

# RoboPlayground: Democratizing Robotic Evaluation through Structured Physical Domains

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Evaluation of robotic manipulation systems has largely relied*  
002 *on fixed benchmarks authored by a small number of experts,*  
003 *where task instances, constraints, and success criteria are*  
004 *predefined and difficult to extend. This paradigm limits who*  
005 *can shape evaluation and obscures how policies respond*  
006 *to user-authored variations in task intent, constraints, and*  
007 *notions of success. We argue that evaluating modern manip-*  
008 *ulation policies requires reframing evaluation as a language-*  
009 *driven process over structured physical domains. We present*  
010 *RoboPlayground, a framework that enables users to au-*  
011 *thor executable manipulation tasks using natural language*  
012 *within a structured physical domain. Natural language in-*  
013 *structions are compiled into reproducible task specifications*  
014 *with explicit asset definitions, initialization distributions,*  
015 *and success predicates. Each instruction defines a struc-*  
016 *tured family of related tasks, enabling controlled semantic*  
017 *and behavioral variation while preserving executability and*  
018 *comparability. We instantiate RoboPlayground in a struc-*  
019 *tured block manipulation domain and evaluate it along three*  
020 *axes. A user study shows that the language-driven interface*  
021 *is easier to use and imposes lower cognitive workload than*  
022 *programming-based and code-assist baselines. Evaluating*  
023 *learned policies on language-defined task families reveals*  
024 *generalization failures that are not apparent under fixed*  
025 *benchmark evaluations. Finally, we show that task diversity*  
026 *scales with contributor diversity rather than task count alone,*  
027 *enabling evaluation spaces to grow continuously through*  
028 *crowd-authored contributions.*

## 029 1. Introduction

030 Who gets to decide what it means for a robot to be com-  
031 petent? Today, robotic manipulation systems are evaluated  
032 almost exclusively through benchmarks designed by a small  
033 number of experts. These benchmarks specify fixed task in-  
034 stances, success conditions, and evaluation protocols, implic-  
035 itly encoding which behaviors matter and which variations

are worth testing. While this paradigm has driven substantial  
progress, it centralizes control over evaluation and limits  
both who can define evaluation tasks and what questions can  
be asked about a system’s behavior.

Similar limitations have appeared in other domains. In vi-  
sual reasoning, diagnostic datasets such as CLEVR reframed  
evaluation around controlled task generation within a struc-  
tured domain [12]. By focusing on interpretable primitives  
rather than maximal realism, CLEVR shifted evaluation  
from static instances to structured families of tasks, enabling  
clearer attribution of failure modes. In contrast, work in  
natural language processing has emphasized the dynamics  
of evaluation over time. Dynamic benchmarking efforts  
such as Dynabench treat evaluation as a human-in-the-loop  
process, where users iteratively generate and refine examples  
to surface model weaknesses [14]. Together, these efforts  
highlight two complementary principles for informative  
evaluation: structured task spaces that make variation  
interpretable, and participatory mechanisms that allow  
evaluation to grow continuously beyond its initial design.

Evaluating modern manipulation policies requires  
embracing both principles. An effective evaluation system  
should satisfy four key desiderata. First, it must be **acces-**  
**sible**, allowing users to express task intent, constraints, and  
success criteria using natural language without expertise in  
simulation internals or benchmark-specific code. Second, it  
should support **continuous growth**, enabling the evaluation  
space to expand over time through contributions from many  
users rather than remaining fixed. Third, it must ensure  
**reproducibility**, so that tasks can be precisely re-executed  
across models and evaluations, enabling fair comparison as  
the evaluation space grows. Finally, it should provide **struc-**  
**tured control**, constraining user-authored instructions to  
remain executable and interpretable while enabling systemat-  
ic variation and meaningful attribution of failure modes.

**We propose reframing robotic manipulation evalua-**  
**tion as a language-driven, user-authored process over**  
**structured physical domains, shifting evaluation from**  
**static expert-defined benchmarks to an accessible, repro-**  
**ducible, and continuously expanding task space.**

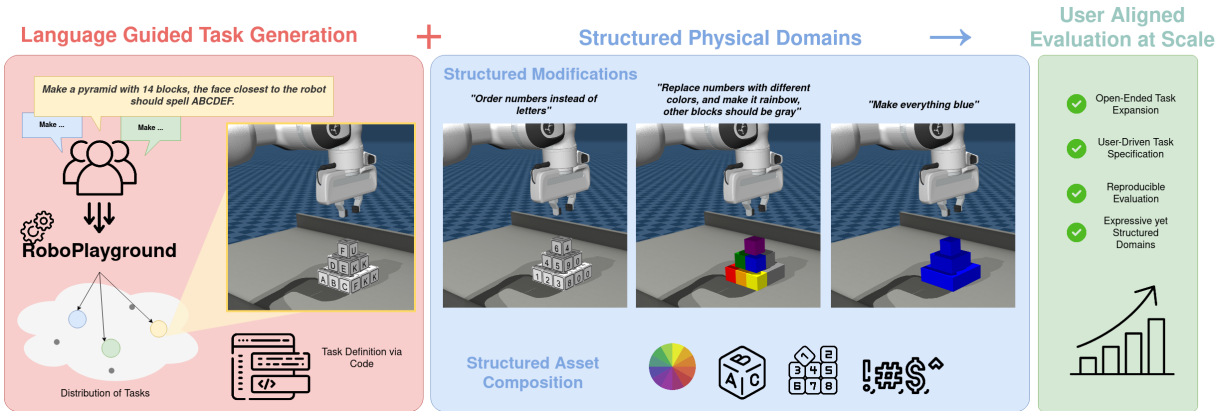


Figure 1. **Language-Guided Task Generation in Structured Physical Domains.** Natural language instructions are compiled into executable task templates within ROBOPLAYGROUND, enabling open-ended task generation while preserving structure and control. Structured domains support systematic steering through controlled variations in task semantics, constraints, and asset composition, enabling expressive, reproducible, and extensible evaluation.

076 In this work, we present RoboPlayground, a language-  
 077 driven framework for defining robotic manipulation tasks in  
 078 a structured physical domain. Natural language serves as the  
 079 primary authoring interface, allowing users to specify task in-  
 080 tent, constraints, and success conditions without interacting  
 081 with benchmark-specific code. Each instruction is compiled  
 082 into an executable task specification with explicit definitions  
 083 of assets, initialization distributions, and success predicates,  
 084 enabling tasks to be precisely re-executed across models and  
 085 evaluations. Rather than yielding a single fixed task instance,  
 086 each instruction defines a family of related tasks, whose  
 087 systematic variations are authored and controlled by users  
 088 through language within the domain’s physical structure,  
 089 allowing the evaluation space to grow continuously over  
 090 time while preserving controlled, comparable structure.

091 In summary, this paper makes three contributions. (1)  
 092 We introduce a language-driven evaluation framework that  
 093 democratizes manipulation evaluation by making task  
 094 specification accessible to non-expert users while maintain-  
 095 ing reproducibility. (2) We empirically demonstrate that  
 096 language-defined task families reveal policy behaviors and  
 097 limitations that are missed by conventional instance-based  
 098 benchmarks. (3) We show that evaluation spaces can grow  
 099 continuously through user- and model-authored instructions  
 100 without sacrificing comparability or scientific rigor.

## 101 2. Related Works

### 102 2.1. Evaluation and Benchmarking for Robotic Ma- 103 nipulation

104 Robotic manipulation systems are most commonly evaluated  
 105 using fixed benchmark suites composed of predefined  
 106 task instances, environments, and success criteria, such as  
 107 RL Bench [11], LIBERO [20], RoboCasa [23], Behaviour-1k  
 108 [16], ManiSkill [22], Colosseum [26], Simpler [17], and  
 109 RoboEval [30]. While these benchmarks have been instru-

mental in standardizing evaluation, they define evaluation  
 over a fixed and finite set of expert-authored tasks, with task  
 structure, constraints, and success criteria encoded procedurally  
 and not exposed for user modification. Although some bench-  
 marks include natural language annotations or language-conditioned  
 tasks, language is typically treated as documentation or policy  
 input rather than as part of the executable task specification,  
 making it difficult to introduce task variations or alternative  
 notions of success without modifying benchmark code. Recent  
 work has explored complementary directions for scaling evalua-  
 tion: RoboArena [3] democratizes who evaluates and where  
 evaluation occurs through crowd-sourced, double-blind pairwise  
 comparisons over unconstrained real-world tasks, while Polaris  
 [10] improves the fidelity and scalability of simulation-based  
 evaluation via real-to-sim scene reconstruction, but retains  
 fixed, expert-authored tasks and success criteria. In contrast,  
 our work focuses on democratizing what is evaluated by en-  
 abling users to author, modify, and refine executable task  
 specifications through language, while preserving structure,  
 reproducibility, and comparability.

### 104 2.2. LLM-Based Task and Environment Generation 131

Recent work has explored using large language models to  
 generate robotic tasks, environments, rewards, or curricula at  
 scale. Systems such as GenSim [28], Gen2Sim [13], Robo-  
 Gen [29], and AnyTask [7] leverage LLMs to synthesize  
 tasks or simulation assets, while Eureka [21] and Eureka-  
 verse [19] use language models to automatically generate  
 reward functions or learning curricula. These approaches  
 primarily target data generation and training diversity, rather  
 than evaluation itself: generated tasks are treated as inputs  
 to learning pipelines, and the task space is not exposed to  
 users as a controllable or interpretable evaluation interface.  
 Task generation is typically decoupled from mechanisms for  
 enforcing reproducibility, tracking task lineage, or sys-

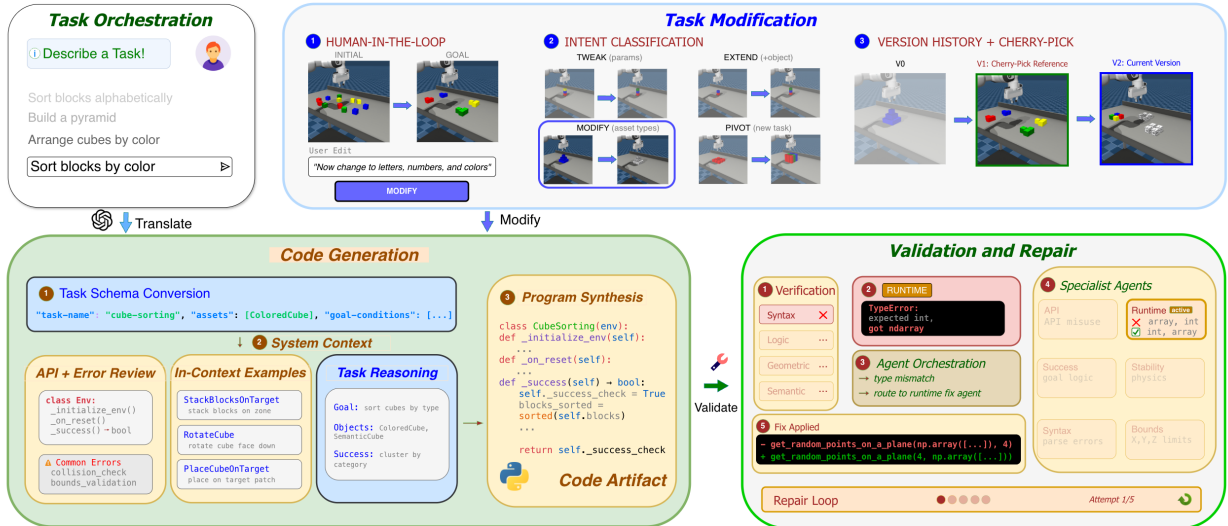


Figure 2. **System overview.** (Left) A user provides a natural language task description, which is compiled into a structured schema specifying assets, initialization logic, and success conditions. The system synthesizes an executable implementation conditioned on this schema and retrieved context (APIs, prior tasks, error patterns). (Bottom right) Tasks are admitted only after passing a multi-stage validation and repair pipeline enforcing software correctness, API consistency, and physical realizability. (Top right) Validated artifacts support context-aware steering: modification requests are classified into intent categories (*tweak*, *extend*, *modify*, *pivot*), with all variants tracked via explicit version history to preserve lineage and comparability.

145 thematically relating task variations to evaluation outcomes.  
 146 Language has also been used to guide robot behavior at ex-  
 147 ecution time, as in SayCan [1] and Code as Policies [18],  
 148 where it serves as a high-level planning or control signal  
 149 while task definitions and success criteria remain fixed and  
 150 externally specified. In contrast, our work treats language  
 151 as a first-class interface for evaluation: natural language  
 152 instructions are compiled into structured, executable task  
 153 specifications with explicit asset definitions, initialization  
 154 distributions, and success predicates, enabling users to author  
 155 reproducible families of semantically related evaluation tasks  
 156 that support controlled variation and systematic comparison.

### 157 3. Methods

158 This section describes how the framework operationalizes  
 159 the core desiderata: accessibility, continuous growth,  
 160 reproducibility, and structured control. We first define the  
 161 design principles and formalize the task representation in  
 162 Section 3.1. We then describe the task orchestration process  
 163 that make language-defined manipulation tasks executable  
 164 and interpretable in Section 3.2. We then describe how  
 165 natural language instructions are compiled into concrete,  
 166 reproducible task artifacts through a validation pipeline that  
 167 enforces physical realizability and consistency Section 3.3.  
 168 Finally, we introduce a context-aware steering mechanism  
 169 that enables users to systematically vary tasks and expand  
 170 the evaluation space over time while preserving explicit  
 171 comparability between task variants in Section 3.4.

### 3.1. Task Representation and Design Principles 172

173 Our framework treats natural language as an executable in-  
 174 terface for task specification. User instructions are compiled  
 175 into concrete task realizations that fully determine assets, ini-  
 176 tialization logic, and success conditions. Rather than produc-  
 177 ing isolated benchmark instances, the system yields reusable  
 178 task artifacts that can be shared, re-executed, and systemati-  
 179 cally varied, enabling evaluation spaces to grow through user  
 180 contribution without sacrificing scientific rigor. Figure 2 pro-  
 181 vides an overview of this process. Natural language instruc-  
 182 tions are compiled into structured task proposals, synthesized  
 183 into executable implementations, and admitted as task arti-  
 184 facts only after passing a multi-stage validation pipeline.  
 185 Validated artifacts can then be iteratively refined through  
 186 context-aware steering, enabling controlled task variation  
 187 while preserving reproducibility and explicit lineage.

188 Formally, we define a manipulation task as a tuple

$$189 T = (\mathcal{A}, \rho_0, G, \ell, \mathcal{V}), \quad (1)$$

190 where  $\mathcal{A}$  denotes the set of task assets,  $\rho_0$  is a distribution  
 191 over initial states,  $G : \mathcal{S} \rightarrow \{0, 1\}$  is a success predicate  
 192 over simulator states,  $\ell$  is the canonical natural language  
 193 instruction, and  $\mathcal{V}$  is a set of paraphrases used for robustness  
 194 testing.

195 This decomposition reflects a deliberate design choice.  
 196 Logical equivalence at the level of language does not  
 197 imply equivalence of task realization: differences in  
 198 tolerances, reset distribution, or success-check timing can  
 199 lead to divergent evaluation outcomes even when tasks  
 200 are described identically. Language alone is therefore  
 201 insufficient as a unit of evaluation.

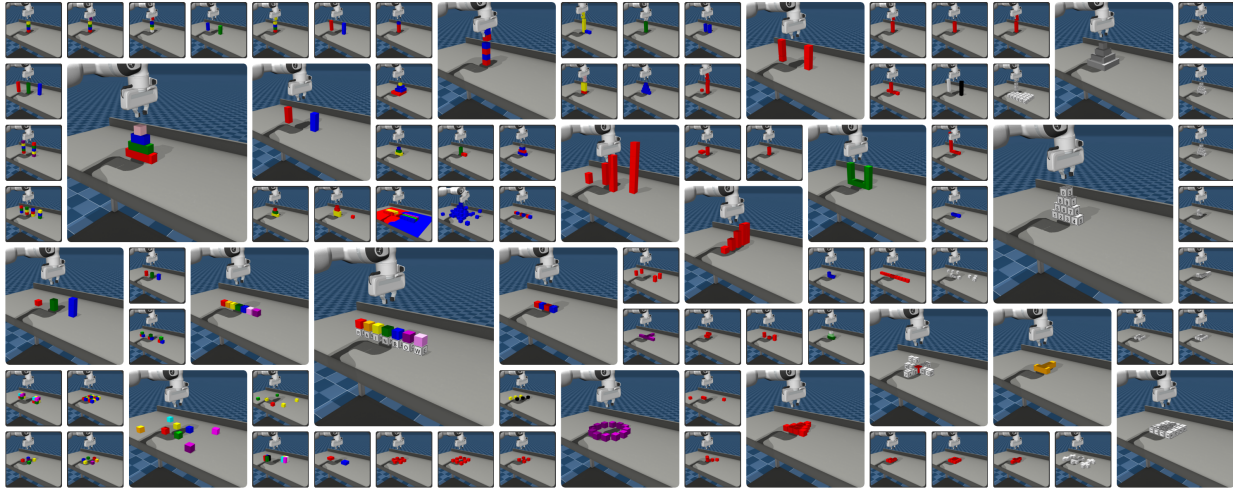


Figure 3. **Sample of Language-Defined Manipulation Tasks.** Executable tasks spanning geometric constructions, spatial alignment, and semantically constrained arrangements, organized by proximity in a learned task-embedding space. Coherent clusters correspond to families with shared attributes, spatial relations, and success definitions (e.g., stacking, ordering, grouping, alignment).

### 202 3.2. Task Orchestration Through Language

203 Given a natural language description  $u$ , task construction be-  
 204 gins by translating language into a structured representation  
 205 of task intent. Specifically, the system infers and populates a  
 206 fixed `TaskSchema` that explicitly specifies the task name,  
 207 relevant assets, goal conditions, and initialization logic. The  
 208 use of a fixed schema ensures that all task-relevant fields  
 209 are present and disambiguated before execution, enabling  
 210 complete interpretation of the instruction and preventing  
 211 underspecified task definitions.

212 Conditioned on the validated task schema, the system  
 213 then synthesizes an executable task implementation. This  
 214 process leverages an LLM with access to relevant envi-  
 215 ronment APIs, prior task implementations, and diagnostic  
 216 error information retrieved based on structural similarity  
 217 to the proposed task. The LLM produces an intermediate  
 218 natural-language task specification that articulates the  
 219 intended objects, goal configuration, and success criteria,  
 220 which is subsequently compiled into executable code.

221 Executable tasks are implemented as classes that extend  
 222 a fixed environment interface. Each task defines methods  
 223 for environment initialization, reset-time sampling from  
 224  $\rho_0$ , and success evaluation corresponding to  $G$ . This  
 225 constrained interface enforces uniform structure across task  
 226 implementations and limits variation arising from authoring  
 227 style, ensuring that differences in evaluation outcomes  
 228 reflect task content rather than implementation artifacts.

### 229 3.3. Validation and Physical Realizability

230 Language-defined tasks are only meaningful if they are  
 231 both executable and physically realizable. Each synthesized  
 232 task implementation is therefore subjected to a multi-stage  
 233 validation pipeline before being admitted as a task artifact.

234 **Basic validation.** Basic validation enforces software

correctness independent of physics simulation. Generated  
 code is subjected to static analysis to detect syntactic errors  
 and forbidden patterns, compiled in an isolated execution  
 environment to detect import and definition errors, and in-  
 stantiated to detect runtime failures during object creation or  
 reset-time sampling.

**Goal-state verification.** To enforce physical realizability  
 of the success predicate, tasks are instantiated directly in  
 the goal configuration and simulated forward under zero  
 action to allow contacts to settle. The success predicate  
 $G$  must evaluate to true after settling and remain true over  
 an extended horizon, ensuring that the goal configuration  
 is both achievable and stable under the simulator’s physics  
 model.

**Iterative repair.** When validation fails, the failure  
 is classified according to its source (e.g., syntax, API  
 usage, runtime instantiation, goal satisfaction, or physical  
 instability), and a corresponding repair operator proposes  
 a localized modification to the task implementation. Repairs  
 may adjust object placements, relax geometric constraints,  
 or rewrite components of the success predicate, depending  
 on the failure type. Validation is then re-run on the repaired  
 implementation. This process repeats until all checks pass  
 or a fixed retry budget is exhausted, ensuring that admitted  
 tasks satisfy executability and physical consistency while  
 remaining faithful to the original language intent.

### 3.4. Controlled Task Modifications

A validated task artifact defines a reference task instance  
 from which a family of related tasks can be derived.  
 To enable systematic task variation without sacrificing  
 comparability, the framework provides a context-aware  
 steering mechanism that interprets user modification  
 requests and constrains how tasks may evolve.

Table 1. User study results ( $N=26$ ). Mean  $\pm$  95% CI margin. SUS: 0–100 scale. TLX: unweighted mean of five NASA-TLX subscales, normalized to 0–100. Rank: 1=best. Preferred: share selecting each system overall.

System	SUS $\uparrow$	TLX workload $\downarrow$	Usability rank $\downarrow$	Preferred (%)
GenSim	52.5 $\pm$ 9.3	41.8 $\pm$ 9.0	2.7 $\pm$ 0.2	8
Cursor	68.8 $\pm$ 7.8	36.7 $\pm$ 10.4	2.0 $\pm$ 0.2	23
RoboPlayground	<b>83.4<math>\pm</math>6.9</b>	<b>18.6<math>\pm</math>7.7</b>	<b>1.3<math>\pm</math>0.3</b>	<b>69</b>

Given a modification request, the system first interprets the intent and extracts structured parameters such as dimensional changes, ordering constraints, or asset-type substitutions. It then classifies the request into one of five steering categories: *Tweak*, *Extend*, *Modify*, *Pivot*, or *Fresh*. Each category specifies explicit preservation guarantees over the task components ( $\mathcal{A}$ ,  $\rho_0$ ,  $G$ ). For example, *Tweak* and *Extend* preserve the original task structure and success predicate, enabling direct comparability with the reference task, while *Modify* and *Pivot* permit progressively broader semantic or structural changes when required by the user intent.

Task evolution is tracked through versioned snapshots that record structured summaries of assets, goals, and code hashes. When a modification requires asset types incompatible with the current version, the system selects a compatible prior snapshot as the reference. This allows coherent multi-step refinement without manual bookkeeping. Each validated variant produces a new snapshot, yielding version-controlled task families with explicit lineage suitable for systematic evaluation and controlled analysis of task variation.

## 4. Results

We evaluate ROBOPLAYGROUND along three axes that are central to its role as a democratized evaluation framework for robotic manipulation: (i) the usability of its task authoring interface, (ii) the diagnostic value of the resulting task set for assessing policy generalization, and (iii) the scalability of task creation under open-world, crowd-driven use. Across all experiments, we focus on whether ROBOPLAYGROUND enables task specifications that are both easier to author and more informative for evaluation than existing alternatives.

### 4.1. Usability of the Task Authoring Interface

**Experimental setting.** We evaluate the usability of RoboPlayground in comparison to two baseline task authoring interfaces, GenSim [28] and Cursor [2], using a within-subjects user study ( $N = 26$ ). Participants were asked to construct an identical manipulation task (build a 3D structure using blocks under various constraints) using each system. All participants interacted with all three systems, enabling paired comparisons of perceived usability, cognitive workload, and user preference. We measure usability using the System Usability Scale (SUS) [5], cognitive workload using NASA-TLX subscales [8], and overall preference through usability and forced-choice

rankings. Results are summarized in Table 1.

**RoboPlayground achieves higher perceived usability than baselines.** Across participants ( $N=26$ ), RoboPlayground attains the highest System Usability Scale (SUS) score ( $83.4 \pm 6.9$ ; mean  $\pm$  95% confidence interval margin), well above the conventional acceptability threshold of 68. GenSim and Cursor achieve substantially lower mean SUS ( $52.5 \pm 9.3$  and  $68.8 \pm 7.8$ , respectively). Paired Wilcoxon signed-rank tests confirm that RoboPlayground significantly outperforms both GenSim ( $p < 0.001$ ) and Cursor ( $p = 0.0017$ ), so the advantage is not limited to the weakest baseline: RoboPlayground is rated more usable than a strong general-purpose assistant interface as well. The interval for GenSim is the widest of the three, consistent with more heterogeneous experiences in that condition, whereas RoboPlayground shows the tightest margin among systems, indicating comparatively consistent high ratings.

**RoboPlayground reduces perceived cognitive workload relative to baselines.** Cognitive workload is summarized as the unweighted mean of five NASA-TLX subscales (Mental Demand, Temporal Demand, Effort, Frustration, and reversed Performance), each normalized to 0–100 and oriented so that lower is better. RoboPlayground yields the lowest mean composite score ( $18.6 \pm 7.7$ ; mean  $\pm$  95% confidence interval margin), compared to  $41.8 \pm 9.0$  for GenSim and  $36.7 \pm 10.4$  for Cursor. Paired Wilcoxon signed-rank tests show that RoboPlayground significantly reduces perceived workload relative to both GenSim ( $p = 0.0007$ ) and Cursor ( $p = 0.0019$ ). GenSim and Cursor do not differ significantly from each other on this composite ( $p = 0.22$ ), whereas RoboPlayground separates clearly from each baseline; Cursor also exhibits the widest TLX margin among the three, indicating somewhat more spread in workload ratings even though the paired comparison to RoboPlayground remains significant.

**Participants consistently prefer RoboPlayground over baseline interfaces.** Subjective measures reinforce the quantitative usability and workload results. RoboPlayground achieves the best mean usability rank ( $1.3 \pm 0.3$ ; lower is better), with GenSim and Cursor at  $2.7 \pm 0.2$  and  $2.0 \pm 0.2$ , respectively. A Friedman test shows strong differences in rankings across systems ( $p < 0.001$ ), and post-hoc paired Wilcoxon tests confirm that RoboPlayground is ranked significantly better than both GenSim ( $p = 0.0001$ ) and Cursor ( $p = 0.0078$ ). In forced-choice overall preference, 69% of participants select RoboPlayground, compared to 23% for Cursor and 8% for GenSim. A chi-square goodness-of-fit test rejects a uniform split across the three options ( $p = 0.0003$ ), consistent with concentration of preference on RoboPlayground. Together, the ranking and preference distributions indicate a stable, statistically supported tilt toward RoboPlayground over both baselines.

Table 2. **Results on in-distribution and generalization tasks.** Success rates (%)  $\pm$  standard errors across six policies on training tasks (top) and held-out generalization tasks (bottom); best per task in **bold**. Generalization tasks perturb training tasks along semantic (**S**), visual (**V**), and behavioural (**B**) axes per [6].

Method	Red Block Front of Yellow	Red Block Right Placement	Red Behind Yellow	Red Block Stacking	Three Block Color Stacking	Red Block Left Placement	Red on Yellow Stack	Yellow on Red Stack	Color Block Alignment	Place Two Blocks on Patch
Pi-0.5	72.0 $\pm$ 6.3	68.0 $\pm$ 6.6	58.0 $\pm$ 7.0	0.0 $\pm$ 0.0	6.0 $\pm$ 3.4	64.0 $\pm$ 6.8	46.0 $\pm$ 7.0	62.0 $\pm$ 6.9	0.0 $\pm$ 0.0	74.0 $\pm$ 6.2
Pi-0.5 (LoRA)	32.0 $\pm$ 6.6	16.0 $\pm$ 5.2	48.0 $\pm$ 7.1	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	12.0 $\pm$ 4.6	12.0 $\pm$ 4.6	36.0 $\pm$ 6.8	0.0 $\pm$ 0.0	28.0 $\pm$ 6.3
Adapter	66.0 $\pm$ 6.7	64.0 $\pm$ 6.8	60.0 $\pm$ 6.9	0.0 $\pm$ 0.0	2.0 $\pm$ 2.0	26.0 $\pm$ 6.2	56.0 $\pm$ 7.0	54.0 $\pm$ 7.0	2.0 $\pm$ 2.0	78.0 $\pm$ 5.9
Dual	<b>88.0 <math>\pm</math> 4.6</b>	76.0 $\pm$ 6.0	<b>74.0 <math>\pm</math> 6.2</b>	2.0 $\pm$ 2.0	12.0 $\pm$ 4.6	<b>84.0 <math>\pm</math> 5.2</b>	64.0 $\pm$ 6.8	54.0 $\pm$ 7.0	6.0 $\pm$ 3.4	84.0 $\pm$ 5.2
GR00T	82.0 $\pm$ 5.4	84.0 $\pm$ 5.2	<b>74.0 <math>\pm</math> 6.2</b>	<b>10.0 <math>\pm</math> 4.2</b>	<b>22.0 <math>\pm</math> 5.9</b>	68.0 $\pm$ 6.6	66.0 $\pm$ 6.7	56.0 $\pm$ 7.0	<b>12.0 <math>\pm</math> 4.6</b>	<b>96.0 <math>\pm</math> 2.8</b>
Qwen-OFT	82.0 $\pm$ 5.4	<b>86.0 <math>\pm</math> 4.9</b>	68.0 $\pm$ 6.6	2.0 $\pm$ 2.0	14.0 $\pm$ 4.9	<b>84.0 <math>\pm</math> 5.2</b>	<b>76.0 <math>\pm</math> 6.0</b>	<b>68.0 <math>\pm</math> 6.6</b>	2.0 $\pm$ 2.0	78.0 $\pm$ 5.9

Perturbation	S+B	V	S	S+V	V+B	S	V	S+B	V	S+B	V	
Method	Red Block Two Tower	Blue Block Stacking	Yellow Block Left Placement	Red Block Left of Blue	Yellow on Red Unstack Restack	Green on Blue Stack	Three Block Perturbed	Three Block Beside	Place Two Blue Blocks on Patch	Place Two Blocks on Green Patch	Stack Two Blocks on Patch	Place Two Blocks on Long Patch
Pi-0.5	4.0 $\pm$ 2.8	0.0 $\pm$ 0.0	76.0 $\pm$ 6.0	50.0 $\pm$ 7.1	0.0 $\pm$ 0.0	60.0 $\pm$ 6.9	10.0 $\pm$ 4.2	12.0 $\pm$ 4.6	46.0 $\pm$ 7.0	80.0 $\pm$ 5.7	0.0 $\pm$ 0.0	<b>68.0 <math>\pm</math> 6.6</b>
Pi-0.5 (LoRA)	2.0 $\pm$ 2.0	0.0 $\pm$ 0.0	10.0 $\pm$ 4.2	32.0 $\pm$ 6.6	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	2.0 $\pm$ 2.0	0.0 $\pm$ 0.0	10.0 $\pm$ 4.2	30.0 $\pm$ 6.5	<b>2.0 <math>\pm</math> 2.0</b>	28.0 $\pm$ 6.3
Adapter	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	4.0 $\pm$ 2.8	36.0 $\pm$ 6.8	0.0 $\pm$ 0.0	0.0 $\pm$ 0.0	8.0 $\pm$ 3.8	0.0 $\pm$ 0.0	46.0 $\pm$ 7.0	38.0 $\pm$ 6.9	0.0 $\pm$ 0.0	26.0 $\pm$ 6.2
Dual	8.0 $\pm$ 3.8	0.0 $\pm$ 0.0	74.0 $\pm$ 6.2	<b>60.0 <math>\pm</math> 6.9</b>	0.0 $\pm$ 0.0	56.0 $\pm$ 7.0	8.0 $\pm$ 3.8	10.0 $\pm$ 4.2	80.0 $\pm$ 5.7	72.0 $\pm$ 6.3	<b>2.0 <math>\pm</math> 2.0</b>	54.0 $\pm$ 7.0
GR00T	<b>10.0 <math>\pm</math> 4.2</b>	0.0 $\pm$ 0.0	<b>78.0 <math>\pm</math> 5.9</b>	52.0 $\pm$ 7.1	<b>2.0 <math>\pm</math> 2.0</b>	62.0 $\pm$ 6.9	<b>20.0 <math>\pm</math> 5.7</b>	4.0 $\pm$ 2.8	<b>86.0 <math>\pm</math> 4.9</b>	<b>90.0 <math>\pm</math> 4.2</b>	0.0 $\pm$ 0.0	<b>68.0 <math>\pm</math> 6.6</b>
Qwen-OFT	6.0 $\pm$ 3.4	0.0 $\pm$ 0.0	54.0 $\pm$ 7.0	48.0 $\pm$ 7.1	0.0 $\pm$ 0.0	<b>66.0 <math>\pm</math> 6.7</b>	14.0 $\pm$ 4.9	<b>14.0 <math>\pm</math> 4.9</b>	54.0 $\pm$ 7.0	66.0 $\pm$ 6.7	0.0 $\pm$ 0.0	56.0 $\pm$ 7.0

363  
364

## 4.2. Evaluating Policies on Training and Generated Generalization Tasks

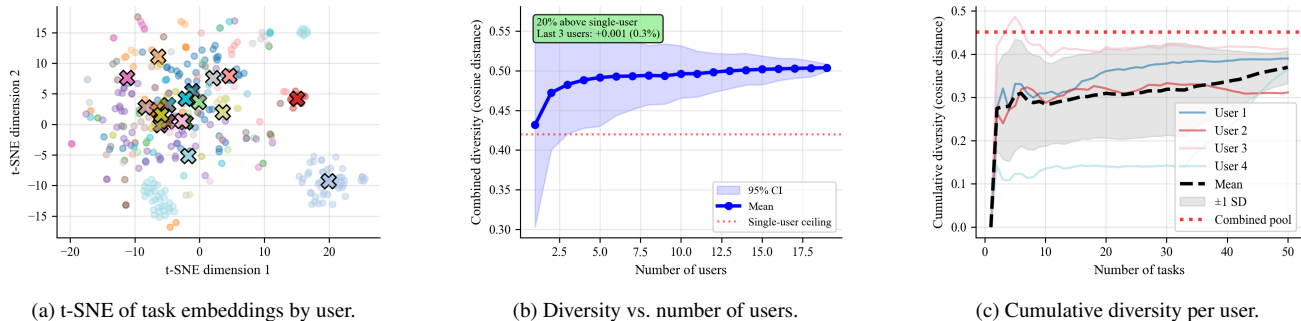
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402

**Tasks.** All policies are trained on a shared set of base manipulation tasks covering spatial relations, stacking, alignment, semantic disambiguation, and targeted placement. For each training task, we generate successful demonstration trajectories using CuTAMP [27] and hold the resulting dataset fixed across policies, ensuring that performance differences reflect policy behavior rather than differences in supervision. Evaluation is performed on (i) tasks drawn from the training distribution, and (ii) user generated generalization tasks that require adaptation beyond the training distribution. Generalization tasks are constructed by applying controlled modifications to base tasks using RoboPlayground, including semantic changes to language instructions, visual changes to object attributes and initial configurations (e.g. partially stacked versus scattered on table), and behavioral changes that alter the required action sequences or temporal structure. These transformations follow the task taxonomy described in [6], and are designed to isolate specific dimensions of generalization while preserving task validity. **Models.** We evaluate six policies spanning different action head architectures and fine-tuning strategies. Four are built on a shared Qwen3-VL-4B-Instruct [4] vision-language backbone and differ in their action decoding mechanism: Adapter appends learnable action query tokens to the VLM sequence and decodes actions via an MLP-ResNet regression head with an L1 objective; GR00T conditions a flow-matching diffusion transformer (DiT) on the VLM’s final hidden states, adopting a dual-system architecture inspired by GR00T N1.5 [24]; Dual extends this flow-matching action head with a secondary DINOv2 [25] visual encoder whose patch features are concatenated with the VLM hidden states before conditioning; and Qwen-OFT regresses actions from special action-token positions via an MLP head, following the OpenVLA-OFT design [15]. As external baselines, we include Pi-0.5 [9], a proprietary VLA with flow-matching action generation, evaluated both with full finetuning and with low-rank adaptation (Pi-0.5 (LoRA)). All models are trained end-to-end on identical demonstration data. Details of training configu-

rations and generalization tasks are outlined in the appendix. 403

**Training-distribution performance.** We first report performance on tasks drawn from the training distribution (top section of Table 2) to provide context for subsequent generalization results. All models achieve moderate to high success on tasks involving simple spatial relations and single-object placement, with GR00T, Dual, and Qwen-OFT consistently outperforming Pi-0.5 and Adapter on most placement and relational tasks. GR00T achieves the highest overall in-distribution performance, reaching 96% on Place Two Blocks on Patch and leading on stacking-related tasks. However, all models exhibit a consistent difficulty gradient: performance degrades sharply on training tasks requiring greater compositional structure or longer-horizon execution, such as multi-block stacking and color block alignment, where even the strongest model does not exceed 22%. Adapter shows notably uneven in-distribution performance, achieving competitive results on some tasks (e.g., 78% on patch placement) while lagging substantially on others (e.g., 26% on left placement). The Pi-0.5 (LoRA) variant underperforms all other models across nearly every training task, suggesting that low-rank adaptation alone is insufficient to retain the base model’s capabilities in this setting. 404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426

**Generalization results reveal a clear asymmetry across perturbation types.** Across held-out evaluation tasks, all models generalize unevenly across perturbation axes, though the degree of degradation varies by architecture. Performance remains relatively strong under visual perturbations that alter perceptual attributes while preserving execution structure: GR00T achieves 90% on Place Two Blocks on Green Patch and 86% on Place Two Blue Blocks on Patch, and Dual similarly transfers well on these tasks (72% and 80%, respectively). Semantic perturbations alone yield mixed but non-zero success for most models, with GR00T reaching 78% on Yellow Block Left Placement and 62% on Green on Blue Stack, and Dual achieving 74% and 56% on the same tasks, indicating meaningful robustness to relational re-specification. Adapter, however, largely fails under se- 427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442



(a) t-SNE of task embeddings by user. (b) Diversity vs. number of users. (c) Cumulative diversity per user.

Figure 4. **Inter-user and intra-user diversity of natural-language manipulation tasks.** (a) t-SNE of sentence embeddings colored by user (crosses = centroids). Tasks cluster by author, indicating systematic conceptual differences. (b) Mean pairwise diversity increases monotonically with number of users (shaded: 95% CI). (c) Intra-user diversity plateaus quickly; the combined pool (red dashed) achieves substantially higher diversity, highlighting complementary contributions.

443 mantic perturbation (e.g., 4% on Yellow Block Left  
444 Placement, 0% on Green on Blue Stack), sug-  
445 gesting that adapter-based finetuning may overfit to surface-  
446 level task features. In contrast, tasks involving behavioural  
447 perturbations consistently expose severe failure modes  
448 across *all* architectures. Tasks requiring multi-stage execu-  
449 tion, non-monotonic progress (e.g., unstack-restore), or com-  
450 positional sequencing yield near-zero success universally—  
451 no model exceeds 2% on Yellow on Red Unstack  
452 Restack or Stack Two Blocks on Patch, and  
453 Blue Block Stacking elicits 0% across the board.  
454 These failures occur despite reasonable performance on simpler  
455 stacking or placement tasks in isolation, suggesting  
456 limited procedural and compositional generalization rather  
457 than a lack of basic manipulation competence.  $\text{Pi-0.5}$   
458 (LoRA) further degrades performance across nearly all per-  
459 turbation axes, confirming increased sensitivity to deviations  
460 from training task structure under constrained adaptation.  
461 Together, these results establish behavioural perturbations  
462 as the dominant source of generalization failure *independent*  
463 *of model architecture* and motivate evaluation protocols that  
464 explicitly probe execution structure.

465 **Implications.** Together, these results demonstrate  
466 that success on a fixed set of training-distribution tasks  
467 can substantially overestimate a policy’s robustness. By  
468 enabling the generation of creative task variants that probe  
469 specific semantic, visual, and behavioral dimensions of  
470 generalization, ROBOPLAYGROUND enables fine-grained  
471 diagnosis of policy capabilities and failure modes that are  
472 obscured by static task definitions.

### 473 4.3. Scalability and Diversity of RoboPlayground

474 For evaluation, diversity matters not as raw task count, but as  
475 coverage of distinct task intents and constraint combinations  
476 that probe different policy behaviors. We evaluate how task  
477 diversity in ROBOPLAYGROUND scales with the number of  
478 contributors and the number of authored tasks. Each contrib-  
479 utor authors up to 50 valid manipulation tasks in the blocks  
480 domain using the same interface and asset set. To quantify

diversity, we compute average pairwise distances between  
task representations using semantic sentence embeddings,  
and analyze both pooled task sets across contributors  
and cumulative task sets authored by individuals. In the  
appendix, we report additional analyses using alternative  
diversity measures, which show consistent trends.

**Inter-user diversity scales with contributors.** Figure  
4(a) shows the distribution of tasks authored by  
individual users. A t-SNE projection reveals that some  
users occupy distinct regions of the embedding space, while  
others exhibit substantial overlap, suggesting systematic  
differences in how contributors conceptualize and describe  
manipulation goals within the domain. When tasks are  
pooled across users (10 tasks per user), the mean pairwise  
diversity increases monotonically with the number of  
contributors (Fig. 4(b)). Even after substantial saturation,  
adding the final three contributors yields a consistent, non-  
zero increase in diversity, indicating that new contributors  
continue to introduce semantically novel task formulations.

**Intra-user diversity exhibits diminishing returns.**  
In contrast, Figure 4(c) shows that when tasks are added  
incrementally from a single user, cumulative diversity  
grows rapidly at first but quickly plateaus. This behavior  
is consistent across most users, suggesting that individual  
authors often explore a limited region of the task space,  
potentially shaped by their preferred abstractions, phrasing,  
and constraint patterns. Even prolific contributors produce  
increasingly redundant task variations over time.

**Complementarity of multiple contributors.** Notably,  
the combined task pool outperforms any individual  
contributor in terms of cumulative diversity, as shown in  
Figure 4(c). This gap highlights the complementary nature  
of crowd-authored task generation: different users introduce  
distinct semantic concepts, compositional structures, and  
constraint combinations that are rarely discovered by a  
single author alone. Qualitative inspection of task clusters  
confirms the presence of novel formulations of spatial  
relations, multi-object constraints, and success conditions  
that are absent from single-author collections.

Table 3. **Ablation Study.** Cumulative addition of pipeline components, starting with all gates disabled per module. All metrics are percentages ( $n = 26$ ); green  $\uparrow$ /red  $\downarrow$  show change from prior row.

Module	Configuration	Task Succ.	Compile	Smoke Test	Human-Ver.	LLM Align.
Task Proposal	None (all disabled)	100.0	100.0	100.0	88.5	74.0
	+ asset inference	100.0	100.0	100.0	92.3 ( $\uparrow$ 3.8)	71.5 ( $\downarrow$ 2.5)
	+ feasibility checking (all gates on)	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b> ( $\uparrow$ 7.7)	<b>73.6</b> ( $\uparrow$ 2.1)
Code Generation	None (all disabled)	100.0	100.0	100.0	96.2	70.8
	+ API review	96.2 ( $\downarrow$ 3.8)	100.0	100.0	100.0 ( $\uparrow$ 3.8)	70.5 ( $\downarrow$ 0.3)
	+ common errors review	100.0 ( $\uparrow$ 3.8)	100.0	100.0	96.2 ( $\downarrow$ 3.8)	72.4 ( $\uparrow$ 1.9)
	+ in-context examples (all gates on)	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b> ( $\uparrow$ 3.8)	<b>73.6</b> ( $\uparrow$ 1.2)
Validation	None (all disabled)	12.0	96.0	12.0	96.2	71.1
	+ text validation	96.2 ( $\uparrow$ 84.2)	100.0 ( $\uparrow$ 4)	100.0 ( $\uparrow$ 88)	96.2	73.8 ( $\uparrow$ 2.7)
	+ compilation	100.0 ( $\uparrow$ 3.8)	100.0	100.0	100.0 ( $\uparrow$ 3.8)	71.5 ( $\downarrow$ 2.3)
	+ instantiation runtime	100.0	100.0	100.0	96.2 ( $\downarrow$ 3.8)	72.3 ( $\uparrow$ 0.8)
	+ success checking	96.2 ( $\downarrow$ 3.8)	100.0	96.2 ( $\downarrow$ 3.8)	96.2	73.0 ( $\uparrow$ 0.7)
	+ bounds checking	100.0 ( $\uparrow$ 3.8)	100.0	100.0 ( $\uparrow$ 3.8)	92.3 ( $\downarrow$ 3.9)	73.1 ( $\uparrow$ 0.1)
	+ specialist agents (all gates on)	<b>100.0</b>	<b>100.0</b>	<b>100.0</b>	<b>100.0</b> ( $\uparrow$ 7.7)	<b>73.6</b> ( $\uparrow$ 0.5)
	None (all disabled)	95.7	100.0	100.0	92.3	73.7
Context Steering	+ intent interpretation	96.2 ( $\uparrow$ 0.5)	100.0	96.2 ( $\downarrow$ 3.8)	100.0 ( $\uparrow$ 7.7)	73.8 ( $\uparrow$ 0.1)
	+ routing classification	96.2	100.0	96.2	76.9 ( $\downarrow$ 23.1)	72.5 ( $\downarrow$ 1.3)
	+ version history tracking	96.2	100.0	96.2	100.0 ( $\uparrow$ 23.1)	73.6 ( $\uparrow$ 1.1)
	+ reference selection (all gates on)	<b>100.0</b> ( $\uparrow$ 3.8)	<b>100.0</b>	<b>100.0</b> ( $\uparrow$ 3.8)	<b>100.0</b>	<b>73.6</b>

520 Overall, these results show that ROBOPLAYGROUND  
521 scales not merely by increasing task count, but by expanding  
522 coverage of the underlying task space through contributor  
523 diversity. By expanding coverage across task intent and  
524 constraint structure, ROBOPLAYGROUND enables evaluation  
525 to reveal brittleness to even seemingly minor semantic  
526 variations.

## 527 5. Ablative Studies

528 We conduct a cumulative ablation study to quantify the  
529 functional contribution of each component in the task  
530 generation pipeline. For each module, we begin with all  
531 components disabled and progressively enable individual  
532 gates. This design disentangles changes in semantic  
533 task specification from improvements in robustness and  
534 correctness under session-level evaluation (Table 3).

535 **Metrics.** We report complementary metrics capturing distinct  
536 failure modes. *Compile* and *Smoke Test* measure code  
537 correctness and execution stability; *Task Success* measures  
538 end-to-end satisfaction of the success predicate; *Human Ver-*  
539 *ification* evaluates perceived task validity; and *LLM Align-*  
540 *ment* measures consistency between the natural language  
541 instruction and the implemented success condition.

542 **Evaluation setting.** Ablations are evaluated on ten bench-  
543 mark testcases, each consisting of multiple related tasks eval-  
544 uated as a single session. Some testcases involve multi-stage  
545 task refinement via context-aware steering; additional details  
546 are provided in the Appendix.

547 **Task Proposal.** Task proposal components primarily affect  
548 semantic grounding rather than executability. Enabling  
549 asset inference improves Human Verification (88.5 to 92.3)  
550 but slightly reduces LLM Alignment (74.0 to 71.5), while

leaving execution metrics unchanged. Adding feasibility  
checking improves both Human Verification (92.3 to 96.2)  
and LLM Alignment (71.5 to 73.6) without affecting execu-  
tability.

551 **Code Generation.** Code generation components primar-  
552 ily improve robustness to systematic implementation errors.  
553 Across ablations, compilation and smoke test success remain  
554 near-perfect. API review, error checks, and in-context exam-  
555 ples incrementally improve LLM Alignment (70.8 to 73.6)  
556 while preserving end-to-end executability.  
557

558 **Validation.** Validation is the dominant determinant of  
559 task correctness. With validation disabled, Task Success  
560 drops to 12.0 despite high compilation rates. Text-level  
561 validation alone recovers Task Success to 96.2, while the  
562 full validation stack achieves perfect Task Success, Compile,  
563 Smoke Test, and Human Verification.  
564

565 **Context Steering.** Context steering influences semantic  
566 coherence across multi-step task sessions. Intent interpre-  
567 tation and version history tracking improve Task Success,  
568 Human Verification, and LLM Alignment, while routing  
569 without history degrades semantic consistency. With full  
570 context steering enabled, execution metrics remain perfect.  
571

## 572 6. Discussions

573 This work explores how robotic manipulation evaluation  
574 changes when task specification is opened to a broader  
575 set of contributors. By treating language as an executable  
576 interface, RoboPlayground allows users to express task  
577 intent, constraints, and success criteria directly, rather than  
578 relying on fixed, expert-authored benchmarks. In doing so, it  
579 reframes evaluation as a process shaped not only by models  
580 and metrics, but by the people defining what is being tested.  
581

582

**References**

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

Amanpreet Singh, Pratik Ringshia, et al. Dynabench: Re-  
thinking benchmarking in nlp, 2021. 1 637  
638[15] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning  
vision-language-action models: Optimizing speed and suc-  
cess, 2025. 6 639  
640  
641[16] Chengshu Li, Ruohan Zhang, Josiah Wong, Cem Gokmen,  
Sanjana Srivastava, Roberto Martín-Martín, Chen Wang,  
Gabrael Levine, Wensi Ai, Benjamin Martinez, Hang Yin,  
Michael Lingelbach, Minjune Hwang, Ayano Hiranaka, Su-  
jay Garlanka, Arman Aydin, Sharon Lee, Jiankai Sun, Mona  
Anvari, Manasi Sharma, Dhruva Bansal, Samuel Hunter,  
Kyu-Young Kim, Alan Lou, Caleb R Matthews, Ivan Villa-  
Renteria, Jerry Huayang Tang, Claire Tang, Fei Xia, Yunzhu  
Li, Silvio Savarese, Hyowon Gweon, C. Karen Liu, Jiajun Wu,  
et al. Behavior-1k: A human-centered, embodied ai bench-  
mark with 1,000 everyday activities and realistic simulation,  
2024. 2 642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653[17] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees,  
Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel  
Sieh, Sean Kirmani, et al. Evaluating real-world robot manip-  
ulation policies in simulation, 2024. 2 654  
655  
656  
657[18] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Haus-  
man, Brian Ichter, Pete Florence, and Andy Zeng. Code as  
policies: Language model programs for embodied control,  
2023. 3 658  
659  
660  
661[19] William Liang, Sam Wang, Hung-Ju Wang, Osbert Bastani,  
Dinesh Jayaraman, and Yecheng Jason Ma. EurekaVerse: En-  
vironment curriculum generation via large language models,  
2024. 2 662  
663  
664  
665[20] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu,  
Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge  
transfer for lifelong robot learning, 2023. 2 666  
667  
668[21] Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An  
Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi  
Fan, and Anima Anandkumar. Eureka: Human-level reward  
design via coding large language models, 2024. 2 669  
670  
671  
672[22] Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xu-  
anlin Li, Stone Tao, Zhiao Huang, Zhiwei Jia, and Hao Su.  
Maniskill: Generalizable manipulation skill benchmark with  
large-scale demonstrations, 2021. 2 673  
674  
675  
676[23] Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, Adeet  
Parikh, Aaron Lo, Abhishek Joshi, Ajay Mandlekar, and Yuke  
Zhu. Robocasa: Large-scale simulation of everyday tasks for  
generalist robots, 2024. 2 677  
678  
679  
680[24] NVIDIA, :, Johan Bjorck, Fernando Castañeda, Nikita Cher-  
niadev, Xingye Da, Runyu Ding, Linxi "Jim" Fan, Yu Fang,  
Dieter Fox, et al. Gr00t n1: An open foundation model for  
generalist humanoid robots, 2025. 6 681  
682  
683  
684[25] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo,  
Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel  
Haziza, Francisco Massa, Alaeldin El-Nouby, et al. Dinov2:  
Learning robust visual features without supervision, 2024. 6 685  
686  
687  
688[26] Wilbert Pumacay, Ishika Singh, Jiafei Duan, Ranjay Krishna,  
Jesse Thomason, and Dieter Fox. The colosseum: A bench-  
mark for evaluating generalization for robotic manipulation,  
2024. 2 689  
690  
691  
692

[27] William Shen, Caelan Garrett, Nishanth Kumar, Ankit Goyal, 693

- 694 Tucker Hermans, Leslie Pack Kaelbling, Tomás Lozano-  
695 Pérez, and Fabio Ramos. Differentiable gpu-parallelized  
696 task and motion planning. *arXiv preprint arXiv:2411.11833*,  
697 2024. 6
- 698 [28] Lirui Wang, Yiyang Ling, Zhecheng Yuan, Mohit Shridhar,  
699 Chen Bao, Yuzhe Qin, Bailin Wang, Huazhe Xu, and Xiao-  
700 long Wang. Gensim: Generating robotic simulation tasks  
701 via large language models. *arXiv preprint arXiv:2310.01361*,  
702 2023. 2, 5
- 703 [29] Yufei Wang, Zhou Xian, Feng Chen, Tsun-Hsuan Wang, Yian  
704 Wang, Katerina Fragkiadaki, Zackory Erickson, David Held,  
705 and Chuang Gan. Robogen: Towards unleashing infinite data  
706 for automated robot learning via generative simulation, 2024.  
707 2
- 708 [30] Yi Ru Wang, Carter Ung, Grant Tannert, Jiafei Duan,  
709 Josephine Li, Amy Le, Rishabh Oswal, Markus Grotz,  
710 Wilbert Pumacay, Yuquan Deng, Ranjay Krishna, Dieter Fox,  
711 and Siddhartha Srinivasa. Roboeval: Where robotic manipu-  
712 lation meets structured and scalable evaluation, 2025. 2