

STEER: Bridging VLMs and Low-Level Control for Adaptable Robotic Manipulation

Anonymous Author(s)

Affiliation

Address

email

Abstract: Recent advances have showcased the opportunity of leveraging the broad semantic understanding learned by vision-language models (VLMs) in robot learning; however, effectively connecting VLMs to robot control remains challenging due to the scarcity of physical robot data compared to internet-scale training data. We propose STEER, a system that learns flexible, low-level manipulation skills, allowing for modulation and adaptation to new situations. By training low-level policies on structured, dense re-annotations of existing robot datasets, we create an intuitive interface for humans or VLMs to guide robots in unfamiliar scenarios and perform new tasks using common-sense reasoning. Our results demonstrate that skills learned through STEER can be synthesized to accomplish held-out tasks without additional training. (Videos¹)

1 Introduction

Designing robots that can handle diverse and nuanced tasks posed by the real world is challenging, as it requires adaptability to complex, dynamic environments. Imitation learning (IL) is a widely-used, data-driven approach that distills expert demonstrations into learned policies, enabling precise manipulation of high-dimensional robot systems in real-world environments [1, 2] and at scale [3, 4, 5, 6, 7, 8, 9]. Despite these advances, robot systems trained with IL remain largely limited to scenarios encountered during training, which are fundamentally narrow as collecting real-world embodied data is costly and constrained by physical limitations.

Humans, on the other hand, can adapt to complex, unfamiliar situations with ease, thanks to "common sense" generalization. Humans effortlessly understand high-level concepts like object affordances, intuitive physics, and compositionality—referred to as 'System 2' processing [10], which involves deliberate, analytical thinking. This contrasts with 'System 1' behaviors, which are reactive and particularly useful in contact-rich manipulation. Natural language serves as a key medium for studying System 2 reasoning, not only as a means through which humans understand and describe the world, but also as the primary input-output modality for vision-language models (VLMs) that demonstrate complex, human-like reasoning capabilities [11, 12, 13]. Connecting high-level System 2 plans to low-level System 1 behaviors, however, is not straightforward. Several methods have been proposed to bridge this gap: some enable System 2 reasoning systems to operate in modalities more easily transferable to robotic policies, such as code or semantic keypoints [14, 15], while others account for the lack of System 2 physical grounding by considering robot affordances during planning [16] or jointly training on both internet and embodied data [6, 5]. These approaches typically view System 1 processing as inflexible, seeking to improve System 2 reasoning outputs to better control a fixed System 1 policy. Instead of augmenting System 2 reasoning outputs, we ask: can we improve the System 1 policy to be more flexible and steerable by System 2 processes? Can this combination enable generalizable, end-to-end control?

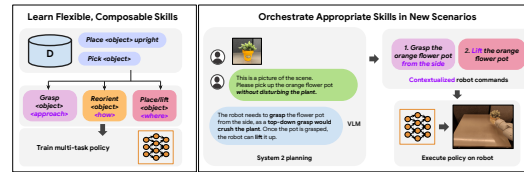


Figure 1: In STEER, we train on a dataset of diverse robot behaviors that is re-annotated to describe the primitive skills used to manipulate objects, with a focus on *how* the robot performed that skill. At inference time, a high-level system (VLM or human) receives a complex language instruction and determines the low-level skills to employ in the given context. Focusing on the “how” enables better contextual behavior.

¹<https://steer-anon.github.io>

43 We introduce STEER: Structured Training for EmbodiEd Reasoning, a framework for training low-
 44 level reactive policies that can be flexibly guided by higher-level reasoning systems like humans or
 45 VLMs. The key insight of STEER is the use of dense language annotations on robot data, allowing
 46 us to train conditional policies based on detailed language instructions. These policies can be con-
 47 ditioned on each step of a plan generated by a high-level model (e.g., VLM or LLM), effectively
 48 combining System 1 and System 2 capabilities. This enables the robot to adapt to new situations by
 49 synthesizing behaviors not explicitly demonstrated during training. We implement this system using
 50 real-world datasets and propose an automated labeling pipeline based on proprioceptive observa-
 51 tions to extract basic, object-centric manipulation skills, which are distilled into a low-level policy.
 52 Additionally, we present a strategy for using a VLM to generate language-based instructions for this
 53 low-level policy. Crucially, our approach enables the repurposing of robot skills in a semantically
 54 meaningful way at test time, allowing robots to autonomously handle novel situations.

55 2 Related Work

56 Imitation learning (IL) has become the dominant paradigm for training robotic manipulation poli-
 57 cies [3, 2, 17]. However, deploying these models in unstructured environments remains challeng-
 58 ing, as robot policies trained on human-collected data struggle in "out-of-distribution" scenarios
 59 where demonstrations are sparse [18]. This limitation arises from the high cost of collecting large-
 60 scale robot data compared to web-scale datasets used for training foundation models [19, 11, 13].
 61 To improve generalization, researchers have leveraged text and vision foundation models to uti-
 62 lize existing datasets. This includes enhancing IL policies for open-world object grasping through
 63 open-vocabulary object detection [20] and relabeling episode-level instructions with models like
 64 CLIP [21, 22, 23]. Our work aligns with these dataset relabeling approaches, aiming to expand
 65 robot capabilities by relabeling behavior modes in existing heterogeneous demonstration datasets.
 66 Previous research has also investigated expressive modalities for policy conditioning, such as goal
 67 target poses [24], images [25, 26], trajectories [27, 28], and code [14]. However, natural language
 68 remains the primary modality for complex planning in state-of-the-art LLMs and VLMs, motivat-
 69 ing STEER to improve language-conditioned action prediction. Additionally, many works explore
 70 learned skills to accelerate new task learning with temporally extended, semantically meaningful ac-
 71 tion sequences [29, 30, 31, 32, 33]. These approaches often employ hierarchical policies that learn
 72 to compose skills through RL [34, 35, 36]. While EXTRACT [33] uses VLMs to label skills for
 73 new tasks, our method leverages VLMs' common-sense reasoning to select appropriate skills with-
 74 out training a separate policy. By reasoning about how humans approach tasks from visual inputs,
 75 our framework enables robots to plan longer-horizon tasks and manage novel object configurations.
 76 This contrasts with prior work on affordances, which typically relies on keypoint representations
 77 in pixel space [37, 38]. Our approach instead reasons about affordances through natural language,
 78 allowing for more nuanced interactions with off-the-shelf VLMs or human operators.

79 3 System Design

80 Our goal is to extract language-indexed, object-centric skills that
 81 facilitate task execution via foundation models. We achieve this
 82 by annotating existing datasets and training a language-conditioned
 83 RT-1 policy [4] using segmented and relabeled instructions. We
 84 extract semantically identifiable categories linked to language de-
 85 scriptions, focusing on shared, object-relational skills like grasping,
 86 lifting, placing, and rotating, originally demonstrated through tem-
 87 plates like pick <object>, move <object1> near <object2>, knock <object>, place
 88 <object> upright, can be executed with varying strategies. Key factors include:

89 *Grasp Angle*. Objects can be grasped in multiple stable positions, and the particular way indeed
 90 impacts the ability to perform downstream tasks. However, grasp positions are rarely labeled apriori,
 91 as they are often implicit. We use a simple approach to label the grasp approach by manually labeling
 92 a relatively small set of 'anchor' grasp poses. We then label an arbitrary grasp with the label of its
 93 nearest neighbor 'anchor' pose as measured by cosine similarity. We represent a grasp pose as a 3D

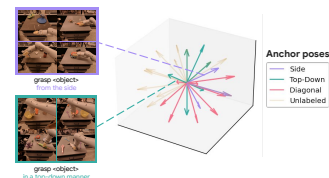


Figure 2: Anchor vectors and their semantic labels. Purple, green, and pink vectors represent side, top-down, and diagonal.

94 unit vector, and we identify the time of a grasp where the gripper changed from fully open to fully
95 closed. To define and label the anchor poses, we took 3D unit vectors that are linear combinations
96 of the elementary 3D basis vectors and visually inspected clusters in order to label them. In the
97 grasp data, we identify three distinct modes via inspecting the grasp images: top-down grasps, side
98 grasps, and diagonal grasps (visualized in Figure 2). The sub-trajectory is relabeled to grasp the
99 `<object>` in a `<grasp approach>`, where `<object>` is from the original instruction and `<grasp`
100 `approach>` is from the anchor’s label.

101 *Reorientation.* Another mode of behavior identified in the dataset is of reorienting objects. In
102 order to identify and label these reorientations, we first label the wrist orientation for every timestep
103 where the gripper is fully closed. Then if the gripper orientation switches between two of the modes
104 (as labeled in *Grasp Angle*), we label the sub-trajectory preceding it as `reorient the <object>`
105 `<direction>`, where `<object>` again is from the original instruction and `<direction>` indicates
106 whether the object is rotated from upright to horizontal or vice versa.

107 *Lifting/Placing.* Complementing grasping, we label whether the object was lifted or placed at the
108 end of completing the original task. If the object is still held at the end of the episode and the gripper
109 moves vertically upward, we label the final sub-trajectory as `hold and lift the <object>`. If
110 not, similar to identifying grasps, we identify the time of placing using the gripper state and label
111 this sub-trajectory as `place the <object>`.

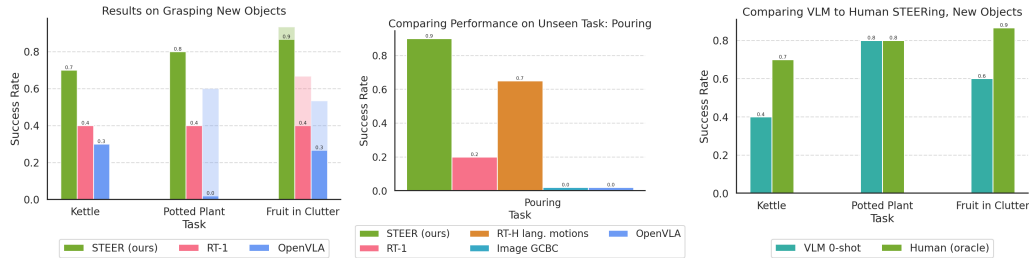
112 **Orchestrating Learned Skills** A key capability of the System 2 component is being able to reason
113 about the visual observation of the scene, the task description, and the robot’s low-level capabili-
114 ties to effectively choose and sequence appropriate skills for the task at hand. To implement an
115 automated System 2 component as a code-writing VLM agent in order to autonomously execute the
116 verbalized plans without additional modules or human effort. To facilitate this, we define an API
117 for the action primitives accessible by the VLM to interface with the System 1, reactive low-level
118 RT-1 policy skills as described in Section 3. The API is based on translating the language com-
119 mands into a simple API that the VLM agent can access. This breakdown is based on *what* the
120 robot should do and *how* to do it. Each primitive skill (i.e. grasping, rotating, lifting, placing) is
121 represented by a function with a keyword argument modifying *how* that primitive is accomplished
122 (i.e. `grasp(object, "top-down")`). Internally, the API translates this code into the corresponding
123 natural language the RT-1 policy was trained on. We use a system prompt to tell the VLM to control
124 its physical embodiment through code, then provide the robot’s visual observation of the scene and
125 a description of the high-level task. The exact system prompt we use in all experiments and example
126 outputs and explanations produced by the model can be found on our project website¹.

127 4 Experiments

128 We evaluate STEER by testing its ability to improve grasping in unseen scenarios and perform novel
129 behaviors that require complex reasoning and motor control. We focus on three main questions:
130 (1) Does learning multiple modes of behavior improve adaptability in new situations? (2) Can
131 combining extracted skills from heterogeneous human demos enable entirely new tasks? (3) To
132 what degree can a state-of-the-art VLM plan orchestrate these skills autonomously?

133 We use a 7 DoF arm, a two-fingered gripper, and a mobile base, as used in RT-1 [4] in a tabletop
134 environment. The experiment involves 70K demonstrations from RT-1’s multi-task dataset and 15K
135 grasping demos from MOO [20]. We choose RT-1 [4] for our System 1 component and Gemini 1.5
136 Pro [39] as our learned System 2 component for those experiments (sample videos¹).

137 **Improving Test-Time Adaptability** We present three challenging, unseen grasping scenarios: a
138 kettle, a potted plant, and grasping in clutter. We compare STEER with the baseline RT-1 [4] (trained
139 on original instructions) and OpenVLA [8], which fine-tunes a VLM on robot data from Open
140 X-Embodiment [40]. We report the success rates in Figure 3a. RT-1 occasionally succeeds, but
141 exhibits different strategies and we observe that failures are often caused by a sub-optimal approach.
142 OpenVLA performed similarly to RT-1, demonstrating that additional web data does not lead to
143 sufficiently strong embodied reasoning about how to grasp in a new scenario where a particular
144 approach is evidently necessary. For example, we find that OpenVLA often picks the potted plant
145 up, but does not respect the language instruction of picking up the flower pot *without disturbing*



(a) Grasping results. Full successes are in solid colors. Partial successes are in light colors. We do 20 trials of Kettle, 10 of Potted Plant, and 15 of Fruit in Clutter. (b) New task results. We run each method 10 times, comparing the low-level capabilities afforded by each model to perform the task usin human guidance. (c) We compare a VLM to a human in STEERING the learned policies. The VLM can effectively recover most of the performance.

Figure 3: Results on grasping in unseen scenarios and performing a new task, with human or VLM guidance. We find that by having access to and being able to reason about extracted low-level strategies enables higher success in OOD scenarios than the baseline RT-1 model and a state-of-the-art VLA.

146 *the plant* and grasps from above around the plant leaves. Decomposing the grasp strategies and
 147 exploiting the most suitable one as we do in STEER reduces this failure mode.

148 **Performing Novel Behaviors** We study whether we can engineer behavior for a new everyday task
 149 without collecting new demonstrations or additional fine-tuning. Pouring is out of the distribution
 150 of demonstrated tasks but should be achievable with the motions that exist in the data. We compare
 151 against the **best-case** version of each of 4 baselines and comparisons: baseline **RT-1** [4], **Language**
 152 **motions from RT-H** [41], defined by narrating end-effector movement to give language like move
 153 arm left and rotate arm right, **a goal-image conditioned variant of RT-1**, which tests whether
 154 language is a better abstraction layer than goal images, and **OpenVLA** [8]. As seen in Figure 3b,
 155 human orchestration with a STEER policy achieves a 90% success rate on pouring as compared to
 156 70% with a policy trained with language motions from RT-H (whose orchestration is significantly
 157 more cumbersome as it requires tight closed-loop guidance). In comparison, baseline RT-1 cannot
 158 complete the task because it is not trained to reorient objects. The goal image conditioned baseline,
 159 despite having demonstration sub goal images from the same starting positions, fails and appears to
 160 mimic the exact arm positions in the subgoals rather than manipulate the object state as prescribed
 161 by the goal image. OpenVLA, despite having access to the same underlying demo data, does not
 162 generalize to the new motion by stitching together the appropriate motions.

163 **VLM Orchestration** Now, we test whether a VLM can effectively select or sequence appropriate
 164 skills afforded by STEER by reasoning about the context, in the visual observation and task descrip-
 165 tion, as well as the skills exposed through the API *without any examples* (i.e. 0-shot). For these
 166 experiments, we compare the VLM to human orchestration of the same low-level policy to serve as
 167 an upper bound on performance. Exact inputs and outputs can be found on the project website¹.

168 **Seen task, new scenarios.** We see that the VLM successfully produces the same high-level plans
 169 as the human expert very reliably for the grasping tasks. However, as shown in Figure 3c we see
 170 that there is a degradation in end-to-end task performance compared to human orchestration when
 171 executing the code produced by the VLM, and we analyze these failures. For the kettle picking
 172 task, we note that the low-level policy appears to be sensitive to the specific naming of objects.
 173 That is, the VLM often produced code to grasp the ‘black and white kettle’ from the top instead of
 174 grasping the ‘black and white object’ from the top, and with further analysis find that this instruction
 175 has a noticeable degradation across all low-level language-conditioned policies. So, while the VLM
 176 reasonably commands the policy to grasp from above, the low-level policy is less reliable. We expect
 177 this to be improved with denser annotation or augmentation on the entity-level, whereas STEER is
 178 concerned with the motion-level. For the Fruit in Clutter grasping task, the VLM did not always
 179 command the appropriate action and we suspect that similar object naming references (‘red apple’
 180 instead of ‘apple’) impact the low-level policy performance.

181 **Seen objects, new task.** Without any examples, the VLM correctly identifies that in order to pour
 182 from the cup, the robot ought to grasp it from the side as if a human were performing the task. It
 183 then recognizes that it must reorient it, then reorient it back in order to place it back upright on the
 184 table. The VLM succeeded in 6 out of 10 trials for zero-shot synthesizing of pouring behavior.

References

- 185
186 [1] T. Z. Zhao, V. Kumar, S. Levine, and C. Finn. Learning fine-grained bimanual manipulation
187 with low-cost hardware. *Robotics: Science and Systems (RSS)*, 2023.
- 188 [2] C. Chi, Z. Xu, S. Feng, E. Cousineau, Y. Du, B. Burchfiel, R. Tedrake, and S. Song. Diffusion
189 policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics*
190 *Research*, 2024.
- 191 [3] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn. BC-z:
192 Zero-shot task generalization with robotic imitation learning. In *5th Annual Conference on*
193 *Robot Learning*, 2021. URL <https://openreview.net/forum?id=8kbp23tSGYv>.
- 194 [4] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Haus-
195 man, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. Joshi, R. Ju-
196 lian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath,
197 I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao,
198 M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran,
199 V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-
200 1: Robotics transformer for real-world control at scale. In *arXiv preprint arXiv:2212.06817*,
201 2022.
- 202 [5] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess,
203 A. Dubey, C. Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to
204 robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- 205 [6] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson,
206 Q. Vuong, T. Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint*
207 *arXiv:2303.03378*, 2023.
- 208 [7] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna,
209 C. Xu, J. Luo, T. Kreiman, Y. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh,
210 C. Finn, and S. Levine. Octo: An open-source generalist robot policy. In *Proceedings of*
211 *Robotics: Science and Systems*, Delft, Netherlands, 2024.
- 212 [8] M. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster,
213 G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine,
214 P. Liang, and C. Finn. Openvla: An open-source vision-language-action model. *arXiv preprint*
215 *arXiv:2406.09246*, 2024.
- 216 [9] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine. Scaling cross-embodied learn-
217 ing: One policy for manipulation, navigation, locomotion and aviation. *arXiv preprint*
218 *arXiv:2408.11812*, 2024.
- 219 [10] D. Kahneman. *Thinking, fast and slow*. Farrar, Straus and Giroux, New York,
220 2011. ISBN 9780374275631 0374275637. URL https://www.amazon.de/Thinking-Fast-Slow-Daniel-Kahneman/dp/0374275637/ref=wl_it_dp_o_pdT1_nS_nC?ie=UTF8&colid=151193SNGKJT9&coliid=I30CESLZCVDFL7.
- 223 [11] OpenAI et al. Gpt-4 technical report, 2024. URL <https://arxiv.org/abs/2303.08774>.
- 224 [12] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. Le, and D. Zhou.
225 Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.
- 226
227 [13] G. Team et al. Gemini: A family of highly capable multimodal models, 2024. URL <https://arxiv.org/abs/2312.11805>.
228

- 229 [14] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code
230 as policies: Language model programs for embodied control. In *2023 IEEE International*
231 *Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023.
- 232 [15] N. Di Palo and E. Johns. Keypoint action tokens enable in-context imitation learning in
233 robotics. *arXiv preprint arXiv:2403.19578*, 2024.
- 234 [16] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakr-
235 ishnan, K. Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances.
236 In *6th Annual Conference on Robot Learning*, 2022.
- 237 [17] S. Lee, Y. Wang, H. Etukuru, H. J. Kim, N. M. M. Shafullah, and L. Pinto. Behavior generation
238 with latent actions, 2024. URL <https://arxiv.org/abs/2403.03181>.
- 239 [18] O. Mees, L. Hermann, and W. Burgard. What matters in language conditioned robotic imitation
240 learning over unstructured data. *IEEE Robotics and Automation Letters (RA-L)*, 7(4):11205–
241 11212, 2022.
- 242 [19] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra,
243 P. Bhargava, S. Bhosale, D. Bikel, L. Blecher, C. C. Ferrer, M. Chen, G. Cucurull, D. Esiobu,
244 J. Fernandes, J. Fu, W. Fu, B. Fuller, C. Gao, V. Goswami, N. Goyal, A. Hartshorn, S. Hosseini,
245 R. Hou, H. Inan, M. Kardas, V. Kerkez, M. Khabsa, I. Kloumann, A. Korenev, P. S. Koura,
246 M.-A. Lachaux, T. Lavril, J. Lee, D. Liskovich, Y. Lu, Y. Mao, X. Martinet, T. Mihaylov,
247 P. Mishra, I. Molybog, Y. Nie, A. Poulton, J. Reizenstein, R. Rungta, K. Saladi, A. Schelten,
248 R. Silva, E. M. Smith, R. Subramanian, X. E. Tan, B. Tang, R. Taylor, A. Williams, J. X. Kuan,
249 P. Xu, Z. Yan, I. Zarov, Y. Zhang, A. Fan, M. Kambadur, S. Narang, A. Rodriguez, R. Stojnic,
250 S. Edunov, and T. Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. URL
251 <https://arxiv.org/abs/2307.09288>.
- 252 [20] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, S. Kirmani,
253 B. Zitkovich, F. Xia, C. Finn, and K. Hausman. Open-world object manipulation using pre-
254 trained vision-language models, 2023.
- 255 [21] T. Xiao, H. Chan, P. Sermanet, A. Wahid, A. Brohan, K. Hausman, S. Levine, and J. Tomp-
256 son. Robotic skill acquisition via instruction augmentation with vision-language models. In
257 *Proceedings of Robotics: Science and Systems*, 2023.
- 258 [22] V. Myers, B. C. Zheng, O. Mees, S. Levine, and K. Fang. Policy adaptation via language
259 optimization: Decomposing tasks for few-shot imitation, 2024. URL <https://arxiv.org/abs/2408.16228>.
- 261 [23] J. Zhang, K. Pertsch, J. Zhang, and J. J. Lim. Sprint: Scalable policy pre-training via language
262 instruction relabeling. In *2024 IEEE International Conference on Robotics and Automation*
263 *(ICRA)*, pages 9168–9175. IEEE, 2024.
- 264 [24] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay policy learning: Solving
265 long horizon tasks via imitation and reinforcement learning. *Conference on Robot Learning*
266 *(CoRL)*, 2019.
- 267 [25] K. Black, M. Nakamoto, P. Atreya, H. Walke, C. Finn, A. Kumar, and S. Levine. Zero-shot
268 robotic manipulation with pretrained image-editing diffusion models. *ArXiv*, abs/2310.10639,
269 2023. URL <https://api.semanticscholar.org/CorpusID:264172455>.
- 270 [26] K. Fang, P. Yin, A. Nair, and S. Levine. Planning to practice: Efficient online fine-tuning
271 by composing goals in latent space. *2022 IEEE/RSJ International Conference on Intelligent*
272 *Robots and Systems (IROS)*, pages 4076–4083, 2022. URL <https://api.semanticscholar.org/CorpusID:248834175>.
- 273

- 274 [27] J. Gu, S. Kirmani, P. Wohlhart, Y. Lu, M. G. Arenas, K. Rao, W. Yu, C. Fu, K. Gopalakrishnan,
275 Z. Xu, P. Sundaresan, P. Xu, H. Su, K. Hausman, C. Finn, Q. Vuong, and T. Xiao. Rt-trajectory:
276 Robotic task generalization via hindsight trajectory sketches. In *Robotics: Science and Systems*
277 (*RSS*), 2024.
- 278 [28] C. Wen, X. Lin, J. So, K. Chen, Q. Dou, Y. Gao, and P. Abbeel. Any-point trajectory modeling
279 for policy learning. *Robotics: Science and Systems (RSS)*, 2024.
- 280 [29] Y. Lee, S.-H. Sun, S. Somasundaram, E. S. Hu, and J. J. Lim. Composing complex skills by
281 learning transition policies. In *International Conference on Learning Representations*, 2019.
- 282 [30] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman. Relay policy learning: Solving
283 long-horizon tasks via imitation and reinforcement learning. *CoRL*, 2019.
- 284 [31] M. Dalal, D. Pathak, and R. R. Salakhutdinov. Accelerating robotic reinforcement learning
285 via parameterized action primitives. *Neural Information Processing Systems (NeurIPS)*, 34:
286 21847–21859, 2021.
- 287 [32] S. Nasiriany, H. Liu, and Y. Zhu. Augmenting reinforcement learning with behavior primitives
288 for diverse manipulation tasks. In *IEEE International Conference on Robotics and Automation*
289 (*ICRA*), 2022.
- 290 [33] J. Zhang, M. Heo, Z. Liu, E. Biyik, J. J. Lim, Y. Liu, and R. Fakoore. Extract: Efficient
291 policy learning by extracting transferrable robot skills from offline data. *arXiv preprint*
292 *arXiv:2406.17768*, 2024.
- 293 [34] P.-L. Bacon, J. Harb, and D. Precup. The option-critic architecture. In *Proceedings of the AAAI*
294 *conference on artificial intelligence*, 2017.
- 295 [35] O. Nachum, S. S. Gu, H. Lee, and S. Levine. Data-efficient hierarchical reinforcement learning.
296 *Neural Information Processing Systems (NeurIPS)*, 31, 2018.
- 297 [36] X. B. Peng, M. Chang, G. Zhang, P. Abbeel, and S. Levine. Mcp: Learning composable hi-
298 erarchical control with multiplicative compositional policies. *Advances in Neural Information*
299 *Processing Systems*, 32, 2019.
- 300 [37] T. Nagarajan, C. Feichtenhofer, and K. Grauman. Grounded human-object interaction hotspots
301 from video. In *International Conference on Computer Vision (ICCV)*, 2019.
- 302 [38] S. Bahl, R. Mendonca, L. Chen, U. Jain, and D. Pathak. Affordances from human videos as a
303 versatile representation for robotics. 2023.
- 304 [39] G. Team. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of con-
305 text, 2024. URL <https://arxiv.org/abs/2403.05530>.
- 306 [40] O. X. E. Collaboration et al. Open x-embodiment: Robotic learning datasets and rt-x models :
307 Open x-embodiment collaboration0. In *2024 IEEE International Conference on Robotics and*
308 *Automation (ICRA)*, 2024.
- 309 [41] S. Belkhale, T. Ding, T. Xiao, P. Sermanet, Q. Vuong, J. Tompson, Y. Chebotar, D. Dwibedi,
310 and D. Sadigh. Rt-h: Action hierarchies using language. *arXiv preprint arXiv:2403.01823*,
311 2024.