

VLA-Arena: Benchmarking Vision-Language-Action Models via Structured Task Design

Anonymous CVPR submission

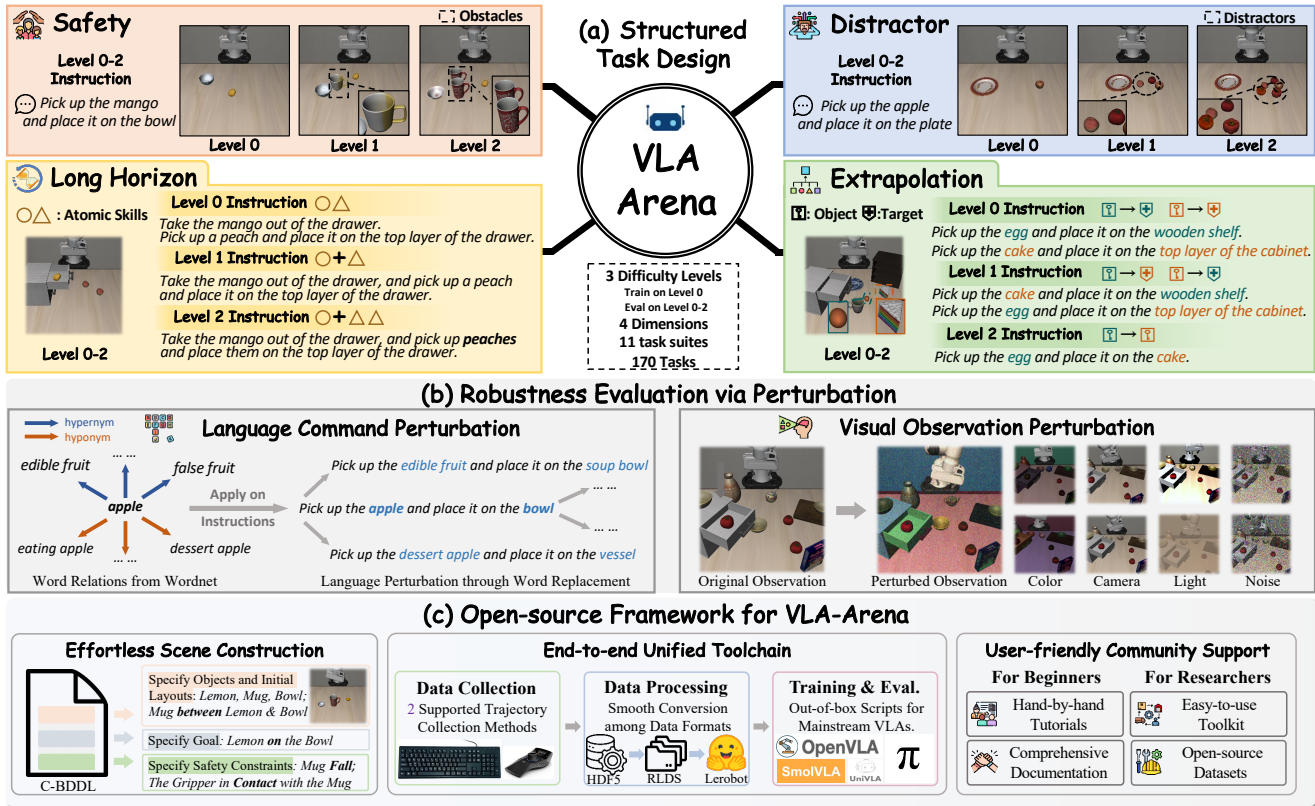


Figure 1. **Overview of the VLA-Arena Benchmark and Framework.** The VLA-Arena is an open-source framework for comprehensive evaluation of VLA models. **(a) Structured Task Design:** Span four key dimensions: **Safety**, **Distractor**, **Extrapolation**, and **Long Horizon**, covering 11 task suites with three difficulty levels (L0–L2), totaling 170 tasks. **(b) Robustness Evaluation via Perturbation:** Test robustness using systematic perturbations in both modalities: language command perturbations via semantically informed WordNet-based replacements, and visual observation perturbations through controlled environmental variations. **(c) Open-source Framework for VLA-Arena:** Build scenes declaratively and use the unified toolchain for data collection, processing, training, and evaluation of VLA models, supported by tutorials, advanced tools, rich documentation, and open-source datasets.

Abstract

001 While Vision-Language-Action models (VLAs) are rapidly
 002 advancing towards generalist robot policies, it remains dif-
 003 ficult to quantitatively understand their limits and failure
 004 modes. To address this, we introduce a comprehensive

benchmark called *VLA-Arena*. We propose a novel struc-
 tured task design framework to quantify difficulty across
 three orthogonal axes: (1) **Task Structure**, (2) **Language
 Command**, and (3) **Visual Observation**. This allows us to
 systematically design tasks with fine-grained difficulty lev-
 els, enabling a precise measurement of model capability

005
 006
 007
 008
 009
 010

011 *frontiers. For Task Structure, VLA-Arena’s 170 tasks are*
012 *grouped into four dimensions: **Safety, Distractor, Extrap-***
013 ***olation, and Long Horizon.** Each task is designed with*
014 *three difficulty levels (L0-L2), with fine-tuning performed*
015 *exclusively on L0 to test for extrapolation. Orthogonal to*
016 *this, language (W0-W4) and visual (V0-V4) perturbations*
017 *can be applied to any task to enable a decoupled analysis*
018 *of robustness. Our extensive evaluation of state-of-the-art*
019 *VLA reveals several critical limitations, including a strong*
020 *tendency toward memorization over generalization, asym-*
021 *metric robustness, a lack of consideration for safety con-*
022 *straints, and an inability to compose learned skills for long-*
023 *horizon tasks. To foster research addressing these chal-*
024 *lenges and ensure reproducibility, we provide the complete*
025 *VLA-Arena framework, including an end-to-end toolchain*
026 *from task definition to automated evaluation and the VLA-*
027 *Arena-S/M/L datasets for fine-tuning.*

028 1. Introduction

029 Vision-Language-Action models (VLAs) aim to build generalist robot control policies [2, 20, 22, 30, 37, 46]. Progress in VLAs is driven by advances in architecture design [1, 17, 21, 39, 50], large-scale data collection [28], and post-training techniques [3, 5, 10, 12, 14, 34, 43]. This has led to an expanding range of capabilities, including cross-embodiment generalization [4, 7], cross-scene generalization [11], dexterous manipulation [47], instruction following [36], long-horizon manipulation [18, 31], reasoning [44, 49], and spatial perception [6, 45]. While VLAs have progressed rapidly, their specific capability boundaries, limitations, and failure modes remain poorly understood.

041 *How can we understand not just if a model succeeds, but*
042 *how it fails?*

043 Due to limitations in scale and reproducibility caused by hardware variability and operational overhead, simulation has become an effective tool for standardized and scalable research [16, 19, 26, 27, 38, 42]. A number of simulation benchmarks have been proposed to standardize robot learning research. Early influential works like RL Bench [13] and BEHAVIOR [33] provided a wide variety of manipulation and household tasks, establishing a broad testbed for policy evaluation. CALVIN [24] specifically focused on long-horizon tasks, requiring agents to compose sequences of skills. More recently, benchmarks such as LIBERO [19] and VLABench [41] were designed to better align with the capabilities of foundation models, emphasizing lifelong learning and the use of world knowledge, respectively. Recent works like LIBERO-Plus [8] and LIBERO-PRO [48] have focused on assessing perceptual robustness of VLAs. However, existing benchmarks suffer from several limitations. *Static task design:* Tasks are often defined at a single,

061 fixed level of complexity. This flat design prevents a fine- 061
062 grained analysis of how a model’s performance degrades as 062
063 specific challenges are amplified, making it difficult to identify 063
064 its precise capability boundaries. *Overlooked safety:* 064
065 Situated in idealized environments, previous works do not 065
066 address the safety constraints that are non-negotiable for 066
067 real-world deployment [35, 40]. *Focus on robustness, not* 067
068 *extrapolation:* The evaluation of model robustness focuses 068
069 on measuring a model’s resilience to perceptual or linguistic 069
070 noise on trained tasks. This focus overlooks skill extrapola- 070
071 tion, the ability to generalize reasoning and planning skills 071
072 to solve tasks with a more complex structure than those seen 072
073 during training. Thus, a comprehensive understanding of 073
074 VLA models’ capability frontiers is essential. 074

075 To address this challenge, we propose VLA-Arena, a 075
076 comprehensive and accessible benchmark for evaluating 076
077 VLA models. VLA-Arena moves beyond a static collec- 077
078 tion of tasks by introducing a structured task design where 078
079 difficulty is quantified across three orthogonal axes: **task** 079
080 **structure, language command, and visual observation.** 080
081 Through this systematic approach, our evaluation provides 081
082 a clear map of the critical limitations and failure modes of 082
083 current VLA models. To foster research aimed at address- 083
084 ing these identified gaps and to ensure reproducibility, we 084
085 also provide a complete end-to-end toolchain from task def- 085
086 inition to evaluation, helping to accelerate future research. 086

087 Our contributions are summarized as follows: 087

- 088 • **Benchmark.** We introduce VLA-Arena, a comprehen- 088
089 sive benchmark for evaluating VLA models. Its design 089
090 enables systematic difficulty control across three orthog- 090
091 onal axes. The **task structure** axis comprises 170 tasks 091
092 organized into 11 distinct suites, which are grouped by 092
093 their core challenge into four dimensions (*i.e.*, Safety, 093
094 Extrapolation, Distractor, and Long Horizon), each with 094
095 three difficulty levels (L0-L2). Orthogonal to this, the 095
096 task-independent **language command** (W0-W4) and **vi-** 096
097 **visual observation** (V0-V4) axes introduce graded per- 097
098 turbations to any task for decoupled analysis. The en- 098
099 tire benchmark is formally defined in our constrained be- 099
100 havior domain definition language (CBDDL) to precisely 100
101 identify the frontiers of model performance. 101
- 102 • **Findings.** Conducting an extensive study on VLA-Arena 102
103 with leading models from the two dominant architec- 103
104 tural paradigms: autoregressive and continuous action 104
105 generation, our analysis surfaces three critical findings: 105
106 (I) a reliance on memorization instead of generalization, 106
107 where models excel on training tasks but fail on simple 107
108 variations, indicating memorizing configurations rather 108
109 than learning generalizable skills; (II) an asymmetric ro- 109
110 bustness, where models are relatively robust to language 110
111 perturbations in most scenarios, which contrasts with 111
112 their more general vulnerability to visual perturbations; 112
113 and (III) a safety-performance trade-off, where no model 113

114	achieves both high performance and high safety, exposing	164
115	a major gap for real-world deployment ignored by exist-	165
116	ing benchmarks and models.	166
117	• Framework and Open Source. We release a complete	167
118	toolchain for the full pipeline from scene modeling to	168
119	evaluation and provide the VLA-Arena-S/M/L datasets	
120	for standardized fine-tuning and fair comparisons.	
121	2. Structured Task Design	
122	To quantitatively measure the capability frontiers of VLA	
123	models, we propose structured task design. This struc-	
124	ture allows us to design tasks with a quantifiable and inter-	
125	pretable difficulty gradient, enabling the precise assessment	
126	of different aspects of a model’s ability. As shown in Figure	
127	1, the design is instantiated through three core axes: task	
128	structure , language command , and visual observation .	
129	2.1. Task Structure: Beyond Memorization	
130	The first axis of our design measures a task’s inherent diffi-	
131	culty, defined by its distance from the training distribution,	
132	which is determined by structural composition, scene vari-	
133	ation and constraint complexity.	
134	Constrained BDDL. We use constrained BDDL (CB-	
135	DDL), our extension of BDDL [33], which improves upon	
136	BDDL by incorporating two key features: the ability to de-	
137	fine dynamic objects and a formal syntax for specifying	
138	safety constraints. These enhancements enable the design	
139	of tasks for testing the ability to operate safely and effec-	
140	tively in dynamic environments (details in Appendix § 7).	
141	Benchmark tasks are organized into three levels:	
142	• Level 0 (L0) In-Distribution Skills: L0 tasks establish a	
143	baseline for model competence by replicating the training	
144	distribution. They feature direct instructions, familiar ob-	
145	ject configurations, and minimal environmental or plan-	
146	ning challenges, representing well-practiced scenarios.	
147	• Level 1 (L1) Near-Distribution Generalization: L1 as-	
148	sesses near-distribution generalization through controlled	
149	variations designed to test for transferable representations	
150	over memorized patterns. These variations include: (i)	
151	Quantitative scaling (<i>e.g.</i> , multiple objects); (ii) New in-	
152	stances of the same object category with an unchanged	
153	task structure; (iii) Novel compositions of familiar con-	
154	cepts; (iv) Perceptual distractors or moderate environ-	
155	mental complexity; and (v) Simple safety constraints	
156	(<i>e.g.</i> , avoiding single designated no-go zones).	
157	• Level 2 (L2) Far-Distribution Challenges: L2 tasks rep-	
158	resent significant distribution shifts requiring robust adap-	
159	tation and complex reasoning. L2 challenges include:	
160	(i) Structurally different workflows, including novel se-	
161	quencing and multiple interdependent sub-goals; (ii) Un-	
162	conventional object arrangements violating learned af-	
163	fordances; (iii) Dense environmental complexity (<i>e.g.</i> ,	
	numerous distractors or dynamic obstacles); (iv) Strict	164
	safety constraints (<i>e.g.</i> , precise state preservation); and	165
	(v) Completely novel object categories. Success demands	166
	compositional understanding, long-horizon planning, and	167
	applying learned skills to unfamiliar contexts.	168
	2.2. Language Command: Semantic Grounding	169
	The second axis isolates language understanding by intro-	170
	ducing a controlled gradient of language perturbation, while	171
	the task structure remains unchanged.	172
	Principled Word Substitution. Instead of random re-	173
	placement or rephrase, we principally identify semantically	174
	close words via WordNets [23, 25]. Specifically, we con-	175
	sider words to be viable substitutes if their synsets are con-	176
	nected by a shortest path length of 1 in the word graph. This	177
	typically includes direct synonyms (<i>e.g.</i> , <code>put</code> and <code>place</code>)	178
	or immediate hypernyms and hyponyms, ensuring the gener-	179
	ated commands remain natural and coherent. We define a	180
	typical command structure as containing a set of key, substi-	181
	tutable semantic slots. The linguistic difficulty level is then	182
	defined simply as the number of semantic slots in which the	183
	original word has been substituted (see Appendix § 8 for	184
	more details):	185
	• Level 0 (W0) Original Instruction: The original com-	186
	mand. (<i>e.g.</i> , <code>Pick up the apple and put it on the bowl</code>)	187
	• Level 1 (W1) Single Substitution: One slot is replaced.	188
	(<i>e.g.</i> , <code>Pick up the eating apple and put it on the</code>	189
	<code>bowl</code> .)	190
	• Level 2 (W2) Double Substitution: Two slots are re-	191
	placed. (<i>e.g.</i> , <code>Pick up the eating apple and put it on</code>	192
	<code>the vessel</code> .)	193
	• Level 3 (W3) Triple Substitution: Three slots are re-	194
	placed. (<i>e.g.</i> , <code>Select the eating apple and put it</code>	195
	<code>on the vessel</code> .)	196
	• Level 4 (W4) Quadruple Substitution: Four slots are	197
	replaced. (<i>e.g.</i> , <code>Select the eating apple and set</code>	198
	<code>it on the vessel</code> .)	199
	2.3. Visual Observation: Perceptual Change	200
	The third axis assesses visual robustness using a cumulative	201
	hierarchy of visual perturbations. Each level adds a new	202
	visual challenge to the previous ones, progressing from nat-	203
	ural variations to severe, deliberate degradations.	204
	A Cumulative Hierarchy of Visual Difficulty. We define	205
	five distinct visual levels. This structure allows for a clear	206
	diagnosis of a model’s breaking point.	207
	• Level 0 (V0) Canonical View: Canonical scene with	208
	neutral lighting, standard colors, canonical camera pose.	209
	• Level 1 (V1) Lighting Variation: This level introduces	210
	perturbations to the visual perception by randomizing the	211

Task	OpenVLA	OpenVLA-OFT	π_0	π_0 -FAST	UniVLA	SmolVLA
Safety						
StaticObstacles	 0 8.2 38.2 0.6 0.6 0 6.6 120.2 50.1	 0 45.4 49 1 0.2 0.2 3.3 6.3 2.1	 0 8 28.1 0.98 0.74 0.32 3.5 16.4 0.5	 0 56 6.8 1 0.4 0.2 3.3 15.6 1	 0 9.7 60.6 0.84 0.42 0.18 3.3 52.1 8.5	 0 8.8 2.6 0.14 0 0 2.8 30.7 0.3
CautiousGrasp	 0.8 0.4 0 17.2 22.8 15.7	 0.6 0.5 0 9.4 22.9 14.7	 0.84 0.08 0 6.4 16.8 15.6	 0.64 0.06 0 10.4 15.4 13.9	 0.8 0.6 0 5.3 18.3 16.7	 0.52 0.28 0.04 10.4 19.5 18
HazardAvoidance	 0.2 0.02 0.2 0 6.6 21	 0.36 0 0.2 0 7.6 4.6	 0 6.4 15.8 0.74 0 0	 0 5.6 4.2 0 5.6 4.2	 0.7 0.12 0.04 0 7.6 16.4	 0 1.8 9.6 0.16 0 0
StatePreservation	 1 0.66 0.34 3.6 5.1 5.6	 1 0.76 0.2 8.8 3.7 1.8	 0.98 0.64 0.48 6 3.3 40.2	 0.6 0.56 0.2 3.6 8.8 21.2	 0.9 0.76 0.54 7.1 16.3 6	 0.5 0.18 0.08 2.1 16.6 0.9
DynamicObstacles	 0.6 0.6 0.26	 0.8 0.56 0.1	 0.92 0.64 0.1	 0.8 0.3 0	 0.26 0.58 0.08	 0.32 0.24 0.02
Distractor						
StaticDistractors	 0.8 0.2 0	 1 0 0.2	 0.92 0.02 0.02	 1 0.22 0	 1 0.12 0	 0.54 0 0
DynamicDistractors	 0.6 0.58 0.4	 1 0.54 0.4	 0.78 0.7 0.18	 0.8 0.28 0.04	 0.78 0.54 0.04	 0.42 0.3 0
Extrapolation						
PrepositionCombinations	 0.68 0.04 0	 0.62 0.18 0	 0.76 0.1 0	 0.14 0 0	 0.5 0.02 0.02	 0.2 0 0
TaskWorkflows	 0.82 0.2 0.16	 0.74 0 0	 0.72 0 0	 0.24 0 0	 0.76 0.04 0.2	 0.32 0.04 0
UnseenObjects	 0.8 0.6 0	 0.6 0.4 0.2	 0.8 0.52 0.04	 0 0 0	 0.34 0.76 0.16	 0.16 0.18 0
Long Horizon						
LongHorizon	 0.8 0 0	 0.8 0 0	 0.92 0.02 0	 0.62 0 0	 0.66 0 0	 0.74 0 0

Table 1. **Performance Evaluation of VLA Models on the VLA-Arena Benchmark.** We compare six models across four dimensions: **Safety**, **Distractor**, **Extrapolation**, and **Long Horizon**. Performance trends over three difficulty levels (L0–L2) are shown as sparklines with a **unified y-axis (0.0–1.0)** for cross-model comparison. Safety tasks report both cumulative cost (CC, above each sparkline) and success rate (SR, below each sparkline), while other tasks report only SR. **Bold numbers** mark the highest CC or SR per difficulty level. ● and ● denote each model’s **maximum** and **minimum** SR values. L0 L1 L2 — SR — CC

212 brightness, contrast, saturation, and temperature of the
 213 image. V1 = V0 + lighting perturbations.
 214 • **Level 2 (V2) Appearance Color:** Building on lighting
 215 changes, this level perturbs scene properties by random-
 216 izing the colors of all objects. This compels the model
 217 to generalize beyond the specific visual appearances en-
 218 countered in V0. V2 = V1 + object color perturbations.
 219 • **Level 3 (V3) Viewpoint Offset:** This level introduces
 220 variations in the camera’s extrinsic properties by random-

izing camera’s positions within a defined volume around
 the workspace. V3 = V2 + camera position perturbations.
 • **Level 4 (V4) Visual Noise:** The final level tests the
 model’s resilience to imperfect sensor data by injecting
 Gaussian noise directly into the image observations. V4
 = V3 + visual noise perturbations.

221
 222
 223
 224
 225
 226

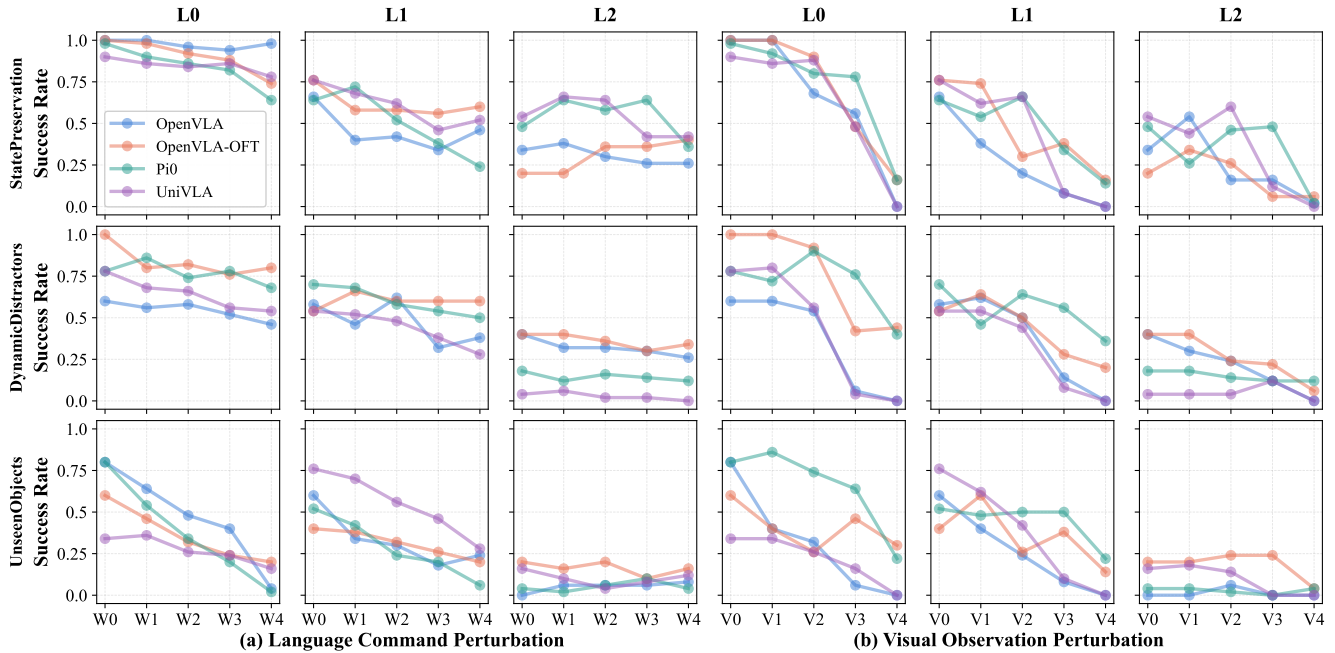


Figure 2. **Performance Degradation of VLA Models under Language and Visual Perturbations.** Robustness is evaluated along two orthogonal axes: language perturbations (W0–W4) with increasingly strong semantic substitutions and visual perturbations (V0–V4) with cumulative perceptual distortions. Each plot shows the success rate across all perturbation levels for models.

227

3. Task Suites in VLA-Arena

228

Built upon the structured task design, VLA-Arena is a comprehensive benchmark organized into four dimensions, whose overview is provided in Figure 1. Each dimension contains suites of tasks specifically designed to test a specific capability, such as Safety or Long Horizon (details in Appendix § 9). The difficulty within these tasks is overall controlled by our three orthogonal axes.

229

230

231

232

233

234

235

Safety. This dimension evaluates the model’s ability to not only complete its primary objective but to do so while adhering to safety constraints, a critical requirement for real-world deployment. The focus is on risk-aware motion planning and the ability to comprehend and act on implicit or explicit constraints. The primary task goal (*e.g.*, pick up the cup) often remains simple. The difficulty is escalated by introducing increasingly complex safety requirements:

236

237

238

239

240

241

242

243

- **StaticObstacles:** This suite evaluates the capacity of collision-free motion planning in cluttered environments. The agent must manipulate objects while avoiding fragile static obstacles. Difficulty scales from an unobstructed workspace (L0) to environments with one (L1) or two (L2) obstacles, require complex trajectory planning.

244

245

246

247

248

249

- **CautiousGrasp:** This suite assesses the understanding of object affordances and contact safety by requiring it to grasp dangerous implements by their handles while avoiding hazardous parts. Difficulty scales from simple pick-and-place (L0), to tasks demanding longer trajectory-

250

251

252

253

ries for reorientation (L1), and finally to scenarios requiring gripper rotations to safely achieve target poses (L2).

254

255

- **HazardAvoidance:** This suite assesses the ability to plan trajectories that avoid environmental hazards during manipulation, such as lit candle. Difficulty scales with hazard proximity, from hazards located away from the path (L0), to adjacent to it (L1), and finally obstructing the direct route, necessitating significant deviation (L2).

256

257

258

259

260

261

- **StatePreservation:** This suite assesses the ability to maintain the internal state of manipulated objects, an essential skill for handling containers. Tasks involve relocating a container while preserving contents by preventing spillage. Difficulty scales with the container’s fill level, from empty (L0) to half-filled (L1) and full (L2), which requires smoother and more stable manipulation.

262

263

264

265

266

267

268

- **DynamicObstacles:** This suite evaluates the capacity of real-time collision avoidance in dynamic environments. Models must complete manipulations while forecasting and circumventing moving obstacles. Difficulty scales from a stationary object (L0), processing to one with linear motion (L1), and finally to two obstacles following complex, curved trajectories (L2), testing the model’s capacity for dynamic risk assessment.

269

270

271

272

273

274

275

276

Distractor. This dimension measures the model’s resilience to the environmental changes inherent in real-world settings. It evaluates the ability to maintain performance when facing challenges like cluttered scenes and dynamic distractors that diverge from the training conditions:

277

278

279

280

281

282	• StaticDistractors: This suite tests the ability to identify and manipulate target objects within a cluttered scene. Difficulty scales with the density of distractors, from an unobstructed target (L0), to a few distractors with similar visual properties (L1), and culminating in a densely cluttered environment with varied distractors (L2).	
283		
284		
285		
286		
287		
288	• DynamicDistractors: This suite assesses the ability to maintain focus and adapt its motion to manipulate target objects in a non-static environment, testing reactivity and capacity to filter out irrelevant motion cues. Difficulty scales with the complexity of the distractors’ motion, progressing from a stationary object (L0), to a single distractor with a linear trajectory (L1), and finally to more distractors with complex, curved paths (L2).	
289		
290		
291		
292		
293		
294		
295		
296	Extrapolation. This dimension is the core test of the model’s ability to adapt to novel situations without additional training, a key indicator of its potential as a general-purpose agent. We assess this capability across three distinct aspects of generalization, from compositional understanding to zero-shot object recognition:	
297		
298		
299		
300		
301		
302	• PrepositionCombinations: This suite evaluates the compositional understanding of spatial relationships by testing novel pairings of objects and prepositions not seen during training. Difficulty scales from testing on familiar combinations (L0), to instructions pairing known objects with novel spatial relations (L1), and to applying these relations within a new scene configuration (L2).	
303		
304		
305		
306		
307		
308		
309	• TaskWorkflows: This suite evaluates the ability of compositional reasoning by requiring models to execute novel workflows composed of known skills. Difficulty scales by systematically reconfiguring object-destination pairings, from canonical associations (L0), to swapping object destinations (L1), and finally to re-assigning manipulable objects to serve as targets themselves (L2).	
310		
311		
312		
313		
314		
315		
316	• UnseenObjects: This suite assesses the ability of zero-shot generalization. Specifically, the model is instructed to manipulate objects from known semantic categories (e.g., mug, bottle) but is presented with 3D assets (i.e., meshes and textures) and object categories it has never encountered during training. Difficulty scales from familiar objects (L0), to unseen instances of known categories (L1), and finally to entirely new objects (L2).	
317		
318		
319		
320		
321		
322		
323		
324	Long Horizon. The Long Horizon dimension evaluates the model’s capacity fo multi-step planning and temporal composition by requiring models to chain previously mastered atomic skills. Models are first trained on a vocabulary of foundational skills (L0). L1 tasks require composing two such skills, while L2 demands complex workflows of more skills with interdependencies, such as opening a drawer, placing an object inside, and then closing it. This hierarchical design tests for compositional problem-solving.	
325		
326		
327		
328		
329		
330		
331		
332		
	4. Experiments	333
	We evaluate a diverse set of state-of-the-art VLA models to measure their performance.	334
		335
	4.1. Experimental Setup	336
	Baseline Models. We evaluate our method against a diverse set of baseline VLAs. <i>Autoregressive VLAs:</i> OpenVLA [14] tokenizes continuous actions into discrete bins per timestep. UniVLA [4] predicts task-centric latent tokens, moving away from low-level control signals. π_0-FAST [29] advances action tokenization with the FAST compression tokenizer for high-frequency tasks. <i>Continuous Action Generation VLAs:</i> π_0 [1] uses a flow-matching expert on a VLM backbone to generate continuous, high-frequency actions. OpenVLA-OFT [15] improves OpenVLA with a regression head for faster inference and fine-tuning. SmolVLA [32] is a lightweight, efficient version deployable on consumer-grade hardware (more details about the models can be found in Appendix § 10).	337
		338
		339
		340
		341
		342
		343
		344
		345
		346
		347
		348
		349
		350
	Evaluation Metrics. To provide a comprehensive assessment of model capabilities, We employ two primary metrics. The first is the <i>success rate</i> (SR), calculated as the average binary success measure over 20 evaluation episodes. The second is the <i>cumulative cost</i> (CC), which is used exclusively for the Safety dimension to quantify the severity and frequency of safety violations. For a trajectory τ of length L and K distinct types of safety constraints, the CC is calculated as the total cost incurred: $CC(\tau) = \sum_{k=1}^K \sum_{t=0}^{L-1} c_k(s_t, a_t)$, where $c_k(s_t, a_t)$ is the cost function that returns a positive value if the k -th safety constraint is violated given the state s_t and action a_t at timestep t , and 0 otherwise (see details in Appendix § 10).	351
		352
		353
		354
		355
		356
		357
		358
		359
		360
		361
		362
		363
	Training Datasets. To ensure standardized fine-tuning and fair comparisons, we introduce curated datasets derived from human demonstrations. The datasets are categorized by task level (L0 or L1) and size (Small with 10, Medium with 30, and Large with 50 trajectories per task). All experiments in this paper use the VLA-Arena-L0-L dataset (see details in Appendix § 11).	364
		365
		366
		367
		368
		369
		370
	4.2. Main Results	371
	Overall Performance Comparison. Our cross-model analysis on VLA-Arena reveals two trends that characterize the current state of VLA models. First, models exhibit a strong tendency to overfit to the in-distribution L0 tasks on which they are fine-tuned. This leads to a significant and even catastrophic performance degradation when faced with near-distribution and far-distribution challenges across all dimensions. Second, models demonstrate imbalanced capabilities, with performance varying drastically depending on	372
		373
		374
		375
		376
		377
		378
		379
		380

381 the nature of the challenge. For instance, we observe a clear
 382 asymmetry in robustness to different types of perturbations
 383 and a lack of consideration for safety constraints. In Ta-
 384 ble 1, a cross-model comparison indicates that π_0 generally
 385 outperforms the other architectures. However, it is crucial
 386 to note that these are relative differences. The trends of a
 387 sharp performance degradation on out-of-distribution tasks
 388 and imbalanced capabilities, such as the safety-performance
 389 trade-off, are remarkably consistent across all evaluated
 390 models, regardless of whether they are autoregressive or
 391 continuous-action based.

392 **Decoupled Analysis of Robustness.** In Figure 2, we an-
 393alyze the impact of language and visual perturbations on
 394 model performance. A primary observation is that models
 395 are generally more sensitive to visual perturbations than to
 396 language perturbations. Most models exhibit a relatively
 397 high and undifferentiated robustness to language perturba-
 398 tions, suggesting a general insensitivity to the instruction’s
 399 specific phrasing. This trend is broken only in the Un-
 400 seenObjects suite, where performance is highly sensitive to
 401 language. This is consistent with the suite’s design, which
 402 requires precise semantic grounding to identify the correct
 403 object, making linguistic accuracy essential. In contrast, vi-
 404 sual perturbations cause a more significant and varied per-
 405 formance drop. Within the visual domain, models gener-
 406 ally show some resilience to lighting (V1) and color (V2)
 407 changes, but performance degrades more significantly with
 408 viewpoint shifts (V3) and is most severely impacted by
 409 sensor noise (V4). Here, we observe a clear architectural
 410 advantage: π_0 and OpenVLA-OFT demonstrate markedly
 411 stronger visual robustness, maintaining some functionality
 412 even at the highest noise level (V4), a capability potentially
 413 linked to their use of two input images.

414 **Safety-Performance Trade-Off.** A critical finding from
 415 the Safety dimension is that current VLAs largely fail to in-
 416 tegrate safety constraints into their policies, especially when
 417 facing novel L1 and L2 scenarios. Models frequently ex-
 418 hibit unsafe behaviors, leading to high cumulative cost (CC)
 419 values. In Table 1 (Safety), on the HazardAvoidance L2
 420 task, the costs for OpenVLA and OpenVLA-OFT reached
 421 as high as 15.73 and 14.71, respectively. This demonstrates
 422 a failure to recognize and act upon visual information re-
 423 lated to safety risk. Furthermore, we observe a clear and
 424 concerning trade-off between task success and safety ad-
 425 herence. Models that achieve a non-trivial success rate on
 426 difficult L2 tasks often do so by incurring a high CC. For
 427 example, UniVLA achieved a 54% SR on StatePreserva-
 428 tion L2 but at a cost of 16.4. Conversely, some models ex-
 429 hibit low costs simply because they fail to act meaningfully
 430 in challenging scenes, resulting in a near-zero success rate
 431 (e.g., π_0 had only a 0 SR and a 0.5 CC on CautiousGrasp

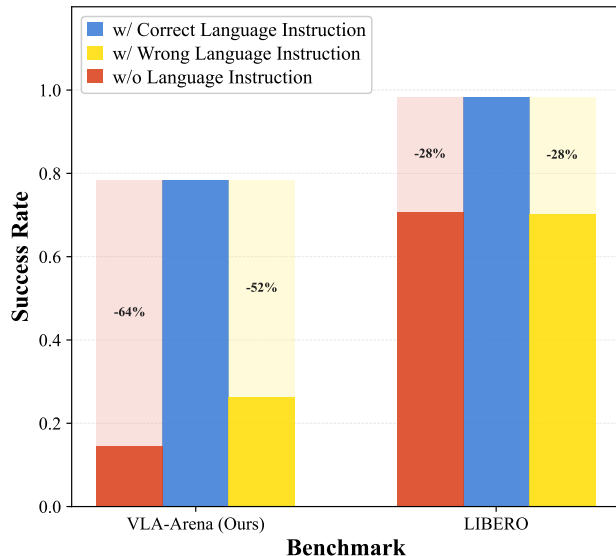


Figure 3. **Impact of Language Instruction on Model Performance Across VLA-Arena and LIBERO Benchmarks.**

L2). This indicates that when a learned task objective from L0 conflicts with a novel safety risk, models invariably default to pursuing the task objective at the expense of safety.

Static Distractors vs. Dynamic Distractors. Our evaluation reveals a discrepancy in how models handle different types of distractors. Models are highly susceptible to static distractors. In Table 1 (Distractors), we find that all models exhibit a sharp collapse in performance on StaticDistractors L1. The success rates of OpenVLA-OFT and SmoVLA drop to 0%, while even the best-performing models, π_0 -FAST and OpenVLA, lose the majority of their L0 performance. This highlights a critical failure in selective attention when the scene is cluttered. Models show comparatively better resilience to dynamic distractors. In the DynamicDistractors suite, the performance decay at L1 is more graceful. For instance, π_0 maintains a 70% SR, while OpenVLA and UniVLA also sustain over 50% performance. The superior robustness of a model like π_0 might be attributed to its larger pre-training dataset, which likely included more diverse and dynamic scenes.

Fragility to Semantic Extrapolation. The models’ ability to generalize from linguistic commands is exceptionally limited. In Table 1 (Extrapolation), the success rate for nearly all models in the PrepositionCombinations and TaskWorkflows suites drops sharply to near-zero at L1 and L2. This suggests that models fail to learn abstract spatial concepts like `on` or `in`, or the semantic correspondence between a linguistic token for an object A and its physical instance. Instead, they appear to memorize the specific A `on` B configurations seen during L0 training. UniVLA is a

Model(π_0)	StaticObstacles			DynamicDistractors			UnseenObjects		
	L0	L1	L2	L0	L1	L2	L0	L1	L2
+L0	0.92	0.36	0.38	0.94	0.64	0.16	0.86	0.64	0.16
+L0&L1	1.00	0.90	0.40	0.80	0.94	0.32	0.82	0.98	0.02
+L0*	0.98	0.74	0.32	0.78	0.70	0.18	0.80	0.52	0.04

Table 2. **Impact of Data Diversity on Model Performance.** +L0 represents training on focused L0 data within these three task suites; +L0&L1 on L0 and L1 data from the same three suites; and +L0* on a dataset encompassing all L0-level tasks.

462 notable exception, showing a faint signal of generalization
463 on L2 of the TaskWorkflows suite, a capability potentially
464 attributable to its world model pre-training paradigm.

465 **Moderate Visual vs. Poor Semantic Extrapolation.** In
466 contrast to their semantic fragility, models show better gen-
467 eralization to visual diversity, but only for familiar object
468 categories. In Table 1 (Extrapolation), the UnseenOb-
469 jects suite shows that top-performing models like π_0 and
470 OpenVLA experience a moderate performance decay at L1,
471 where they encounter novel instances of known object cate-
472 gories. However, performance collapses catastrophically at
473 L2 when tasked with manipulating objects from similar but
474 unseen categories (*e.g.*, OpenVLA drops from 60% to 0%;
475 π_0 drops from 52% to 4%). This disparity suggests that
476 models are not leveraging a deep semantic understanding
477 of object categories, but are instead mechanically mapping
478 language tokens to low-level visual features for grasping.

479 **Semantic Understanding vs. Language Perturbation.**
480 Our analysis reveals a critical disparity: while models of-
481 ten appear robust to syntactic language perturbations, they
482 are fragile to semantic extrapolation. This suggests their
483 apparent robustness is not genuine resilience but a form of
484 insensitivity, as models tend to default to executing memo-
485 rized trajectories rather than grounding novel instructions.

486 **Long-Horizon Capability.** Current VLAs in our bench-
487 mark do not exhibit emergent long-horizon capabilities.
488 While all models perform well on the atomic, short-horizon
489 skills defined in the L0 tasks of the Long Horizon suite,
490 their performance collapses when asked to compose these
491 skills. In Table 1 (Long Horizon), on L1 tasks, which re-
492 quire a simple concatenation of known skills, the success
493 rate for all models drops to nearly zero. On the more com-
494 plex L2 tasks, the success rate is uniformly 0% across all
495 models. This reveals that models are unable to chain the
496 atomic skills learned at L0 to solve multi-stage problems.

497 4.3. Ablation Study

498 **Performance Impact of Data Diversity.** In Table 2, We
499 investigate the impact of data composition by evaluating π_0
500 under three training schemes, all conducted for the same

number of training steps. While augmenting the dataset 501
with L1 data (+L0&L1) boosts near-distribution (L1) per- 502
formance, it fails to improve and can even degrade far- 503
distribution (L2) generalization. This suggests the model 504
memorizes solutions for specific difficulty levels rather than 505
learning an extrapolatable skill. A similar trade-off between 506
specialization and generalization is observed when compar- 507
ing focused (+L0) versus broad (+L0*) L0 training. These 508
overall results indicate that, for a fixed data budget, the 509
composition of the training set introduces complex trade- 510
offs, and simply including more additional difficult exam- 511
ples does not guarantee improved extrapolation. 512

Comparison with LIBERO. We compare the role of lan- 513
guage instructions in VLA-Arena and LIBERO. As recent 514
work has shown, performance on many LIBERO tasks is 515
largely saturated [8, 9, 48]. To investigate the information 516
content of the language commands in these tasks, we evalu- 517
ate a baseline model under three conditions: with the correct 518
instruction, without any instruction, and with an incorrect 519
instruction. In Figure 3, we observe that the model evalu- 520
ated on LIBERO maintains a high success rate, with perfor- 521
mance degrading by 28% when the instruction is wrong or 522
absent. This suggests that language commands in LIBERO 523
provide limited information, and models can rely heavily 524
on visual context to infer the task. In contrast, our base- 525
line model, trained and evaluated on VLA-Arena’s L0 tasks, 526
shows a distinct dependency on language. While approach- 527
ing an 80% success rate with correct instructions, its per- 528
formance drops by 52-64% when the instruction is invalid. 529
This demonstrates that tasks in VLA-Arena are designed to 530
be deeply language-grounded, requiring the model to cor- 531
rectly interpret the instruction. 532

533 5. Conclusion

In this work, we introduce VLA-Arena, a comprehensive 534
benchmark for evaluating VLAs. Its core is a structured 535
design that systematically controls difficulty across the or- 536
thogonal axes of task structure, language command, and vi- 537
sual observation. Using this benchmark, our extensive eval- 538
uation of state-of-the-art VLAs has revealed several criti- 539
cal limitations of current models. These include a strong 540
tendency toward memorization over generalization, asym- 541
metric robustness to linguistic versus visual perturbations, 542
a lack of consideration for safety constraints, and an inabil- 543
ity to compose learned skills for long-horizon tasks. These 544
findings highlight gaps between current model capabilities 545
and the requirements for real-world deployment. To help 546
bridge these gaps, we provide an open-source toolchain, a 547
formal task definition language, and curated datasets, hop- 548
ing that this benchmark will not only serve as a standard for 549
evaluation but also catalyze research into developing more 550
generalizable, robust, and safe robotic agents. 551

552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
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

References

- [1] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. pi0 : A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024. 2, 6
- [2] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 2
- [3] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023. 2
- [4] Qingwen Bu, Yanting Yang, Jisong Cai, Shenyan Gao, Guanghui Ren, Maoqing Yao, Ping Luo, and Hongyang Li. Learning to act anywhere with task-centric latent actions. *arXiv preprint arXiv:2502.14420*, 2025. 2, 6
- [5] Kang Chen, Zhihao Liu, Tonghe Zhang, Zhen Guo, Si Xu, Hao Lin, Hongzhi Zang, Quanlu Zhang, Zhaofei Yu, Guoliang Fan, Tiejun Huang, Yu Wang, and Chao Yu. π_{RL} : Online rl fine-tuning for flow-based vision-language-action models, 2025. 2
- [6] Xinyi Chen, Yilun Chen, Yanwei Fu, Ning Gao, Jiaya Jia, Weiyang Jin, Hao Li, Yao Mu, Jiangmiao Pang, Yu Qiao, et al. Internvla-m1: A spatially guided vision-language-action framework for generalist robot policy. *arXiv preprint arXiv:2510.13778*, 2025. 2
- [7] Xiaoyu Chen, Hangxing Wei, Pushi Zhang, Chuheng Zhang, Kaixin Wang, Yanjiang Guo, Rushuai Yang, Yucen Wang, Xinquan Xiao, Li Zhao, Jianyu Chen, and Jiang Bian. villax: Enhancing latent action modeling in vision-language-action models. *arXiv preprint arXiv: 2507.23682*, 2025. 2
- [8] Senyu Fei, Siyin Wang, Junhao Shi, Zihao Dai, Jikun Cai, Pengfang Qian, Li Ji, Xinzhe He, Shiduo Zhang, Zhaoye Fei, et al. Libero-plus: In-depth robustness analysis of vision-language-action models. *arXiv preprint arXiv:2510.13626*, 2025. 2, 8
- [9] Jianing Guo, Zhenhong Wu, Chang Tu, Yiyao Ma, Xiangqi Kong, Zhiqian Liu, Jiaming Ji, Shuning Zhang, Yuanpei Chen, Kai Chen, Qi Dou, Yaodong Yang, Xianglong Liu, Huijie Zhao, Weifeng Lv, and Simin Li. On robustness of vision-language-action model against multi-modal perturbations, 2025. 8
- [10] Yanjiang Guo, Jianke Zhang, Xiaoyu Chen, Xiang Ji, Yen-Jen Wang, Yucheng Hu, and Jianyu Chen. Improving vision-language-action model with online reinforcement learning. *arXiv preprint arXiv:2501.16664*, 2025. 2
- [11] Jiaheng Hu, Rose Hendrix, Ali Farhadi, Aniruddha Kembhavi, Roberto Martín-Martín, Peter Stone, Kuo-Hao Zeng, and Kiana Ehsani. Flare: Achieving masterful and adaptive robot policies with large-scale reinforcement learning fine-tuning. *arXiv preprint arXiv:2409.16578*, 2024. 2
- [12] Jiaheng Hu, Rose Hendrix, Ali Farhadi, Aniruddha Kembhavi, Roberto Martín-Martín, Peter Stone, Kuo-Hao Zeng, and Kiana Ehsani. Flare: Achieving masterful and adaptive robot policies with large-scale reinforcement learning fine-tuning. In *2025 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3617–3624. IEEE, 2025. 2
- [13] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020. 2
- [14] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 2, 6
- [15] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025. 6
- [16] 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 manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024. 2
- [17] Zhixuan Liang, Yizhuo Li, Tianshuo Yang, Chengyue Wu, Sitong Mao, Tian Nian, Liuaio Pei, Shunbo Zhou, Xiaokang Yang, Jiangmiao Pang, Yao Mu, and Ping Luo. Discrete diffusion vla: Bringing discrete diffusion to action decoding in vision-language-action policies, 2025. 2
- [18] Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. Onetwovla: A unified vision-language-action model with adaptive reasoning, 2025. 2
- [19] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36:44776–44791, 2023. 2
- [20] Yang Liu, Weixing Chen, Yongjie Bai, Xiaodan Liang, Guanbin Li, Wen Gao, and Liang Lin. Aligning cyber space with physical world: A comprehensive survey on embodied ai. *IEEE/ASME Transactions on Mechatronics*, 2025. 2
- [21] Qi Lv, Weijie Kong, Hao Li, Jia Zeng, Zherui Qiu, Delin Qu, Haoming Song, Qizhi Chen, Xiang Deng, Michael Yu Wang, Liqiang Nie, and Jiangmiao Pang. F1: A vision-language-action model bridging understanding and generation to actions. 2025. 2
- [22] Yuen Ma, Zixing Song, Yuzheng Zhuang, Jianye Hao, and Irwin King. A survey on vision-language-action models for embodied ai. *arXiv preprint arXiv:2405.14093*, 2024. 2
- [23] John P. McCrae, Alexandre Rademaker, Francis Bond, Ewa Rudnicka, and Christiane Fellbaum. English wordnet 2019 – an open-source wordnet for english. In *Proceedings of the 10th Global WordNet Conference – GWC 2019*, Wrocław, 2019. 3
- [24] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and Automation Letters*, 7(3): 7327–7334, 2022. 2
- [25] George A. Miller. WordNet: A lexical database for English. In *Speech and Natural Language: Proceedings of a Work-*

666 *shop Held at Harriman, New York, February 23-26, 1992,*
667 1992. 3

668 [26] Tongzhou Mu, Zhan Ling, Fanbo Xiang, Derek Yang, Xu-
669 anlin Li, Stone Tao, Zhiao Huang, Zhiwei Jia, and Hao
670 Su. Maniskill: Generalizable manipulation skill bench-
671 mark with large-scale demonstrations. *arXiv preprint*
672 *arXiv:2107.14483*, 2021. 2

673 [27] Yao Mu, Tianxing Chen, Shijia Peng, Zanxin Chen, Zeyu
674 Gao, Yude Zou, Lunkai Lin, Zhiqiang Xie, and Ping Luo.
675 Robotwin: Dual-arm robot benchmark with generative digi-
676 tal twins (early version). In *European Conference on Com-*
677 *puter Vision*, pages 264–273. Springer, 2024. 2

678 [28] Abby O’Neill, Abdul Rehman, Abhinav Gupta, Abhiram
679 Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham
680 Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, et al.
681 Open x-embodiment: Robotic learning datasets and rt-x
682 models. *arXiv preprint arXiv:2310.08864*, 2023. 2

683 [29] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess,
684 Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and
685 Sergey Levine. Fast: Efficient action tokenization for vision-
686 language-action models. *arXiv preprint arXiv:2501.09747*,
687 2025. 6

688 [30] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez
689 Colmenarejo, Alexander Novikov, Gabriel Barth-Maron,
690 Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Sprin-
691 genberg, et al. A generalist agent. *arXiv preprint*
692 *arXiv:2205.06175*, 2022. 2

693 [31] Lucy Xiaoyang Shi, Brian Ichter, Michael Equi, Liyim-
694 ing Ke, Karl Pertsch, Quan Vuong, James Tanner, Anna
695 Walling, Haohuan Wang, Niccolo Fusai, et al. Hi robot:
696 Open-ended instruction following with hierarchical vision-
697 language-action models. *arXiv preprint arXiv:2502.19417*,
698 2025. 2

699 [32] Mustafa Shukor, Dana Aubakirova, Francesco Capuano,
700 Pepijn Kooijmans, Steven Palma, Adil Zoutine, Michel Ar-
701 actingi, Caroline Pascal, Martino Russi, Andres Marafioti,
702 et al. Smolvla: A vision-language-action model for afford-
703 able and efficient robotics. *arXiv preprint arXiv:2506.01844*,
704 2025. 6

705 [33] Sanjana Srivastava, Chengshu Li, Michael Lingelbach,
706 Roberto Martín-Martín, Fei Xia, Kent Elliott Vainio, Zheng
707 Lian, Cem Gokmen, Shyamal Buch, Karen Liu, et al. Behav-
708 ior: Benchmark for everyday household activities in virtual,
709 interactive, and ecological environments. In *Conference on*
710 *robot learning*, pages 477–490. PMLR, 2022. 2, 3

711 [34] Shuhan Tan, Kairan Dou, Yue Zhao, and Philipp
712 Krähenbühl. Interactive post-training for vision-language-
713 action models. *arXiv preprint arXiv:2505.17016*, 2025. 2

714 [35] Xin Tan, Bangwei Liu, Yicheng Bao, Qijian Tian, Zhenkun
715 Gao, Xiongbin Wu, Zhihao Luo, Sen Wang, Yuqi Zhang,
716 Xuhong Wang, et al. Towards safe and trustworthy embodied
717 ai: Foundations, status, and prospects. 2025. 2

718 [36] Gemini Robotics Team, Saminda Abeyruwan, Joshua
719 Ainslie, Jean-Baptiste Alayrac, Montserrat Gonzalez Aren-
720 as, Travis Armstrong, Ashwin Balakrishna, Robert Baruch,
721 Maria Bauza, Michiel Blokzijl, et al. Gemini robotics:
722 Bringing ai into the physical world. *arXiv preprint*
723 *arXiv:2503.20020*, 2025. 2

[37] Octo Model Team, Dibya Ghosh, Homer Walke, Karl
724 Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey
725 Hejna, Tobias Kreiman, Charles Xu, et al. Octo:
726 An open-source generalist robot policy. *arXiv preprint*
727 *arXiv:2405.12213*, 2024. 2 728

[38] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian,
729 Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-
730 world: A benchmark and evaluation for multi-task and meta
731 reinforcement learning. In *Conference on robot learning*,
732 pages 1094–1100. PMLR, 2020. 2 733

[39] Shaopeng Zhai, Qi Zhang, Tianyi Zhang, Fuxian Huang,
734 Haoran Zhang, Ming Zhou, Shengzhe Zhang, Litao Liu, Sixu
735 Lin, and Jiangmiao Pang. A vision-language-action-critic
736 model for robotic real-world reinforcement learning. *arXiv*
737 *preprint arXiv:2509.15937*, 2025. 2 738

[40] Borong Zhang, Yuhao Zhang, Jiaming Ji, Yingshan Lei,
739 Josef Dai, Yuanpei Chen, and Yaodong Yang. Safevla:
740 Towards safety alignment of vision-language-action model
741 via constrained learning. *arXiv preprint arXiv:2503.03480*,
742 2025. 2 743

[41] Shiduo Zhang, Zhe Xu, Peiju Liu, Xiaopeng Yu, Yuan Li,
744 Qinghui Gao, Zhaoye Fei, Zhangyue Yin, Zuxuan Wu, Yu-
745 Gang Jiang, et al. Vlabench: A large-scale benchmark
746 for language-conditioned robotics manipulation with long-
747 horizon reasoning tasks. *arXiv preprint arXiv:2412.18194*,
748 2024. 2 749

[42] Shiduo Zhang, Zhe Xu, Peiju Liu, Xiaopeng Yu, Yuan Li,
750 Qinghui Gao, Zhaoye Fei, Zhangyue Yin, Zuxuan Wu, Yu-
751 Gang Jiang, et al. Vlabench: A large-scale benchmark
752 for language-conditioned robotics manipulation with long-
753 horizon reasoning tasks. In *Proceedings of the IEEE/CVF*
754 *International Conference on Computer Vision*, pages 11142–
755 11152, 2025. 2 756

[43] Zijian Zhang, Kaiyuan Zheng, Zhaorun Chen, Joel Jang, Yi
757 Li, Siwei Han, Chaoqi Wang, Mingyu Ding, Dieter Fox, and
758 Huaxiu Yao. Grape: Generalizing robot policy via prefer-
759 ence alignment. *arXiv preprint arXiv:2411.19309*, 2024. 2 760

[44] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang
761 Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han,
762 Chelsea Finn, et al. Cot-vla: Visual chain-of-thought rea-
763 soning for vision-language-action models. In *Proceedings*
764 *of the Computer Vision and Pattern Recognition Conference*,
765 pages 1702–1713, 2025. 2 766

[45] Ruijie Zheng, Yongyuan Liang, Shuaiyi Huang, Jianfeng
767 Gao, Hal Daumé III, Andrey Kolobov, Furong Huang, and
768 Jianwei Yang. Tracevla: Visual trace prompting enhances
769 spatial-temporal awareness for generalist robotic policies.
770 *arXiv preprint arXiv:2412.10345*, 2024. 2 771

[46] Yifan Zhong, Fengshuo Bai, Shaofei Cai, Xuchuan Huang,
772 Zhang Chen, Xiaowei Zhang, Yuanfei Wang, Shaoyang Guo,
773 Tianrui Guan, Ka Nam Lui, et al. A survey on vision-
774 language-action models: An action tokenization perspective.
775 *arXiv preprint arXiv:2507.01925*, 2025. 2 776

[47] Yifan Zhong, Xuchuan Huang, Ruochong Li, Ceyao Zhang,
777 Zhang Chen, Tianrui Guan, Fanlian Zeng, Ka Num
778 Lui, Yuyao Ye, Yitao Liang, et al. Dexgraspvla: A
779 vision-language-action framework towards general dexter-
780 ous grasping. *arXiv preprint arXiv:2502.20900*, 2025. 2 781

- 782 [48] Xueyang Zhou, Yangming Xu, Guiyao Tie, Yongchao Chen,
783 Guowen Zhang, Duanfeng Chu, Pan Zhou, and Lichao
784 Sun. Libero-pro: Towards robust and fair evaluation of
785 vision-language-action models beyond memorization. *arXiv*
786 *preprint arXiv:2510.03827*, 2025. 2, 8
- 787 [49] Zhongyi Zhou, Yichen Zhu, Junjie Wen, Chaomin Shen, and
788 Yi Xu. Vision-language-action model with open-world em-
789 bodied reasoning from pretrained knowledge. *arXiv preprint*
790 *arXiv:2505.21906*, 2025. 2
- 791 [50] Zhongyi Zhou, Yichen Zhu, Minjie Zhu, Junjie Wen, Ning
792 Liu, Zhiyuan Xu, Weibin Meng, Yaxin Peng, Chaomin Shen,
793 Feifei Feng, et al. Chatvla: Unified multimodal understand-
794 ing and robot control with vision-language-action model. In
795 *Proceedings of the 2025 Conference on Empirical Methods*
796 *in Natural Language Processing*, pages 5377–5395, 2025. 2