

# Toward Motion Diverse Community Datasets for Scaling VLA Pretraining

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Vision-Language-Action (VLA) models have recently*  
002 *benefited from large-scale pretraining strategies inspired*  
003 *by advances in language and vision. However, unlike text*  
004 *and image domains, robotics remains constrained by lim-*  
005 *ited embodied interaction data. While community efforts*  
006 *such as LeRobot community datasets have improved ac-*  
007 *cessibility and standardization, existing datasets are often*  
008 *dominated by structurally similar, mostly pick and place*  
009 *tasks, limiting diversity in motion primitives. In this work,*  
010 *we investigate the role of motion diversity in scaling com-*  
011 *munity collected VLA pretraining data. We introduce a*  
012 *manipulation dataset, called HeteroMotion, comprising 15*  
013 *tasks across five behavior categories and 1,050 trajectories*  
014 *designed to expand action space coverage and reasoning*  
015 *complexity. Through controlled scaling HeteroMotion im-*  
016 *proves downstream real-world performance compared to*  
017 *direct fine-tuning or small-scale pretraining. A joint vari-*  
018 *ance analysis further reveals that HeteroMotion, provides*  
019 *broader motion coverage across all robot joints relative to*  
020 *existing community datasets.*

## 021 1. Introduction

022 In recent years, there has been growing interest in trans-  
023 ferring techniques across natural language processing and  
024 computer vision through the adoption of large-scale pre-  
025 trained models [16, 36]. These models, trained on vast cor-  
026 pora of text [9], images [34], image-text pairs [36], and au-  
027 dio data [20, 43], have demonstrated strong cross-domain  
028 transfer capabilities [1, 37]. The availability of large, di-  
029 verse datasets has enabled scaling laws to emerge, allow-  
030 ing data-intensive algorithms to achieve remarkable perfor-  
031 mance improvements [16, 22].

032 In contrast, robotics has not benefited from the same  
033 abundance of large-scale data. The limited availability of  
034 diverse robotic interaction data has created a significant bot-  
035 tleneck for adopting data-intensive approaches that scale  
036 well but require substantial training data [12, 35].

037 This issue has largely arisen due to the time-consuming

and costly nature of robotic data collection, which primar- 038  
ily relies on humans teleoperating expensive robots. Al- 039  
though it is possible to train a model with only an action 040  
expert [39, 45], recent work typically uses a pretrained 041  
Vision-Language Model (VLM) as the backbone of VLA 042  
systems to leverage large-scale visual and linguistic pre- 043  
training [3, 4, 6, 8, 14, 17, 19, 23, 25, 26, 29, 38]. While 044  
this approach has led to remarkable improvements, prior 045  
work has noted that robotic performance still remains fun- 046  
damentally constrained by the scale and diversity of em- 047  
bodied interaction data available during training [8, 35]. In 048  
other words, while the VLM’s perceptual abilities transfer 049  
effectively, the ability to generalize actions is still limited 050  
by the available robotic training data. 051

The open-source community has approached the chal- 052  
lenge of obtaining large-scale robotic data through two 053  
primary strategies. The first approach involves multi- 054  
institutional collaborations that collect teleoperated data 055  
across different embodiments. Prominent examples include 056  
Open X-Embodiment (OXE) [35], with more recent efforts 057  
such as DROID [24] and Bridge [15]. These datasets have 058  
powered many recent VLA models; however, they exhibit 059  
three main limitations. First, they are not highly extensi- 060  
ble. After the initial data collection phase, further additions 061  
are rare, leading to the formation of “data islands.” Sec- 062  
ond, they are often not collected through a unified frame- 063  
work, resulting in significant data cleaning, episode filter- 064  
ing, and engineering efforts before they can be effectively 065  
used. This issue is particularly evident in OXE, where prior 066  
works have selected their own subsets [25]. Although newer 067  
datasets such as DROID mitigate some of these concerns, 068  
standardization remains an important consideration for pub- 069  
lic datasets. Third, the cost of entry is high: these datasets 070  
typically rely on expensive robotic platforms (often exceed- 071  
ing \$30,000 [2]), which naturally limits broader community 072  
participation. 073

The second approach focuses on community-driven 074  
datasets, with LeRobot [11] being the most prominent ex- 075  
ample. This approach addresses many of the limitations of 076  
the first strategy. It provides a unified and extensible data 077  
collection framework, supports integration of new robots, 078

079	and remains fully open source. Additionally, it leverages af-	<b>Limitations of Existing Large-Scale VLA Datasets</b>	129
080	fordable robotic platforms (approximately \$300 [10]), sig-	There are examples of closed-source VLA models where	130
081	nificantly lowering the barrier to entry for contributors. Pre-	the training data is not public [18], preventing future work	131
082	vious work has demonstrated that models trained solely on	from leveraging these datasets as part of their pretraining	132
083	community-collected data can achieve competitive perfor-	regimen. However, there are also notable efforts aimed at	133
084	mance [38].	gathering large-scale pretraining data. Prominent examples	134
085	However, this accessibility introduces its own trade-offs.	include OXE [35], DROID [24], and Bridge [40], all of	135
086	First, affordable platforms such as SO101 have a more lim-	which have been repeatedly used in recent work to develop	136
087	ited range of motion compared to higher-end systems [2].	foundation VLA models. Nevertheless, three main chal-	137
088	Second, many community-collected datasets consist pri-	lenges remain. First, some of the data was not collected us-	138
089	marily of simple, short-horizon pick-and-place tasks. While	ing a unified platform, resulting in inconsistencies in control	139
090	such tasks are useful for adapting models to new environ-	frequency and data quality. Consequently, substantial data	140
091	ments and setups, they do not sufficiently address the fun-	engineering is required to extract high-quality pretraining	141
092	damental issue of motion diversity. As a result, scaling pre-	data. Second, this line of work is not easily extensible, as	142
093	training on structurally similar tasks may yield diminishing	it primarily relies on academic-industrial collaborations to	143
094	performance gains [41].	create large-scale datasets, leading to the formation of data	144
095	Our work aims to address this second limitation. We	islands [21]. Finally, data collection in these datasets de-	145
096	first analytically demonstrate the prevalence of similar tasks	pends on expensive robotic infrastructure, which creates a	146
097	within the existing dataset. We then introduce HeteroMo-	barrier to entry for broader community participation.	147
098	tion, a motion-diverse robotic manipulation dataset that ex-		
099	pands action-space coverage. It includes heterogeneous	<b>LeRobot Community Dataset</b> LeRobot aims to address	148
100	motion primitives and interactions beyond simple tasks. Fi-	two of these main challenges. First, it provides an extensi-	149
101	nally, through real-world experiments, we show that pre-	ble framework that supports the integration of new datasets	150
102	training on this dataset leads to measurable improvements	and robotic platforms. Second, it lowers the barrier to entry	151
103	compared to models trained without diverse pretraining.	by offering support for affordable robots such as SO100,	152
		SO101, and Koch v1.1 [10]. Previous work has demon-	153
104	<b>2. Related Work</b>	strated that pretraining on the available LeRobot dataset	154
		yields measurable performance improvements [38]. How-	155
105	Recent progress in VLA models [19, 25, 38] has been	ever, upon closer inspection of the latest version of the	156
106	closely tied to advances in large-scale robotic data collec-	dataset, we find that out of 791 total tasks, approximately	157
107	tion [7, 8]. Existing approaches primarily differ in how data	500 correspond to variations of pick-and-place behaviors	158
108	is gathered, standardized, and shared across institutions and	that are highly similar in their underlying motion primi-	159
109	embodiments [24, 35, 40]. While these efforts have signif-	tives. Our work explicitly aims to address these issues and	160
110	icantly improved policy learning, they introduce trade-offs	propose a path toward higher-quality, community-collected	161
111	between accessibility, diversity, and extensibility. In this	VLA pretraining data with increased motion diversity.	162
112	section, we position our work within these data paradigms		
113	filling the remaining gaps in action-space diversity.	<b>3. Methodology</b>	163
114	<b>Data Requirements in VLA Training</b> In the recent de-	Our work directly addresses the limitations of motion ho-	164
115	velopment of VLA models, the general training recipe has	mogeneity in existing community datasets. We introduce a	165
116	followed the pretraining/post-training paradigm established	structured collection of 15 tasks spanning five behavior cat-	166
117	in large language and vision models [9, 16, 36]. As a result,	egories and, more importantly, through a controlled quan-	167
118	two categories of data are required for this type of train-	titative scaling study, demonstrate that increasing struc-	168
119	ing. The first is large-scale pretraining data, which usu-	tured motion diversity in community-collected data leads	169
120	ally consists of multiple embodiments performing diverse	to measurable improvements in real-world downstream per-	170
121	tasks across different environments [24, 35]. This helps	formance across diverse manipulation settings.	171
122	the model generalize, recover from mistakes, and adapt to		
123	different embodiments [7]. The second category is high-	<b>3.1. Dataset</b>	172
124	quality, task-specific fine-tuning data. Fine-tuning data is	Figure 1 illustrates HeteroMotion, which consists of 15	173
125	typically embodiment and task specific; therefore, much of	tasks grouped into five behavior categories, with three tasks	174
126	the community effort has focused on creating large pretrain-	per category and 50 demonstrations per task.	175
127	ing datasets, or incorporating fine-tuning data into larger	To increase motion diversity and reasoning complexity	176
128	community datasets.	beyond repetitive short-horizon pick-and-place, we design	177
		the following categories:	178

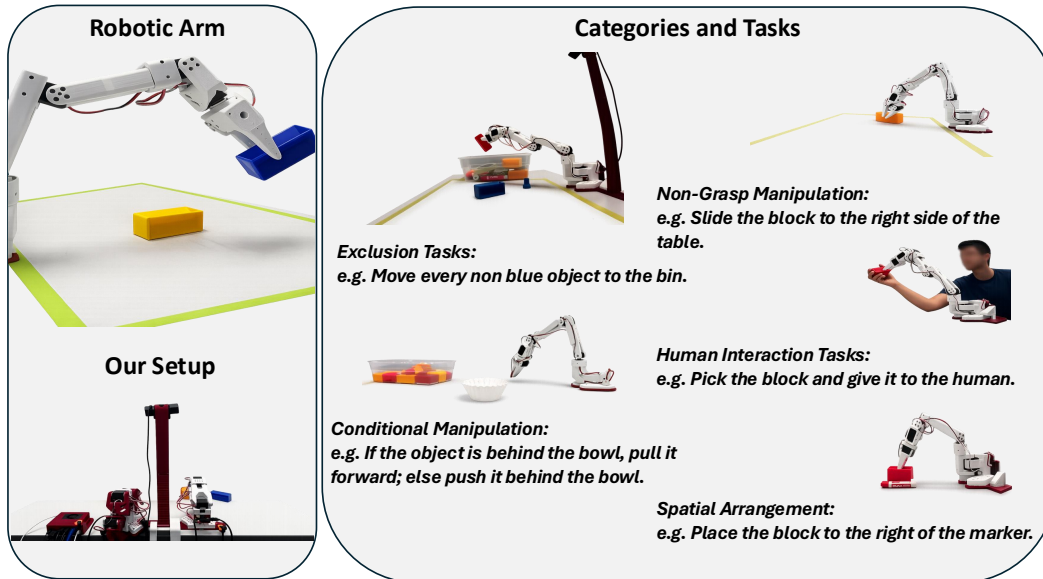


Figure 1. Overview of the proposed dataset, HeteroMotion, highlighting the robotic platform and five task categories designed to promote diverse manipulation skills and complex reasoning.

- 179 1. **Non-Grasp Manipulation:** Focuses on actions such as  
180 pushing or sliding that do not involve gripper actuation.
- 181 2. **Conditional Manipulation:** Tasks that require reason-  
182 ing over conditional instructions prior to execution.
- 183 3. **Human Interaction:** Includes direct human–robot in-  
184 teraction scenarios, such as object handover. The robot  
185 must account for human presence and coordinate its mo-  
186 tion accordingly.
- 187 4. **Spatial Arrangement:** Involves arranging objects ac-  
188 cording to specified spatial relationships (e.g., left of,  
189 behind, inside).
- 190 5. **Exclusion Tasks:** This category requires the robot to  
191 deliberately ignore specified objects while manipulating  
192 others.

193 A complete breakdown of all 15 tasks and their descriptions  
194 is provided in [Appendix A](#) (Tab. 2).

195 These five categories form the primary pretraining  
196 dataset. Additionally, we provide two supplementary  
197 datasets: (1) 150 autonomously generated, human labeled  
198 trajectories to support reinforcement learning methods such  
199 as HIL-SERL [31], and (2) 150 simple manipulation tra-  
200 jectories from downstream tasks introduced in experiments.  
201 HeteroMotion contains approximately 1,050 manipulation  
202 trajectories for VLA pretraining.

### 203 3.2. Training

204 To validate the effectiveness of HeteroMotion in improving  
205 downstream VLA performance and to study scaling effects,  
206 we trained five model variants fine-tuned on downstream  
207 tasks. These include the original SmolVLA [38] model and

four additional variants of the same model pretrained on  
10%, 25%, 50%, and 100% of HeteroMotion.

We ensured that the smaller pretraining subsets were  
sampled uniformly and maintained task diversity. Each  
variant was then fine-tuned for three downstream tasks: (1)  
pick-and-place, (2) sorting, and (3) stacking. All mod-  
els were trained for a proportional number of steps, cor-  
responding to an equal number of epochs, followed by  
20k steps for task adaptation. We used the LeRobot train-  
ing framework [11], initializing from *lerobot/smolvla\_base*.  
During training, only the action expert was trained while the  
backbone VLM remained frozen.

The action expert models the distribution  $p(\mathbf{A}_t|\mathbf{o}_t)$   
where  $\mathbf{A}_t$  corresponds to an action chunk [45] for  $H$  actions  
and hence  $\mathbf{A}_t = [\mathbf{a}_t, \mathbf{a}_{t+1}, \dots, \mathbf{a}_{t+H-1}]$ . Our observation  
on the other hand is  $\mathbf{o}_t = [\mathbf{I}_1, \mathbf{I}_2, l_t, \mathbf{q}_t]$  where  $l_t$  correspond  
to the language prompt,  $\mathbf{q}_t$  to the robot’s joints and  $\mathbf{I}_1$  and  
 $\mathbf{I}_2$  correspond to top and side camera respectively. We used  
the following flow matching loss, introduced in [27, 30] and  
later adopted in [5]:

$$L^\tau(\theta) = \mathbb{E}_{p(\mathbf{A}_t|\mathbf{o}_t), q(\mathbf{A}_t^\tau|\mathbf{A}_t)} \|\mathbf{v}_\theta(\mathbf{A}_t^\tau, \mathbf{o}_t) - \mathbf{u}(\mathbf{A}_t^\tau|\mathbf{A}_t)\|^2 \quad 228$$

where the subscripts correspond to trajectory timesteps  
( $t$ ) and the superscripts correspond to the flow-matching  
timesteps ( $\tau$ ). The network works by sampling an  $\epsilon \sim$   
 $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , computing the noisy actions  $\mathbf{A}_t^\tau = \tau\mathbf{A}_t + (1-\tau)\epsilon$ ,  
and then training  $\mathbf{v}_\theta(\mathbf{A}_t^\tau, \mathbf{o}_t)$  to match the denoising vector  
field  $\mathbf{u}(\mathbf{A}_t^\tau|\mathbf{A}_t) = \mathbf{A}_t - \epsilon$ .

All experiments were conducted on a single A100 GPU.  
Training required approximately one hour per 5k steps, to-

237 taling roughly 100 GPU hours for the project.

## 238 4. Evaluation and Results

239 All experiments were performed on the SO101 robotic arm  
240 platform. The system utilized two RGB cameras: a top-  
241 down camera for global scene observation and a side-view  
242 camera to capture interaction dynamics. Both cameras op-  
243 erated at a resolution of  $640 \times 480$  pixels. This dual-view  
244 configuration provides complementary spatial perspectives,  
245 enhancing manipulation robustness and mitigating errors  
246 caused by occlusions. Experiments were conducted on  
247 a Jetson Orin Nano (8GB) [33], an Ampere-based edge  
248 AI device with CUDA and Tensor Core acceleration, en-  
249 abling GPU-accelerated inference in a compact and energy-  
250 efficient form factor suitable for robotics applications.

### 251 4.1. Experiment Design and Evaluation Protocol

252 We evaluate our method on three commonly used down-  
253 stream manipulation tasks, following [28, 38]: pick-and-  
254 place, sorting, and stacking. These tasks are particularly  
255 suitable because they are prevalent in SmolVLA’s pretrain-  
256 ing data. Improvements on them therefore reflect stronger  
257 underlying learning, rather than gains from introducing en-  
258 tirely new tasks absent from pretraining.

259 Performance was assessed using task-specific scoring  
260 schemes designed to capture partial and complete task exe-  
261 cution. For the pick-and-place tasks, we employed a three-  
262 level scoring system 0, 0.5, 1, corresponding respectively to  
263 task failure, successful grasping of the cube without place-  
264 ment, and successful grasping followed by correct place-  
265 ment into the designated cup. The stacking tasks, defined  
266 as grasping one cube and placing it on top of another, were  
267 evaluated using the same discrete scale: 0 indicates failure  
268 to grasp, 0.5 indicates successful grasp without successful  
269 stacking, and 1 denotes complete execution of the stacking  
270 behavior. The sorting tasks, which require placing cubes  
271 of different colors into their assigned cups, were evaluated  
272 using a four-level scale 0.25, 0.5, 0.75, 1 to capture progres-  
273 sive task completion: grasping the first cube (0.25), correct  
274 placement of the first cube (0.5), correct placement of the  
275 first cube followed by grasping the second cube (0.75), and  
276 successful placement of both cubes (1). Each task was ex-  
277 ecuted twice per configuration. Final task scores were ag-  
278 gregated across 10 repetitions and normalized by a factor of  
279 ten to yield the reported success rates.

### 280 4.2. Real World Experimental Results

281 Table 1 summarizes the real-world in-distribution perfor-  
282 mance across our downstream tasks. The base model cor-  
283 responds to SmolVLA fine-tuned on each downstream task  
284 without any pretraining on our proposed dataset. In con-  
285 trast, the remaining variants were first pretrained on 10%,  
286 25%, 50%, or 100% of HeteroMotion and then fine-tuned

Table 1. Results of real-world robotic manipulation experiments. The percentages denote the fraction of fine-tuning data employed.

Model Configuration	Task Success Rate (%)			
	Pick and Place	Stacking	Sorting	Average
SmolVLA (0% Pre-trained)	30.0	20.0	2.5	17.50
SmolVLA (10% Pre-trained)	35.0	15.0	0.0	16.67
SmolVLA (25% Pre-trained)	45.0	2.0	5.0	17.33
SmolVLA (50% Pre-trained)	<b>50.0</b>	25.0	12.5	29.17
SmolVLA (100% Pre-trained)	45.0	<b>30.0</b>	<b>15.0</b>	<b>30.00</b>

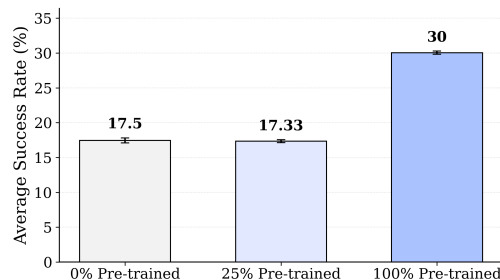


Figure 2. Average success rate across down-stream tasks. Error bars denote 95% confidence intervals.

287 on the downstream tasks under identical training conditions.  
288 Figure 2 presents the average task success rate across the  
289 three representative configurations with 95% confidence in-  
290 tervals. The results show that pretraining on HeteroMotion  
291 yields a clear improvement over direct fine-tuning, with the  
292 100% pretrained model achieving a substantial performance  
293 gain compared to both the 0% and 25% variants.

294 Notably, smaller pretraining subsets (e.g., 10% and 25%)  
295 provide limited or inconsistent improvements over the base  
296 model, which is in line with previous findings that pretrain-  
297 ing must surpass a certain threshold before yielding notice-  
298 able performance gains [42]. The 50% and especially the  
299 100% pretraining variants show more pronounced improve-  
300 ments, indicating that meaningful benefits arise when suffi-  
301 cient data is incorporated.

302 The observed performance gains can be attributed to the  
303 exposure of the model during pretraining to a broader range  
304 of motions and task complexities. This broader exposure  
305 increases coverage of target distribution, rather than simply  
306 increasing sample density within a restricted subset of the  
307 action space. Domain adaptation results show that the ex-  
308 pected target error  $\epsilon_P(h)$  can be bounded as:

$$309 \epsilon_P(h) \leq \epsilon_Q(h) + \text{disc}(Q, P) + \lambda$$

310 where  $\epsilon_Q(h)$  denotes the source (pretraining) error,  
311  $\text{disc}(Q, P)$  measures the divergence between the source and  
312 target distributions, and  $\lambda$  captures the irreducible joint er-  
313 ror [13, 32]. Importantly, reducing  $\epsilon_Q(h)$  by scaling data  
314 within a narrow region does not reduce the divergence term

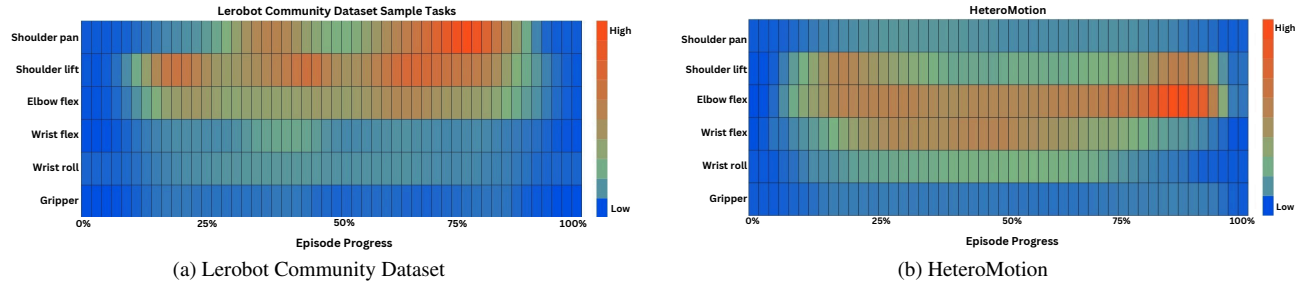


Figure 3. Comparison between the variance of Lerobot Community dataset and our proposed dataset, HeteroMotion.

disc(Q, P) if substantial portions of the target support remain uncovered. In contrast, expanding pretraining coverage reduces this distributional discrepancy. Our results are therefore consistent with improved support coverage during pretraining, rather than gains arising solely from increased data volume within a limited portion of the action space.

### 4.3. Joint Variance Analysis

To demonstrate the presence of motion diversity in HeteroMotion, we analyzed the variance across the six joints of the SO101. We compared samples from HeteroMotion to a merged subset of three different datasets from the LeRobot Community dataset, ensuring that they represented distinct tasks. This variance analysis is important, as prior work has suggested that low-variance regions may lack the exploratory coverage necessary for model generalization [44].

As shown in Figure 3a, the existing community dataset exhibits reasonable variance along the first three joints. This aligns with our observation of the task distribution: most pick-and-place tasks involve lifting the arm (primarily moving the shoulder lift and elbow flex joints), rotating toward the object (shoulder pan), and then moving toward a bin to place the item. The prevalence of pick-and-place behaviors is reflected in Figure 3a, where the high-variance regions showcase this movement. For example, the shoulder pan joint shows two distinct high-variance regions: one during rotation toward the object and another during rotation back toward the bin, with relatively constant motion in between.

In contrast, HeteroMotion dataset, shown in Figure 3b, exhibits sustained and more evenly distributed variance throughout the entire episode. Notably, it captures significantly greater diversity in wrist flex and wrist roll, while also maintaining variability across all joints. This broader distribution of motion better reflects diverse manipulation skills and supports our claim that increasing motion diversity is critical for scaling community-collected VLA pretraining data.

## 5. Conclusion

In this work, we investigated the motion diversity present in the LeRobot community dataset. Through a systematic

analysis of task descriptions, we observed a strong concentration of pick-and-place tasks, suggesting limited diversity in motion patterns. We further examined this through a joint variance analysis, which further suggested that HeteroMotion dataset is dominated by a narrow set of motion primitives. To address this limitation, we collected a structured dataset spanning five behavior categories, comprising 15 tasks and over 1,050 trajectories. To evaluate the impact of increased motion diversity, we conducted real-world experiments using the baseline SmolVLA model alongside variants pretrained on different portions of our dataset. The results demonstrate clear performance improvements when larger scale pretraining is present.

Overall, our findings suggest that community-collected datasets would benefit from more intentional guidance toward complex and physically diverse tasks. Rather than focusing primarily on scaling the number of similar demonstrations, future data collection efforts should prioritize tasks that expand motion patterns, interaction complexity, and reasoning requirements. Encouraging contributors to record more compositionally and physically varied behaviors may provide a more effective path toward scalable and robust VLA models.

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022. 1
- [2] Jorge Aldaco, Travis Armstrong, Robert Baruch, Jeff Bingham, Sanky Chan, Kenneth Draper, Debidatta Dwibedi, Chelsea Finn, Pete Florence, Spencer Goodrich, et al. Aloha 2: An enhanced low-cost hardware for bimanual teleoperation. *arXiv preprint arXiv:2405.02292*, 2024. 1, 2
- [3] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 1
- [4] Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox,

- 395 Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open  
396 foundation model for generalist humanoid robots. *arXiv*  
397 *preprint arXiv:2503.14734*, 2025. 1
- 398 [5] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail,  
399 Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom,  
400 Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim  
401 Jones, Liyiming Ke, Sergey Levine, Adrian Li-Bell, Mohith  
402 Mothukuri, Suraj Nair, Karl Pertsch, Lucy Xiaoyang Shi,  
403 James Tanner, Quan Vuong, Anna Walling, Haohuan Wang,  
404 and Ury Zhilinsky.  $\pi_0$ : A vision-language-action flow model  
405 for general robot control, 2026. 3
- 406 [6] Florian Bordes, Richard Yuanzhe Pang, Anurag Ajay,  
407 Alexander C Li, Adrien Bardes, Suzanne Petryk, Oscar  
408 Mañas, Zhiqiu Lin, Anas Mahmoud, Bargav Jayaraman,  
409 et al. An introduction to vision-language modeling. *arXiv*  
410 *preprint arXiv:2405.17247*, 2024. 1
- 411 [7] Konstantinos Bousmalis, Giulia Vezzani, Dushyant Rao, Co-  
412 line Devin, Alex X Lee, Maria Bauzá, Todor Davchev, Yuxi-  
413 ang Zhou, Agrim Gupta, Akhil Raju, et al. Robocat: A self-  
414 improving generalist agent for robotic manipulation. *arXiv*  
415 *preprint arXiv:2306.11706*, 2023. 2
- 416 [8] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen  
417 Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakr-  
418 ishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al.  
419 Rt-1: Robotics transformer for real-world control at scale.  
420 *arXiv preprint arXiv:2212.06817*, 2022. 1, 2
- 421 [9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Sub-  
422 biah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakan-  
423 tan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Lan-  
424 guage models are few-shot learners. *Advances in neural in-*  
425 *formation processing systems*, 33:1877–1901, 2020. 1, 2
- 426 [10] Remi Cadene, Simon Alibert, Francesco Capuano, Michel  
427 Aractingi, Adil Zouitine, Pepijn Kooijmans, Jade Choghari,  
428 Martino Russi, Caroline Pascal, Steven Palma, et al. Lerobot:  
429 An open-source library for end-to-end robot learning. In *The*  
430 *Fourteenth International Conference on Learning Representa-*  
431 *tions*. 2
- 432 [11] Remi Cadene, Simon Alibert, Alexander Soare, Quentin  
433 Gallouedec, Adil Zouitine, Steven Palma, Pepijn Kooij-  
434 mans, Michel Aractingi, Mustafa Shukor, Dana Aubakirova,  
435 Martino Russi, Francesco Capuano, Caroline Pascal, Jade  
436 Choghari, Jess Moss, and Thomas Wolf. Lerobot: State-of-  
437 the-art machine learning for real-world robotics in pytorch.  
438 <https://github.com/huggingface/lerobot>,  
439 2024. 1, 3
- 440 [12] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair,  
441 Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh,  
442 Sergey Levine, and Chelsea Finn. Robonet: Large-scale  
443 multi-robot learning. *arXiv preprint arXiv:1910.11215*,  
444 2019. 1
- 445 [13] Shai Ben David, Tyler Lu, Teresa Luu, and Dávid Pál. Im-  
446 possibility theorems for domain adaptation. In *Proceedings*  
447 *of the Thirteenth International Conference on Artificial In-*  
448 *telligence and Statistics*, pages 129–136. JMLR Workshop  
449 and Conference Proceedings, 2010. 4
- 450 [14] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch,  
451 Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid,  
Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-  
e: An embodied multimodal language model. *arXiv preprint*  
*arXiv:2303.03378*, 2023. 1
- [15] Frederik Ebert, Yanlai Yang, Karl Schmeckpeper, Bernadette  
Bucher, Georgios Georgakis, Kostas Daniilidis, Chelsea  
Finn, and Sergey Levine. Bridge data: Boosting general-  
ization of robotic skills with cross-domain datasets. *arXiv*  
*preprint arXiv:2109.13396*, 2021. 1
- [16] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch,  
Elena Buchatskaya, Trevor Cai, Eliza Rutherford, DDL  
Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark,  
et al. Training compute-optimal large language models.  
*arXiv preprint arXiv:2203.15556*, 10, 2022. 1, 2
- [17] Chi-Pin Huang, Yunze Man, Zhiding Yu, Min-Hung Chen,  
Jan Kautz, Yu-Chiang Frank Wang, and Fu-En Yang. Fast-  
thinkact: Efficient vision-language-action reasoning via ver-  
balizable latent planning. *arXiv preprint arXiv:2601.09708*,  
2026. 1
- [18] Physical Intelligence, Ali Amin, Raichelle Aniceto, Ash-  
win Balakrishna, Kevin Black, Ken Conley, Grace Con-  
nors, James Darpinian, Karan Dhabalia, Jared DiCarlo, et al.  
 $\pi_{0.6}$ : a vla that learns from experience. *arXiv preprint*  
*arXiv:2511.14759*, 2025. 2
- [19] Physical Intelligence, Kevin Black, Noah Brown, James  
Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail,  
Michael Equi, Chelsea Finn, Niccolo Fusai, et al.  $\pi_{0.5}$ : a  
vision-language-action model with open-world generaliza-  
tion. *arXiv preprint arXiv:2504.16054*, 2025. 1, 2
- [20] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh,  
Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom  
Duerig. Scaling up visual and vision-language representa-  
tion learning with noisy text supervision. In *International*  
*conference on machine learning*, pages 4904–4916. PMLR,  
2021. 1
- [21] Peter Kairouz and H Brendan McMahan. Advances and open  
problems in federated learning. *Foundations and trends in*  
*machine learning*, 14(1-2):1–210, 2021. 2
- [22] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B  
Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec  
Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for  
neural language models. *arXiv preprint arXiv:2001.08361*,  
2020. 1
- [23] Kento Kawaharazuka, Jihoon Oh, Jun Yamada, Ingmar Pos-  
ner, and Yuke Zhu. Vision-language-action models for  
robotics: A review towards real-world applications. *IEEE*  
*Access*, 2025. 1
- [24] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ash-  
win Balakrishna, Sudeep Dasari, Siddharth Karamcheti,  
Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yun-  
liang Chen, Kirsty Ellis, et al. Droid: A large-scale  
in-the-wild robot manipulation dataset. *arXiv preprint*  
*arXiv:2403.12945*, 2024. 1, 2
- [25] 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. 1, 2

- 509 [26] Jason Lee, Jiafei Duan, Haoquan Fang, Yuquan Deng, Shuo  
510 Liu, Boyang Li, Bohan Fang, Jieyu Zhang, Yi Ru Wang,  
511 Sangho Lee, et al. Molmoact: Action reasoning models that  
512 can reason in space. *arXiv preprint arXiv:2508.07917*, 2025.  
513 1
- 514 [27] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian  
515 Nickel, and Matt Le. Flow matching for generative modeling.  
516 *arXiv preprint arXiv:2210.02747*, 2022. 3
- 517 [28] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu,  
518 Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge  
519 transfer for lifelong robot learning, 2023. 4
- 520 [29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee.  
521 Visual instruction tuning. *Advances in neural information  
522 processing systems*, 36:34892–34916, 2023. 1
- 523 [30] Qiang Liu. Rectified flow: A marginal preserving approach  
524 to optimal transport. *arXiv preprint arXiv:2209.14577*, 2022.  
525 3
- 526 [31] Jianlan Luo, Charles Xu, Jeffrey Wu, and Sergey Levine.  
527 Precise and dexterous robotic manipulation via human-in-  
528 the-loop reinforcement learning. *Science Robotics*, 10(105):  
529 eads5033, 2025. 3
- 530 [32] Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh.  
531 Domain adaptation: Learning bounds and algorithms. *arXiv  
532 preprint arXiv:0902.3430*, 2009. 4
- 533 [33] NVIDIA Corporation. Nvidia jetson orin nano super devel-  
534 oper kit, 2025. Accessed: 2026-02. 4
- 535 [34] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy  
536 Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez,  
537 Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al.  
538 Dinov2: Learning robust visual features without supervision.  
539 *arXiv preprint arXiv:2304.07193*, 2023. 1
- 540 [35] Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek  
541 Gupta, Abhishek Padalkar, Abraham Lee, Acorn Poo-  
542 ley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open  
543 x-embodiment: Robotic learning datasets and rt-x models:  
544 Open x-embodiment collaboration 0. In *2024 IEEE Inter-  
545 national Conference on Robotics and Automation (ICRA)*,  
546 pages 6892–6903. IEEE, 2024. 1, 2
- 547 [36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya  
548 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,  
549 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning  
550 transferable visual models from natural language supervi-  
551 sion. In *International conference on machine learning*, pages  
552 8748–8763. PmLR, 2021. 1, 2
- 553 [37] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez  
554 Colmenarejo, Alexander Novikov, Gabriel Barth-Maron,  
555 Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Spring-  
556 enberg, et al. A generalist agent. *arXiv preprint  
557 arXiv:2205.06175*, 2022. 1
- 558 [38] Mustafa Shukor, Dana Aubakirova, Francesco Capuano,  
559 Pepijn Kooijmans, Steven Palma, Adil Zouitine, Michel Ar-  
560 actingi, Caroline Pascal, Martino Russi, Andres Marafioti,  
561 et al. Smolvla: A vision-language-action model for afford-  
562 able and efficient robotics. *arXiv preprint arXiv:2506.01844*,  
563 2025. 1, 2, 3, 4
- 564 [39] Octo Model Team, Dibya Ghosh, Homer Walke, Karl  
565 Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey  
Hejna, Tobias Kreiman, Charles Xu, et al. Octo:  
An open-source generalist robot policy. *arXiv preprint  
arXiv:2405.12213*, 2024. 1
- [40] Homer Rich Walke, Kevin Black, Tony Z Zhao, Quan  
Vuong, Chongyi Zheng, Philippe Hansen-Estruch, An-  
drew Wang He, Vivek Myers, Moo Jin Kim, Max Du, et al.  
Bridgedata v2: A dataset for robot learning at scale. In *Con-  
ference on Robot Learning*, pages 1723–1736. PMLR, 2023.  
2
- [41] Yixiao Wang. Generalization capability for imitation learn-  
ing. *arXiv preprint arXiv:2504.18538*, 2025. 2
- [42] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret  
Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma,  
Denny Zhou, Donald Metzler, et al. Emergent abilities of  
large language models. *arXiv preprint arXiv:2206.07682*,  
2022. 4
- [43] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor  
Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale con-  
trastive language-audio pretraining with feature fusion and  
keyword-to-caption augmentation. In *ICASSP 2023-2023  
IEEE International Conference on Acoustics, Speech and  
Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 1
- [44] Thomas T. Zhang, Daniel Pfrommer, Chaoyi Pan, Niko-  
lai Matni, and Max Simchowitz. Action chunking and ex-  
ploratory data collection yield exponential improvements in  
behavior cloning for continuous control, 2025. 5
- [45] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea  
Finn. Learning fine-grained bimanual manipulation with  
low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023.  
1, 3

596

**Appendix**

597

**A. Additional Dataset Details**

Table 2. Overview of the 15 manipulation tasks by category.

Category	Task	Description
<b>Conditional Manipulation</b>	• If the block is near the sides, move it inward; else move it toward the sides.	• Requires conditional decision-making based on spatial proximity to table boundaries.
	• If the object is closer to the bin than the box, move it to the box; else move it to the bin.	• Involves comparative spatial reasoning between two target containers.
	• If the object is behind the bowl, pull it forward; else push it behind the bowl	• Tests relational spatial understanding based on a reference object.
<b>Spatial Arrangement</b>	• Place the marker between two blocks	• Evaluates precise spatial placement within a constrained configuration.
	• Arrange blocks in a straight horizontal line	• Requires structured alignment across multiple objects.
	• Place the block to the right of the marker	• Assesses directional spatial reasoning relative to a fixed reference.
<b>Non-Grasp Manipulation</b>	• Slide the block to the right side of the table	• Tests non-prehensile manipulation without lifting or grasping.
	• Push the block into the marked region	• Requires accurate planar control to move the object into a predefined region.
	• Push the block until it touches the box	• Evaluates contact-aware motion execution and physical interaction.
<b>Exclusion Tasks</b>	• Put all blocks except the red one into the bin	• Tests selective reasoning by filtering objects by attribute.
	• Move everything except the block to	• Requires identifying and preserving a designated object while relocating others.
	• Move every non blue object to the bin	• Evaluates attribute-based filtering and selective object transfer.
<b>Human Interaction Tasks</b>	• Pick the block and give it to the human	• Assesses coordinated grasping and safe object handover.
	• Take the block from the human hand and put it into the bin	• Requires bidirectional human-robot interaction and placement.
	• Pick the cracker and give it to the human	• Evaluates safe manipulation of fragile objects during interaction.

598