

GHIL-Glue: Hierarchical Control with Filtered Subgoal Images

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Image and video generative models that are pre-trained on Internet-
2 scale data can increase the generalization capacity of robot learning systems.
3 These models can function as high-level planners, generating intermediate sub-
4 goals for low-level goal-conditioned policies to reach. However, the performance
5 of these systems can be bottlenecked by the interface between generative mod-
6 els and low-level controllers. Generative models may predict photorealistic yet
7 physically infeasible frames. Low-level policies may also be sensitive to subtle
8 visual artifacts in generated goal images. This paper addresses these facets
9 of generalization, providing an interface to “glue together” language-conditioned
10 image or video prediction models with low-level goal-conditioned policies. Our
11 method, Generative Hierarchical Imitation Learning-Glue (GHIL-Glue), filters
12 out subgoals that do not lead to task progress and improves the robustness of goal-
13 conditioned policies to generated subgoals with harmful visual artifacts. GHIL-
14 Glue achieves a new state-of-the-art on the CALVIN simulation benchmark for
15 policies using observations from a single RGB camera. GHIL-Glue also outper-
16 forms other generalist robot policies across 3/4 language-conditioned manipula-
17 tion tasks testing zero-shot generalization on a physical robot. Additional details
18 are available at <https://generative-hierarchical-glue.github.io>.

19 **Keywords:** Hierarchical Imitation Learning, Image Generation, Video Prediction

20 1 Introduction

21 As Internet-scale foundation models achieve success in computer vision and natural language pro-
22 cessing, a central question arises for robot learning: how can Internet-scale models enable embodied
23 behavior generalization? While one approach is to collect increasingly large action-labeled robot
24 manipulation training datasets [1, 2, 3], video datasets (without actions) from the Internet are vastly
25 larger. However, while videos may be useful for inferring the steps in a task, such as how the objects
26 should be moved, or which parts of an object to manipulate (e.g., grabbing a cup by the handle), they
27 are less useful for learning details about low-level control. For example, it is difficult to infer the
28 action commands for controlling a robot’s fingers from videos of humans performing manipulation
29 tasks. One promising solution to this challenge is to employ a hierarchical approach [4, 5]: infer
30 high-level subgoals in the form of goal images using models trained on Internet-scale videos, and
31 then fill in the fine-grained motions with low-level policies trained on robot data.

32 While this general approach has seen success in prior robotic manipulation work [6, 4, 7, 5, 8, 9],
33 the interface between the high-level planner generating subgoals and the low-level policy that must
34 reach these subgoals can be brittle. First, generative models may occasionally sample subgoals that
35 do not progress towards completing a given language instruction. If one such “off-task” subgoal is
36 followed, it can have a compounding errors effect, leading to subsequent subgoals being increasingly
37 “off-task.” Second, even if the generated subgoals lead to task progress, they can contain subtle
38 visual artifacts that degrade the performance of a naively trained low-level policy.



Figure 1: GHIL-Glue. We consider language-conditioned image and video prediction models that can generate multiple subgoals. GHIL-Glue has two components: augmentation de-synchronization (top) and subgoal filtering (bottom). **Subgoal filtering:** We train a classifier to identify which subgoal is most likely to progress towards completing the language instruction. This subgoal and the image observation are then passed to the low-level policy to choose a robot action. **Augmentation de-synchronization:** The distribution shift between subgoals sampled from the robot dataset during training and those sampled from the generative model during inference can degrade low-level policy and subgoal classifier performance. To robustify the low-level policy and subgoal classifier to artifacts in generated subgoals, we explicitly de-synchronize the image-augmentations applied to the current state (State Aug) and the sampled goal (Subgoal Aug).

- 39 To address these issues, we propose Generative Hierarchical Imitation Learning-Glue (GHIL-Glue)
40 (fig. 1), a method to *robustly* “Glue” together image or video generative models to a low-level robotic
41 control policy. **First**, we filter out “off-task” subgoals that are physically inconsistent with the com-
42 manded language instruction. We do this by training a subgoal classifier to predict the likelihood
43 of the transition between the current state and a given subgoal resulting in progress towards com-
44 pleting the provided language instruction. We then sample a number of candidate subgoals from the
45 generative model and choose the subgoal with the highest classifier ranking. **Second**, we identify
46 a simple yet non-obvious data augmentation practice to robustify the low-level policy and subgoal
47 classifier to visual artifacts in the generated subgoals. While image augmentations are ubiquitous in
48 robot learning methods, our key finding is that the standard way of applying image augmentations
49 does not make low-level policies robust to visual artifacts in generated subgoal images.
50 Experiments on the CALVIN [10] simulation benchmark and four language-conditioned tasks on
51 the Bridge V2 physical robot platform [11] suggest that GHIL-Glue improves upon prior SOTA
52 methods for zero-shot generalization while adding minimal additional algorithmic complexity.

53 2 GHIL-Glue

54 2.1 Subgoal Filtering

55 The image and video generative models we consider are first pre-trained on general Internet-scale
56 image and video data, and then fine-tuned on a modest amount of robot data. A common failure
57 mode we observe across different models is that, over the course of executing a task, the model
58 begins to go “off-task,” meaning that it starts generating subgoals that are consistent with the current
59 image observation but that do not progress towards completing the language instruction l . We hy-
60 pothesize that this is due to the distribution shift between the Internet image and video pre-training
61 data and the robot data they are fine-tuned on.

62 To address this challenge, we train a subgoal classifier $f_\theta(s, g, l)$ on a language-conditioned dataset
63 of trajectories \mathcal{D}_l that predicts the probability that the transition between the current image ob-
64 servation s and the next subgoal g makes progress towards completing language instruction l .
65 During training, we sample positive examples of state-goal transitions for l from the set of tra-
66 jectories that successfully complete the instruction. We construct negative examples in the fol-
67 lowing three ways: **1) Wrong Instruction:** (s, g, l') where l' is sampled from a different transi-
68 tion than s and g , **2) Wrong Goal Image:** (s, g', l) where g' is sampled from a different transi-
69 tion than s and l , and **3) Reverse Direction:** (g, s, l) , where the order of the current image
70 observation and the subgoal image have been switched. We refer to this dataset of negative ex-
71 amples constructed from \mathcal{D}_l as \mathcal{D}_l^- . We then train the subgoal classifier by minimizing the bi-
72 nary cross entropy loss between the positive examples and the constructed negative examples:

73 $\mathcal{J}(\theta) = \mathbb{E}_{(s,g,l) \sim \mathcal{D}_l} [\log(f_\theta(s, g, l))] + \mathbb{E}_{(s^-, g^-, l^-) \sim \mathcal{D}_l^-} [\log(1 - f_\theta(s^-, g^-, l^-))]$. At inference,
74 given a set of K subgoals predicted by the image or video model, GHIL-Glue selects the subgoal
75 with the highest classifier ranking to the low-level policy for conditioning.

76 2.2 Image Augmentation De-Synchronization

77 For both the low-level goal-conditioned policy and the subgoal classifier, each training sample in-
78 cludes two images: the current state s and the corresponding goal g . Applying image augmentation
79 procedures such as random cropping or color jitter during training is a standard approach in image-
80 based robot learning methods [12] to improve the robustness of learned models to distribution shifts
81 between their training and evaluation domains. Standard practice is to sample augmentation parame-
82 ters $\hat{\phi}$ and apply them to all images in a given training sample [4, 13], which corresponds to applying
83 the same $\hat{\phi}$ to both s and g . In a non-hierarchical policy setting, this makes sense, because at infer-
84 ence time s and g will both be sampled from the camera observations of the current environment
85 instantiation. However, when using an image or video prediction model for subgoal generation, at
86 inference time the low-level policy and subgoal classifier will see states from the camera observa-
87 tions, but the goals will be generated by the image or video prediction model. There will often be
88 differences in the visual artifacts between a camera observation s and the corresponding generated
89 subgoal image g , such as differences in color, contrast, blurriness, and the shapes of objects, which
90 can degrade the performance of low-level policies and subgoal classifiers.

91 To encourage robustness to this distribution shift, we sample separate augmentation parameters for
92 s and g , denoted by $\hat{\phi}_s$ and $\hat{\phi}_g$ (i.e., we de-synchronize the image augmentations applied to s and
93 g). Concretely, for each s and g pair sampled during training, a different random crop, brightness,
94 contrast, saturation, and hue shift are applied to s than are applied to g . This forces the low-level
95 policy and the subgoal classifier to be robust to differences in visual artifacts between s and g .

96 3 Experiments

97 3.1 Experimental Domains

98 **Simulation Experiment Setup:** Simulation experiments are performed in the CALVIN [10] bench-
99 mark, which focuses on long-horizon language-conditioned robot manipulation. We follow the same
100 protocol as in [4], and train on data from three environments (A, B, and C) and test policies on a
101 fully unseen environment (D). The held-out environment (D) contains unseen desk and object colors,
102 positions, and shapes. See appendix A for a visualization of the CALVIN environment.

103 **Physical Experiment Setup:** Physical experiments are performed with the Bridge V2 [11] ex-
104 periment setup with a WidowX250 robot. We use the same datasets as in [4] for training both the
105 high-level image prediction model and the low-level goal-conditioned policy. The Bridge V2 dataset
106 contains 45K language-annotated trajectories, which are used for the language-labeled robot dataset
107 $\mathcal{D}_{l,a}$. The remaining 15K trajectories are used for the action-only dataset \mathcal{D}_a . As in [4], we use a
108 filtered version of the Something-Something V2 dataset [14] with the same filtering scheme as in [4]
109 (resulting in 75K video clips) as our video-only dataset \mathcal{D}_l . We test our policies on four tasks on four
110 different cluttered table top scenes (fig. 2) on the Bridge V2 physical robot platform. These environ-
111 ments require generalizing to novel scenes, with novel objects, and with novel language commands
112 that are not seen in the Bridge V2 dataset. See appendix A for visualizations of the evaluation set-up.

113 3.2 Comparison Algorithms

114 To evaluate GHIL-Glue’s performance, we study the impact of applying it to two SOTA hierarchical
115 imitation learning algorithms: SuSIE [4] and UniPi [5]. We use either 4 or 8 candidate subgoals for
116 subgoal filtering (see appendix E for details). To evaluate the importance of hierarchy more gener-
117 ally, we also compare GHIL-Glue to a flat language-conditioned diffusion policy (LCBC Diffusion
118 Policy). Finally, we consider ablations where we separately study the impact of each of our proposed

119 contributions: subgoal filtering (section 2.1) and de-synchronizing augmentations (section 2.2). For
 120 physical experiments, we additionally consider a comparison to OpenVLA [15], which is trained on
 121 the Open X-Embodiment dataset [2] (which includes the Bridge V2 dataset).

122 **3.3 Experimental Results**

Method	Tasks completed in a row					
	1	2	3	4	5	Avg. Len.
LCBC Diffusion Policy	68.5%	43.0%	22.5%	11.0%	6.8%	1.52
SuSIE [4]	89.8%	75.0%	57.5%	41.8%	29.8%	2.94
GHIL-Glue (SuSIE) - Aug De-sync Only	95.2%	84.0%	69.5%	56.0%	46.2%	3.51
GHIL-Glue (SuSIE) - Subgoal Filtering Only	88.5%	75.5%	56.2%	43.0%	32.5%	2.96
GHIL-Glue (SuSIE)	95.2%	88.5%	73.2%	62.5%	49.8%	3.69
UniPi [5]	56.8%	28.3%	12.0%	3.5%	1.5%	1.02
GHIL-Glue (UniPi) - Aug De-sync Only	60.2%	29.5%	12.5%	5.5%	1.8%	1.1
GHIL-Glue (UniPi) - Subgoal Filtering Only	69.5%	40.0%	15.8%	6.5%	4.2%	1.36
GHIL-Glue (UniPi)	75.2%	44.8%	19.7%	11.2%	5.5%	1.56

Table 1: CALVIN: Simulation Results. Success rates on the validation tasks from the held-out D environment of the CALVIN zero-shot generalization challenge averaged across 4 random seeds. Applying GHIL-Glue to SuSIE and UniPi significantly improves performance over their respective base methods. GHIL-Glue (SuSIE) significantly outperforms all other methods, achieving a new state-of-the-art on the CALVIN benchmark for policies using observations from a single RGB camera.

	Task	OpenVLA [15]	SuSIE [4]	GHIL-Glue (SuSIE)
Scene A	Put Sushi On Towel	22/30	19/30	28/30
Scene B	Put Red Bell Pepper in Bowl	14/30	12/30	16/30
Scene C	Open Drawer	23/30	19/30	22/30
Scene D	Put Sushi in Bowl	15/30	15/30	18/30

Table 2: Bridge V2 Physical Experiments Results. Success rates across four tasks on four physical robot scenes (pictured in fig. 2) that test zero-shot generalization to novel objects, novel language commands, and novel scene configurations. GHIL-Glue applied to SuSIE outperforms SuSIE across all tasks and outperforms OpenVLA on 3 out of 4 tasks.

123 **Simulation Experiments:** We present results on the CALVIN benchmark in table 1. Applying
 124 GHIL-Glue yields significant improvements for SuSIE and UniPi, increasing the average successful
 125 task sequence length from **2.94** to **3.69** for SuSIE and from **1.02** to **1.56** for UniPi. **GHIL-Glue**
 126 (**SuSIE**) achieves a new SOTA on CALVIN for policies that use single RGB camera observations.

127 **Physical Experiments:** We present results (table 2) comparing GHIL-Glue (SuSIE) to OpenVLA
 128 and SuSIE across four environments on the Bridge V2 robot platform that require interacting with
 129 a number of objects on a cluttered table (fig. 2). GHIL-Glue applied to SuSIE outperforms SuSIE
 130 across all tasks and outperforms OpenVLA, a 7-billion parameter SOTA VLA, on 3 out of 4 tasks.
 131 Significantly, the baseline SuSIE implementation does not outperform OpenVLA on a single task,
 132 whereas **GHIL-Glue (SuSIE) outperforms OpenVLA on 3 out of 4 tasks**, demonstrating that hi-
 133 erarchical goal conditioned architectures with well-tuned interfaces between the high and low-level
 134 policies can outperform SOTA VLA methods on zero-shot generalization tasks. See appendix D for
 135 qualitative examples of success and failure cases of GHIL-Glue in physical experiments.

136 **4 Conclusion**

137 We present GHIL-Glue, a method for better aligning image and video prediction models and low-
 138 level control policies for hierarchical imitation learning. Our key insight is that while image and
 139 video foundation models can generate highly realistic subgoals for goal-conditioned policy learn-
 140 ing, when generalizing to novel environments, the generated images are prone to containing visual
 141 artifacts and can be inconsistent with the task the robot is commanded to perform. GHIL-Glue pro-
 142 vides two simple ideas to address these challenges, significantly improving zero-shot generalization
 143 performance over prior work in the CALVIN simulation benchmark and in physical experiments.

144 **References**

- 145 [1] S. Dasari, F. Ebert, S. Tian, S. Nair, B. Bucher, K. Schmeckpeper, S. Singh, S. Levine, and
146 C. Finn. Robonet: Large-scale multi-robot learning. In *Conference on Robot Learning (CoRL)*,
147 2019.
- 148 [2] O. X.-E. Collaboration, A. O'Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee,
149 A. Pooley, A. Gupta, A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Irpan,
150 A. Khazatsky, A. Rai, A. Gupta, A. Wang, A. Kolobov, A. Singh, A. Garg, A. Kembhavi,
151 A. Xie, A. Brohan, A. Raffin, A. Sharma, A. Yavary, A. Jain, A. Balakrishna, A. Wahid,
152 B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Wulfe, B. Ichter, C. Lu, C. Xu, C. Le, C. Finn,
153 C. Wang, C. Xu, C. Chi, C. Huang, C. Chan, C. Agia, C. Pan, C. Fu, C. Devin, D. Xu,
154 D. Morton, D. Driess, D. Chen, D. Pathak, D. Shah, D. Büchler, D. Jayaraman, D. Kalash-
155 nikov, D. Sadigh, E. Johns, E. Foster, F. Liu, F. Ceola, F. Xia, F. Zhao, F. V. Frujeri, F. Stulp,
156 G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Feng, G. Schiavi, G. Berseth, G. Kahn,
157 G. Yang, G. Wang, H. Su, H.-S. Fang, H. Shi, H. Bao, H. B. Amor, H. I. Christensen, H. Fu-
158 ruta, H. Walke, H. Fang, H. Ha, I. Mordatch, I. Radosavovic, I. Leal, J. Liang, J. Abou-Chakra,
159 J. Kim, J. Drake, J. Peters, J. Schneider, J. Hsu, J. Bohg, J. Bingham, J. Wu, J. Gao, J. Hu,
160 J. Wu, J. Wu, J. Sun, J. Luo, J. Gu, J. Tan, J. Oh, J. Wu, J. Lu, J. Yang, J. Malik, J. Silvério,
161 J. Hejna, J. Booher, J. Tompson, J. Yang, J. Salvador, J. J. Lim, J. Han, K. Wang, K. Rao,
162 K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund,
163 K. Kawaharazuka, K. Black, K. Lin, K. Zhang, K. Ehsani, K. Lekkala, K. Ellis, K. Rana,
164 K. Srinivasan, K. Fang, K. P. Singh, K.-H. Zeng, K. Hatch, K. Hsu, L. Itti, L. Y. Chen, L. Pinto,
165 L. Fei-Fei, L. Tan, L. J. Fan, L. Ott, L. Lee, L. Weihs, M. Chen, M. Lepert, M. Memmel,
166 M. Tomizuka, M. Itkina, M. G. Castro, M. Spero, M. Du, M. Ahn, M. C. Yip, M. Zhang,
167 M. Ding, M. Heo, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess,
168 N. J. Joshi, N. Suenderhauf, N. Liu, N. D. Palo, N. M. M. Shafiuallah, O. Mees, O. Kroemer,
169 O. Bastani, P. R. Sanketi, P. T. Miller, P. Yin, P. Wohlhart, P. Xu, P. D. Fagan, P. Mitrano,
170 P. Sermanet, P. Abbeel, P. Sundaresan, Q. Chen, Q. Vuong, R. Rafailov, R. Tian, R. Doshi,
171 R. Mart'in-Mart'in, R. Baijal, R. Scalise, R. Hendrix, R. Lin, R. Qian, R. Zhang, R. Men-
172 donca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Lin, S. Moore,
173 S. Bahl, S. Dass, S. Sonawani, S. Song, S. Xu, S. Haldar, S. Karamcheti, S. Adebola, S. Guist,
174 S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Ramamoorthy, S. Dasari, S. Belkhale, S. Park,
175 S. Nair, S. Mirchandani, T. Osa, T. Gupta, T. Harada, T. Matsushima, T. Xiao, T. Kollar, T. Yu,
176 T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, T. Chung, V. Jain, V. Vanhoucke,
177 W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Chen, X. Wang, X. Zhu, X. Geng, X. Liu,
178 X. Liangwei, X. Li, Y. Pang, Y. Lu, Y. J. Ma, Y. Kim, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Wu,
179 Y. Xu, Y. Wang, Y. Bisk, Y. Dou, Y. Cho, Y. Lee, Y. Cui, Y. Cao, Y.-H. Wu, Y. Tang, Y. Zhu,
180 Y. Zhang, Y. Jiang, Y. Li, Y. Li, Y. Iwasawa, Y. Matsuo, Z. Ma, Z. Xu, Z. J. Cui, Z. Zhang,
181 Z. Fu, and Z. Lin. Open X-Embodiment: Robotic learning datasets and RT-X models. 2024.
- 182 [3] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany,
183 M. K. Srirama, L. Y. Chen, K. Ellis, P. D. Fagan, J. Hejna, M. Itkina, M. Lepert, Y. J. Ma,
184 P. T. Miller, J. Wu, S. Belkhale, S. Dass, H. Ha, A. Jain, A. Lee, Y. Lee, M. Memmel, S. Park,
185 I. Radosavovic, K. Wang, A. Zhan, K. Black, C. Chi, K. B. Hatch, S. Lin, J. Lu, J. Mer-
186 cat, A. Rehman, P. R. Sanketi, A. Sharma, C. Simpson, Q. Vuong, H. R. Walke, B. Wulfe,
187 T. Xiao, J. H. Yang, A. Yavary, T. Z. Zhao, C. Agia, R. Baijal, M. G. Castro, D. Chen, Q. Chen,
188 T. Chung, J. Drake, E. P. Foster, J. Gao, D. A. Herrera, M. Heo, K. Hsu, J. Hu, D. Jackson,
189 C. Le, Y. Li, K. Lin, R. Lin, Z. Ma, A. Maddukuri, S. Mirchandani, D. Morton, T. Nguyen,
190 A. O'Neill, R. Scalise, D. Seale, V. Son, S. Tian, E. Tran, A. E. Wang, Y. Wu, A. Xie, J. Yang,
191 P. Yin, Y. Zhang, O. Bastani, G. Berseth, J. Bohg, K. Goldberg, A. Gupta, A. Gupta, D. Ja-
192 yaraman, J. J. Lim, J. Malik, R. Martín-Martín, S. Ramamoorthy, D. Sadigh, S. Song, J. Wu,
193 M. C. Yip, Y. Zhu, T. Kollar, S. Levine, and C. Finn. Droid: A large-scale in-the-wild robot
194 manipulation dataset. 2024.

- 195 [4] K. Black, M. Nakamoto, P. Atreya, H. Walke, C. Finn, A. Kumar, and S. Levine. Zero-
196 shot robotic manipulation with pretrained image-editing diffusion models. *arXiv preprint*
197 *arXiv:2310.10639*, 2023.
- 198 [5] Y. Du, S. Yang, B. Dai, H. Dai, O. Nachum, J. Tenenbaum, D. Schuurmans, and P. Abbeel.
199 Learning universal policies via text-guided video generation. *Advances in Neural Information*
200 *Processing Systems*, 36, 2024.
- 201 [6] I. Kapelyukh, V. Vosylius, and E. Johns. Dall-e-bot: Introducing web-scale diffusion models
202 to robotics. *IEEE Robotics and Automation Letters*, 2023.
- 203 [7] Y. Du, M. Yang, P. Florence, F. Xia, A. Wahid, B. Ichter, P. Sermanet, T. Yu, P. Abbeel, J. B.
204 Tenenbaum, et al. Video language planning. *arXiv preprint arXiv:2310.10625*, 2023.
- 205 [8] A. Ajay, S. Han, Y. Du, S. Li, A. Gupta, T. Jaakkola, J. Tenenbaum, L. Kaelbling, A. Srivastava,
206 and P. Agrawal. Compositional foundation models for hierarchical planning. *Advances in*
207 *Neural Information Processing Systems*, 36, 2024.
- 208 [9] J. Gao, K. Hu, G. Xu, and H. Xu. Can pre-trained text-to-image models generate visual goals
209 for reinforcement learning? *Advances in Neural Information Processing Systems*, 36, 2024.
- 210 [10] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard. Calvin: A benchmark for language-
211 conditioned policy learning for long-horizon robot manipulation tasks. In *IEEE Robotics and*
212 *Automation Letters (RAL)*, 2021.
- 213 [11] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch,
214 Q. Vuong, A. He, V. Myers, K. Fang, C. Finn, and S. Levine. Bridgedata v2: A dataset
215 for robot learning at scale. In *Conference on Robot Learning (CoRL)*, 2023.
- 216 [12] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for
217 transferring deep neural networks from simulation to the real world. *International Conference*
218 *on Intelligent Robots and Systems*, 2017.
- 219 [13] C. Zheng, B. Eysenbach, H. Walke, P. Yin, K. Fang, R. Salakhutdinov, and S. Levine. Stabi-
220 lizing contrastive rl: Techniques for offline goal reaching. *arXiv preprint arXiv:2306.03346*,
221 2023.
- 222 [14] R. Goyal, S. E. Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fru-
223 end, P. Yianilos, M. Mueller-Freitag, and et al. The “something something” video database for
224 learning and evaluating visual common sense. In *IEEE international conference on computer*
225 *vision (ICCV)*, 2017.
- 226 [15] M. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster,
227 G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine,
228 P. Liang, and C. Finn. Openvla: An open-source vision-language-action model. *arXiv preprint*
229 *arXiv:2406.09246*, 2024.
- 230 [16] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville. Film: Visual reasoning with
231 a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*,
232 volume 32, 2018.
- 233 [17] J. Ho and T. Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
234 2022.

235 **A Experimental Domains**

236 We study the degree to which GHIL-Glue improves existing hierarchical imitation learning algo-
 237 rithms across a number of tasks in simulation and physical experiments that assess zero-shot gen-
 238 eralization. We evaluate our method on the CALVIN [10] simulation benchmark and the Bridge
 239 V2 [11] physical experiment setup with a WidowX250 robot. The experimental domains are visual-
 240 ized in fig. 2.

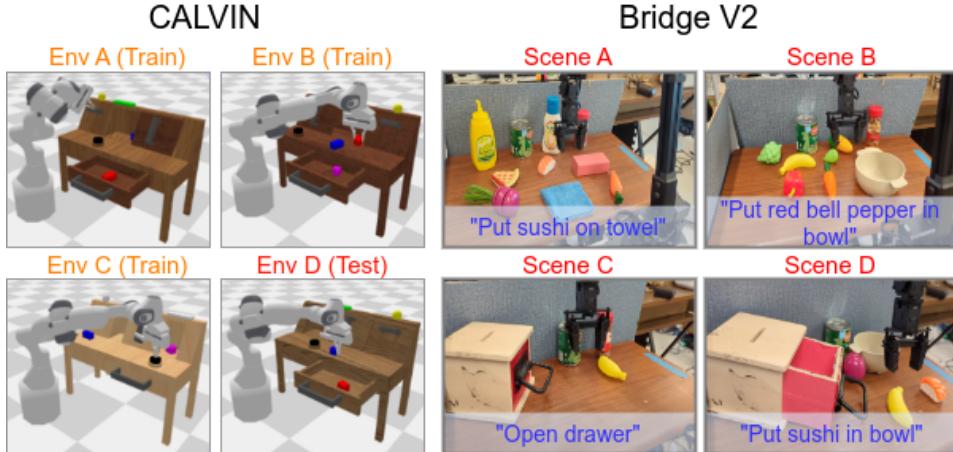


Figure 2: Experimental Domains. Simulation Environments (Left): Train/test environments in the CALVIN simulation benchmark. The environments each have different table textures, furniture positions, and initial configurations of the colored blocks. Each environment contains 34 tasks, each with an associated language instruction. To test zero-shot generalization, environment D is held out for evaluation. Physical Environments (Right): We consider four test scenes in the Bridge V2 robot platform with four total language instructions. To test zero-shot generalization, these test scenes contain novel objects, language commands, and object configurations not seen in the training data.

241 **B Classifier Training**

242 **Training objective:** The classifier is trained using binary cross-entropy loss:

$$\mathcal{J}(\theta) = \mathbb{E}_{(s,g,l) \sim D_l} [\log(f_\theta(s, g, l))] + \mathbb{E}_{(s',g',l') \sim N(D_l)} [1 - \log(f_\theta(s', g', l'))],$$

243 where D_l is the language-annotated dataset that consists of trajectory and language task pairs, and
 244 N is a function for generating negative examples from the dataset. Given a dataset D_l , N generates
 245 negatives from D_l in the following ways:

- 246 1. **Wrong Instruction:** (s, g, l') where l' is sampled from a different transition than s and g .
- 247 2. **Wrong Goal Image:** (s, g', l) where g' is sampled from a different transition than s and l .
- 248 3. **Reverse Direction:** (g, s, l) , where the order of the current image observation and the
 249 subgoal image have been switched.

250 Across all our experiments, we sample 50% of each training batch to be positive examples and
 251 50% of each training batch to be negative examples. Of the negative examples, 40% are “wrong
 252 instruction”, 40% are “reverse direction”, and 20% are “wrong goal image”.

253 **Goal sampling:** In a given training tuple (s_t, g, l) , g is sampled by taking the goal image from the
 254 s_{t+k} , where k is a uniformly sampled integer from 16 to 24.

255 **Network architecture and training hyperparameters:** The classifier network architecture consists
 256 of a ResNet-34 encoder from [11], followed by a two-layer MLP with layers of dimension 256.

257 Separate encoders are used to encode the image observations and the goal images (parameters are
258 not shared between the two). Both of these encoders use FiLM conditioning [16] after each residual
259 block to condition on the language instruction. Classifier networks are trained using a learning rate
260 of 3×10^{-4} and a batch size of 256 for 100,000 gradient steps. A dropout rate of 0.1 is used.

261 C Image Augmentations

262 During training of low-level policy networks and classifier networks, we apply the following aug-
263 mentations to the image observations and the goal images, in the following order:

- 264 1. Random Resized Crop:
 - 265 • scale: (0.8, 1.0)
 - 266 • ratio: (0.9, 1.1)
- 267 2. Random Brightness Shift:
 - 268 • shift ratio: 0.2
- 269 3. Random Contrast:
 - 270 • Contrast range: (0.8, 1.2)
- 271 4. Random Saturation:
 - 272 • Saturation range: (0.8, 1.2)
- 273 5. Random Hue:
 - 274 • shift ratio: 0.1

275 Figure 3 visualizes examples from the Bridge dataset before and after augmentations are applied:



Figure 3: Image augmentation examples Examples of images from the Bridge dataset before and after having the image augmentations applied to them that are used during policy and classifier training.

276 D Qualitative Analysis

277 D.1 Effect of subgoal filtering

278 Although we use classifier-free guidance (CfG) [17] on the image or video generative model with
279 respect to the language-prompt at inference in our experiments, we find that producing “off-task”
280 subgoals is still a common failure mode that is not solved by increasing the guidance parameter
281 alone. In fig. 4, we visualize how subgoal filtering can prevent “off-task” subgoals generated by the
282 image or video model from being passed to the low-level control policy.

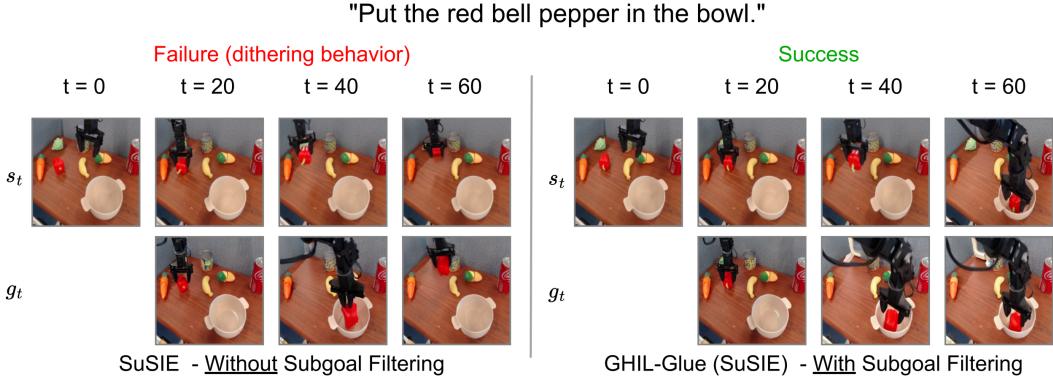
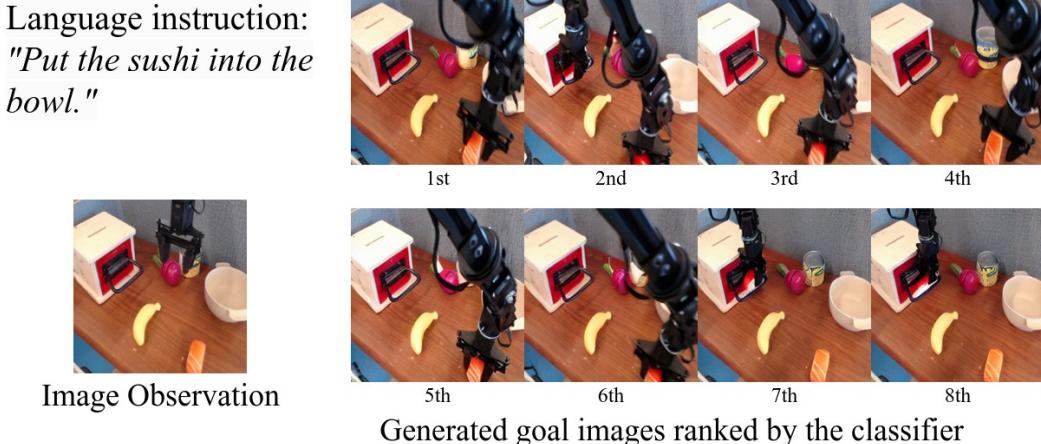


Figure 4: GHIL-Glue Subgoal Filtering. We visualize policy rollouts of SuSIE without subgoal filtering vs. GHIL-Glue SuSIE with subgoal filtering. We show the states reached every 20 timesteps (top row) and the corresponding predicted subgoals (bottom row). Without subgoal filtering, the subgoal at $t = 60$ is not consistent with making progress towards placing the pepper in the bowl, causing the robot to dither and drop the pepper. When subgoal filtering is used, the selected subgoals make iterative progress towards a successful task completion.

283 D.2 Classifier rankings

284 We show examples of how the classifier network ranks generated goal images on tasks from Scene D
 285 of our physical experimental domain. Figures 5a, 5b, 5c show examples of the classifier correctly
 286 ranking the generated goal images (highly ranked images correspond to making progress towards
 287 correctly completing the language instruction), while fig. 5d shows an example of the classifier
 288 erroneously giving high rankings to goal images that do not make progress towards completing the
 289 language instruction. Note that while the classifier scores can be close across various goal images,
 290 so long as the relative ranking of the generated goal images is correct, then incorrect subgoal images
 291 will be rejected and correct subgoal images will be passed to the low-level policy.

Figure 5: Classifier ranking examples Examples of the classifier network rankings on 8 generated candidate subgoals given an observation from Scene D of the physical experiments and a language instruction. Note that during GHIL-Glue inference, only the first-ranked subgoal is passed to the low-level policy.

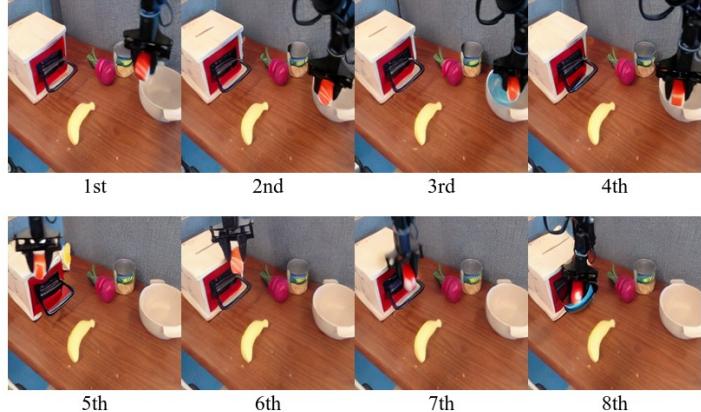


(a) Correct Example of Classifier Filtering The classifier correctly ranks the subgoal images where the robot is grasping the sushi higher than the subgoal images where the robot is grasping the drawer handle.

Language instruction:
"Put the sushi into the bowl."



Image Observation



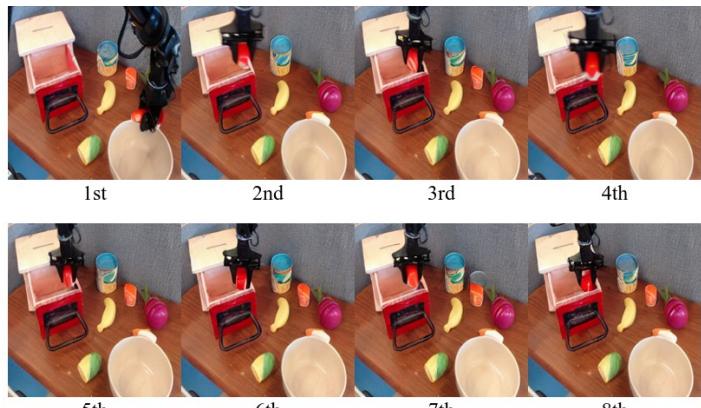
Generated goal images ranked by the classifier

(b) Correct Example of Classifier Filtering The classifier correctly ranks the subgoal images where the robot moves to place the grasped sushi into the bowl higher than the subgoal images where the robot moves its gripper towards the drawer handle. It ranks the subgoal image with the hallucinated blue bowl-like artifact last.

Language instruction:
"Put the sushi into the bowl."



Image Observation



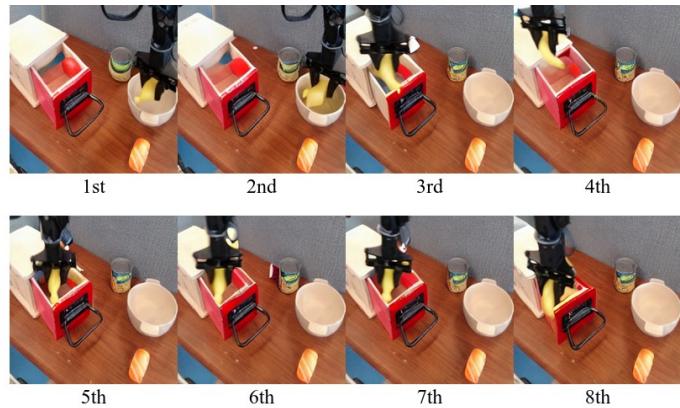
Generated goal images ranked by the classifier

(c) Correct Example of Classifier Filtering The classifier correctly ranks the subgoal image highest that shows the robot completing the correct task – only a single generated subgoal image shows the robot placing the sushi into the bowl, while all other generated subgoal images show the robot placing the sushi into the drawer.

Language instruction:
"Put the banana into the drawer."



Image Observation



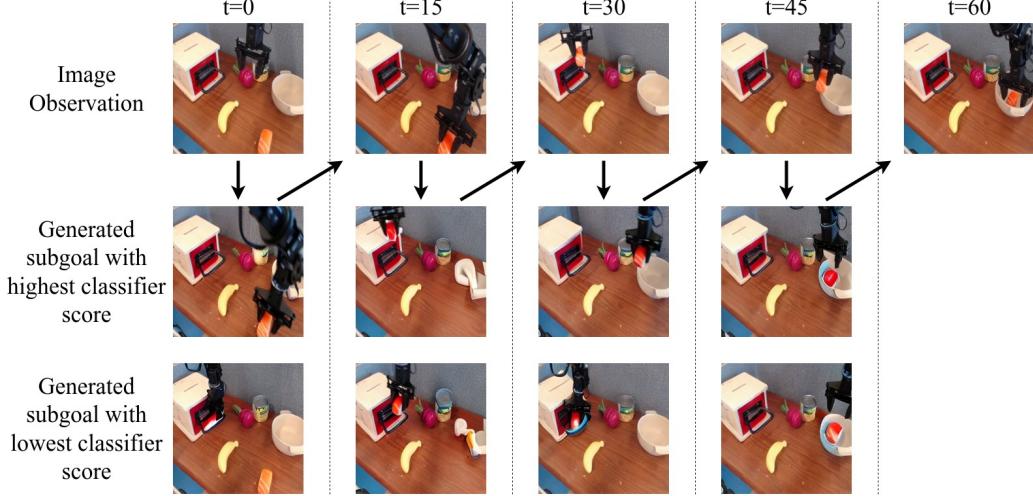
Generated goal images ranked by the classifier

(d) The classifier incorrectly ranks the subgoal images higher where the robot is placing the banana into the bowl than it ranks the subgoal images where the robot is placing the banana into the drawer. This could be due to there being a strong bias for placing objects in bowls in the Bridge V2 training data.

292 **D.3 Trajectory Visualizations**

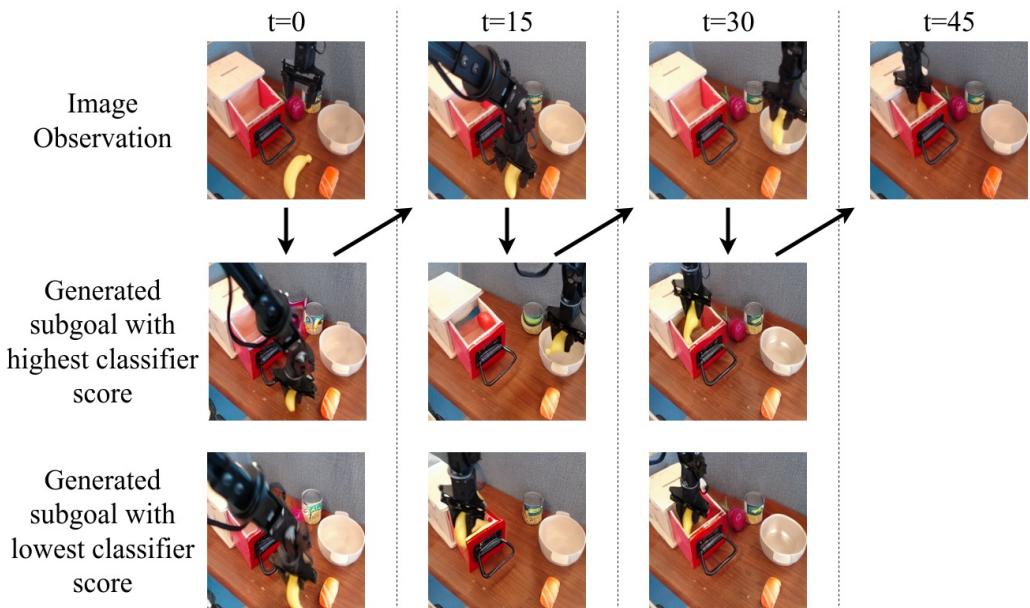
293 We show examples of rollouts of GHIL-Glue (SuSIE) on our physical experiment set up. These
 294 examples showcase when GHIL-Glue successfully filters out off-task subgoal images (Figure 6a),
 295 as well as an instance of when GHIL-Glue nearly causes a failure (Figure 6b).

Figure 6: GHIL-Glue (SuSIE) Trajectory Visualization Visualization of a rollout of GHIL-Glue (SuSIE) on Scene D in the physical experiments set up. The top row shows the current image observation at every timestep at which the video prediction model is queried. The second and third rows show the highest and lowest ranked generated subgoal images out of the 8 generated subgoal images, as ranked by the classifier. Note that during GHIL-Glue inference, only the first-ranked subgoal is passed to the low-level policy.



Language instruction: "Put the sushi into the bowl."

(a) **"Put the sushi into the bowl."** This rollout shows two examples of the classifier filtering preventing the policy from going off-task: at $t = 0$, the lowest ranked generated subgoal shows the gripper grasping the drawer handle instead of moving to grasp the sushi; at $t = 30$, the lowest ranked generated subgoal shows the gripper moving towards the drawer handle instead of towards placing the sushi into the bowl. Note the hallucinated objects and artifacts visible in the goal images at $t = 15, 30, 45$. Augmentation de-synchronization helps to make the low-level policy and classifier robust to hallucinated artifacts such as these.



Language instruction: "*Put the banana into the drawer.*"

(b) “Put the banana into the drawer.” In this rollout, classifier filtering fails and causes a near-miss. At $t = 15$, the classifier ranks a subgoal image highest that shows the robot placing the banana into the bowl instead of the drawer. However, at $t = 30$, when the robot reaches the state specified by this subgoal image, the subsequent generated subgoals all show the robot correctly placing the banana into the drawer. Although, as in this example, the classifier network can occasionally rank incorrect subgoal images higher than correct subgoal images, such errors occur infrequently as GHIL-Glue (SuSIE/UniPi) outperforms base-SuSIE/UniPi across all of our physical and simulated experiments.

296 **D.4 Qualitative Analysis of Augmentation De-synchronization**

297 We see that when applying aug-
 298 mentation de-synchronization, the
 299 number of failures due to low-
 300 level policy errors (missed grasps,
 301 dropping held objects, etc.) de-
 302 creases, indicating that augmentation
 303 de-synchronization is important for
 304 the low-level policy to be able to cor-
 305 rectly interpret and follow the sub-
 306 goal images generated by the video
 307 prediction model. This is particularly
 308 important in domains where there is
 309 a large visual generalization gap be-
 310 tween the training data and the eval-
 311 uation tasks. For example, in the
 312 CALVIN benchmark, the colors and
 313 shapes of objects differ between the
 314 training and evaluation scenes. This
 315 difference causes the subgoals gener-
 316 ated by the video prediction model to
 317 often contain objects with incorrect shapes and colors (Figure 7). Augmentation de-synchronization
 318 seems to be critical to allowing the low-level policy to be robust to these hallucinations and artifacts.

Language Instruction: "Go push the red block left."

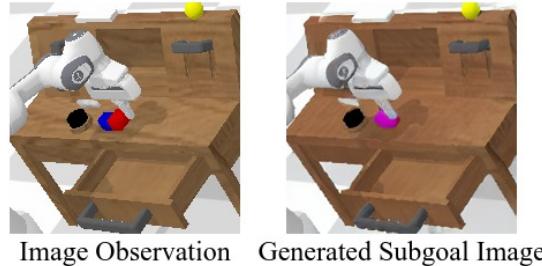


Figure 7: Generated Subgoal Image on CALVIN A subgoal image generated by the SuSIE video model on the unseen environment D of the CALVIN benchmark. The colors and shapes of objects are different in each of the four CALVIN environments, and since the model was not trained on data from environment D, it often generates images with incorrect shapes and colors. Augmentation de-synchronization is important for the low-level policy and classifier to be able to handle these mismatches between image observations and corresponding generated subgoal images.

319 **E Number of Candidate Subgoals**

320 We conduct an ablation over the number of candidate subgoals used for subgoal filtering in GHIL-
 321 Glue (SuSIE) in the CALVIN benchmark. We find that GHIL-Glue (SuSIE) achieves similar per-
 322 formance whether 4, 8, or 16 candidate subgoals are used. In our main results (section 3.3), we
 323 report the performance of GHIL-Glue (SuSIE) on the CALVIN benchmark when using 8 candidate
 324 subgoals for filtering. For GHIL-Glue (UniPi) on the CALVIN benchmark, we use 4 candidate sub-
 325 goals for filtering, due to the increased computation burden of generating video subgoals with the
 326 UniPi video model vs. generating image subgoals with the SuSIE image model. In our physical
 327 experiments, we run GHIL-Glue (SuSIE) using 4 candidate subgoals for filtering.

Method	Tasks completed in a row					
	1	2	3	4	5	Avg. Len.
GHIL-Glue (SuSIE) - 4 samples	95.2%	86.0%	71.2%	60.5%	50.0%	3.63
GHIL-Glue (SuSIE) - 8 samples	95.2%	88.5%	73.2%	62.5%	49.8%	3.69
GHIL-Glue (SuSIE) - 16 samples	95.0%	86.5%	72.8%	60.8%	48.0%	3.63

Table 3: Effect of Number of Candidate Goal Images Sampled in GHIL-Glue (SuSIE) Success rates on the validation tasks from environment D of the CALVIN Challenge when using GHIL-Glue (SuSIE) when using 4, 8, or 16 candidate goal images with classifier filtering. Results are averaged across 4 random seeds. Results are similar across all numbers of samples, with 8 samples performing the best by a slight margin.