligned with the failure
527 scenarios in the fine-tuning data. However, there is an opportunity to output more open-ended failures,
528 to cover those arising from modes outside of the ones included in the failure taxonomy. Additionally,
529 while FailGen systematically curates failure data from simulations, distilling large pretrained
530 policies to perform diverse tasks in simulation and sampling failure modes would allow us to generate
531 more open-ended failure examples, potentially enhancing the instruction-tuned performance of AHA.

532 **Conclusion.** We introduce AHA, an open-source vision-language model that significantly enhances
533 robots’ ability to detect and reason about manipulation task failures using natural language. By
534 framing failure detection as a free-form reasoning task, AHA not only identifies failures but also
535 provides detailed explanations adaptable to various robots, tasks, and environments. Leveraging
536 FailGen and the curated AHA dataset, we trained AHA on a diverse set of robotic failure trajectories.
537 Our evaluations show that AHA outperforms existing models across multiple metrics and datasets.
538 When integrated into manipulation frameworks, its natural language feedback greatly improves error
539 recovery and policy performance compared to GPT-4 models. These results demonstrate AHA’s
effectiveness in enhancing task performance through accurate error detection and correction.

REFERENCES

- 540
541
542 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
543 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
544 *arXiv preprint arXiv:2303.08774*, 2023.
- 545 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
546 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language
547 model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736,
548 2022.
- 549 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang
550 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities.
551 *arXiv preprint arXiv:2308.12966*, 2023.
- 552 Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski,
553 Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action
554 models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- 555 Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh,
556 Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial
557 reasoning capabilities. *arXiv preprint arXiv:2401.12168*, 2024.
- 558 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared
559 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large
560 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 561 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
562 reinforcement learning from human preferences. *Advances in neural information processing
563 systems*, 30, 2017.
- 564 Aidan Curtis, Nishanth Kumar, Jing Cao, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. Trust the
565 proc3s: Solving long-horizon robotics problems with llms and constraint satisfaction, 2024. URL
566 <https://arxiv.org/abs/2406.05572>.
- 567 Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan
568 Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal
569 language model. *arXiv preprint arXiv:2303.03378*, 2023.
- 570 Yuqing Du, Ksenia Konyushkova, Misha Denil, Akhil Raju, Jessica Landon, Felix Hill, Nando
571 de Freitas, and Serkan Cabi. Vision-language models as success detectors. *arXiv preprint
572 arXiv:2303.07280*, 2023.
- 573 Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and Cheston Tan. A survey of embodied ai:
574 From simulators to research tasks. *IEEE Transactions on Emerging Topics in Computational
575 Intelligence*, 6(2):230–244, 2022.
- 576 Jiafei Duan, Wentao Yuan, Wilbert Pumacay, Yi Ru Wang, Kiana Ehsani, Dieter Fox, and Ranjay
577 Krishna. Manipulate-anything: Automating real-world robots using vision-language models. *arXiv
578 preprint arXiv:2406.18915*, 2024.
- 579 Roya Firoozi, Johnathan Tucker, Stephen Tian, Anirudha Majumdar, Jiankai Sun, Weiyu Liu, Yuke
580 Zhu, Shuran Song, Ashish Kapoor, Karol Hausman, et al. Foundation models in robotics: Applica-
581 tions, challenges, and the future. *arXiv preprint arXiv:2312.07843*, 2023.
- 582 Caelan Reed Garrett, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. Pddlstream: Integrating
583 symbolic planners and blackbox samplers via optimistic adaptive planning. In *Proceedings of the
584 international conference on automated planning and scheduling*, volume 30, pp. 440–448, 2020.
- 585 Alison Gopnik. Childhood as a solution to explore–exploit tensions. *Philosophical Transactions of
586 the Royal Society B*, 375(1803):20190502, 2020.
- 587 Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance
588 segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
589 recognition*, pp. 5356–5364, 2019.

- 594 Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and
595 Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In
596 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3608–3617,
597 2018.
- 598 Huy Ha, Pete Florence, and Shuran Song. Scaling up and distilling down: Language-guided robot
599 skill acquisition. In *Conference on Robot Learning*, pp. 3766–3777. PMLR, 2023.
- 600 Gail D Heyman. Children’s critical thinking when learning from others. *Current directions in
601 psychological science*, 17(5):344–347, 2008.
- 602 Ryan Hoque, Ajay Mandlekar, Caelan Garrett, Ken Goldberg, and Dieter Fox. Intervengen: Inter-
603 ventional data generation for robust and data-efficient robot imitation learning. *arXiv preprint
604 arXiv:2405.01472*, 2024.
- 605 Yafei Hu, Quanting Xie, Vidhi Jain, Jonathan Francis, Jay Patrikar, Nikhil Keetha, Seungchan Kim,
606 Yaqi Xie, Tianyi Zhang, Zhibo Zhao, et al. Toward general-purpose robots via foundation models:
607 A survey and meta-analysis. *arXiv preprint arXiv:2312.08782*, 2023.
- 608 Haoxu Huang, Fanqi Lin, Yingdong Hu, Shengjie Wang, and Yang Gao. Copa: General robotic
609 manipulation through spatial constraints of parts with foundation models. *arXiv preprint
610 arXiv:2403.08248*, 2024a.
- 611 Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer:
612 Composable 3d value maps for robotic manipulation with language models. *arXiv preprint
613 arXiv:2307.05973*, 2023.
- 614 Wenlong Huang, Chen Wang, Yunzhu Li, Ruohan Zhang, and Li Fei-Fei. Rekep: Spatio-temporal rea-
615 soning of relational keypoint constraints for robotic manipulation. *arXiv preprint arXiv:2409.01652*,
616 2024b.
- 617 Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot
618 learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–
619 3026, 2020.
- 620 Parag Khanna, Elmira Yadollahi, Mårten Björkman, Iolanda Leite, and Christian Smith. User study
621 exploring the role of explanation of failures by robots in human robot collaboration tasks. *arXiv
622 preprint arXiv:2303.16010*, 2023.
- 623 Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth
624 Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis,
625 et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*,
626 2024.
- 627 Xiang Li, Cristina Mata, Jongwoo Park, Kumara Kahatapitiya, Yoo Sung Jang, Jinghuan Shang,
628 Kanchana Ranasinghe, Ryan Burgert, Mu Cai, Yong Jae Lee, et al. Llara: Supercharging robot
629 learning data for vision-language policy. *arXiv preprint arXiv:2406.20095*, 2024.
- 630 Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object
631 hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023.
- 632 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human
633 falsehoods. *arXiv preprint arXiv:2109.07958*, 2021.
- 634 Fangchen Liu, Kuan Fang, Pieter Abbeel, and Sergey Levine. Moka: Open-vocabulary robotic
635 manipulation through mark-based visual prompting. *arXiv preprint arXiv:2403.03174*, 2024a.
- 636 Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction
637 tuning, 2023a.
- 638 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b.
- 639 Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
640 Llava-next: Improved reasoning, ocr, and world knowledge, 2024b.

- 648 Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi
649 Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player?
650 *arXiv preprint arXiv:2307.06281*, 2023c.
- 651 Zeyi Liu, Arpit Bahety, and Shuran Song. Reflect: Summarizing robot experiences for failure
652 explanation and correction. *arXiv preprint arXiv:2306.15724*, 2023d.
- 653
- 654 Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
655 Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
656 science question answering. In *The 36th Conference on Neural Information Processing Systems*
657 *(NeurIPS)*, 2022.
- 658
- 659 Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy
660 Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training.
661 *arXiv preprint arXiv:2210.00030*, 2022.
- 662
- 663 Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman,
664 Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding
665 large language models. *arXiv preprint arXiv:2310.12931*, 2023.
- 666
- 667 Yecheng Jason Ma, William Liang, Hung-Ju Wang, Sam Wang, Yuke Zhu, Linxi Fan, Osbert Bastani,
668 and Dinesh Jayaraman. Dreureka: Language model guided sim-to-real transfer. *arXiv preprint*
669 *arXiv:2406.01967*, 2024.
- 670
- 671 Ajay Mandlekar, Soroush Nasiriany, Bowen Wen, Iretiayo Akinola, Yashraj Narang, Linxi Fan, Yuke
672 Zhu, and Dieter Fox. Mimicgen: A data generation system for scalable robot learning using human
673 demonstrations. *arXiv preprint arXiv:2310.17596*, 2023.
- 674
- 675 OpenAI. Hello gpt-4o, May 2024. URL <https://openai.com/index/hello-gpt-4o>.
- 676
- 677 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
678 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
679 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
680 27744, 2022.
- 681
- 682 Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander
683 Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic
684 learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.
- 685
- 686 Wilbert Pumacay, Ishika Singh, Jiafei Duan, Ranjay Krishna, Jesse Thomason, and Dieter Fox. The
687 colosseum: A benchmark for evaluating generalization for robotic manipulation. *arXiv preprint*
688 *arXiv:2402.08191*, 2024.
- 689
- 690 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
691 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini
692 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint*
693 *arXiv:2403.05530*, 2024.
- 694
- 695 Amanpreet Singh, Vivek Natarjan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus
696 Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer*
697 *Vision and Pattern Recognition*, pp. 8317–8326, 2019.
- 698
- 699 Marta Skreta, Zihan Zhou, Jia Lin Yuan, Kouros Darvish, Alán Aspuru-Guzik, and Animesh
700 Garg. Replan: Robotic replanning with perception and language models. *arXiv preprint*
701 *arXiv:2401.04157*, 2024.
- 702
- 703 Stone Tao, Fanbo Xiang, Arth Shukla, Yuzhe Qin, Xander Hinrichsen, Xiaodi Yuan, Chen Bao,
704 Xinsong Lin, Yulin Liu, Tse kai Chan, Yuan Gao, Xuanlin Li, Tongzhou Mu, Nan Xiao, Arnab
705 Gurha, Zhiao Huang, Roberto Calandra, Rui Chen, Shan Luo, and Hao Su. Maniskill3: Gpu
706 parallelized robotics simulation and rendering for generalizable embodied ai, 2024. URL <https://arxiv.org/abs/2410.00425>.

- 702 Julen Urain, Ajay Mandlekar, Yilun Du, Mahi Shafiullah, Danfei Xu, Katerina Fragkiadaki, Georgia
703 Chalvatzaki, and Jan Peters. Deep generative models in robotics: A survey on learning from
704 multimodal demonstrations. *arXiv preprint arXiv:2408.04380*, 2024.
705
- 706 Lirui Wang, Yiyang Ling, Zhecheng Yuan, Mohit Shridhar, Chen Bao, Yuzhe Qin, Bailin Wang,
707 Huazhe Xu, and Xiaolong Wang. Gensim: Generating robotic simulation tasks via large language
708 models. *arXiv preprint arXiv:2310.01361*, 2023a.
- 709 Yi Ru Wang, Jiafei Duan, Dieter Fox, and Siddhartha Srinivasa. Newton: Are large language models
710 capable of physical reasoning? *arXiv preprint arXiv:2310.07018*, 2023b.
711
- 712 Annie Xie, Lisa Lee, Ted Xiao, and Chelsea Finn. Decomposing the generalization gap in imitation
713 learning for visual robotic manipulation. In *2024 IEEE International Conference on Robotics and
714 Automation (ICRA)*, pp. 3153–3160. IEEE, 2024.
- 715 Sean Ye, Glen Neville, Mariah Schrum, Matthew Gombolay, Sonia Chernova, and Ayanna Howard.
716 Human trust after robot mistakes: Study of the effects of different forms of robot communication.
717 In *2019 28th IEEE International Conference on Robot and Human Interactive Communication
718 (RO-MAN)*, pp. 1–7. IEEE, 2019.
- 719 H Peyton Young. Learning by trial and error. *Games and economic behavior*, 65(2):626–643, 2009.
720
- 721 Samson Yu, Kelvin Lin, Anxing Xiao, Jiafei Duan, and Harold Soh. Octopi: Object property
722 reasoning with large tactile-language models. *arXiv preprint arXiv:2405.02794*, 2024.
723
- 724 Wentao Yuan, Jiafei Duan, Valts Blukis, Wilbert Pumacay, Ranjay Krishna, Adithyavairavan Murali,
725 Arsalan Mousavian, and Dieter Fox. Robopoint: A vision-language model for spatial affordance
726 prediction for robotics. *arXiv preprint arXiv:2406.10721*, 2024.
- 727 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
728 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
729 pp. 3836–3847, 2023.
- 730 Zhi Zheng, Qian Feng, Hang Li, Alois Knoll, and Jianxiang Feng. Evaluating uncertainty-based
731 failure detection for closed-loop llm planners, 2024. URL [https://arxiv.org/abs/2406.
732 00430](https://arxiv.org/abs/2406.00430).
- 733
734 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
735 Tianyue Ou, Yonatan Bisk, Daniel Fried, et al. Webarena: A realistic web environment for building
736 autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023.
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755