

GHIL-Glue: Hierarchical Control with Filtered Subgoal Images

Kyle B. Hatch¹ Ashwin Balakrishna¹ Oier Mees² Suraj Nair¹ Seohong Park² Blake Wulfe¹
Masha Itkina¹ Benjamin Eysenbach³ Sergey Levine² Thomas Kollar¹ Benjamin Burchfiel¹

Abstract—Image and video generative models that are pre-trained on Internet-scale data can greatly increase the generalization capacity of robot learning systems. These models can function as high-level planners, generating intermediate subgoals for low-level goal-conditioned policies to reach. However, the performance of these systems can be greatly bottlenecked by the interface between generative models and low-level controllers. For example, generative models may predict photo-realistic yet physically infeasible frames that confuse low-level policies. Low-level policies may also be sensitive to subtle visual artifacts in generated goal images. This paper addresses these two facets of generalization, providing an interface to effectively “glue together” language-conditioned image or video prediction models with low-level goal-conditioned policies. Our method, Generative Hierarchical Imitation Learning-Glue (GHIL-Glue), filters out subgoals that do not lead to task progress and improves the robustness of goal-conditioned policies to generated subgoals with harmful visual artifacts. GHIL-Glue achieves a new state-of-the-art on the CALVIN simulation benchmark for policies using observations from a single RGB camera. GHIL-Glue also outperforms other generalist robot policies across 3/4 language-conditioned manipulation tasks testing zero-shot generalization in physical experiments. Additional details are available at <https://ghil-glue.github.io>.

I. INTRODUCTION

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

Modern hierarchical imitation learning algorithms [4], [5] typically use an image or video generative model trained on

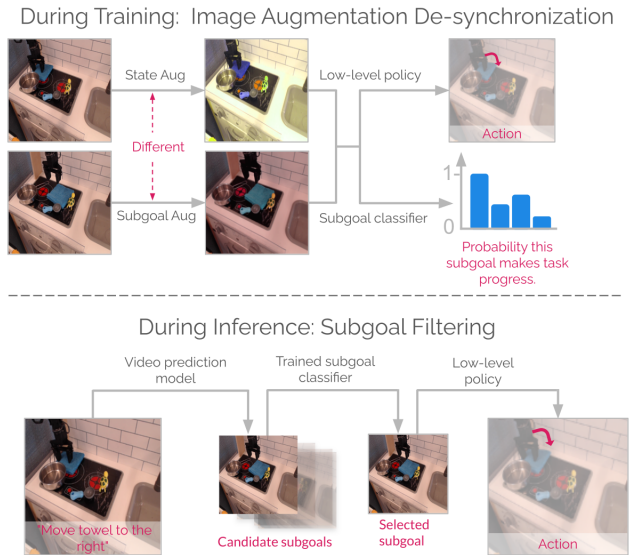


Fig. 1: **GHIL-Glue**. We consider language-conditioned image and video prediction models that can generate multiple subgoals. GHIL-Glue has two components: augmentation de-synchronization (top) and subgoal filtering (bottom). **Subgoal filtering**: We train a classifier to identify which subgoal is 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).

Internet-scale data to predict subgoal images, and then use a low-level control policy to translate these subgoal images into a sequence of motor commands [4], [5]. This approach allows the generative model to shoulder the hardest aspects of robotic generalization, such as generalizing to novel scenes, objects, and tasks. The low-level policy is then left with the comparatively easy task of choosing actions to reach these goals over short time horizons, which can be learned from a modest amount of robot data.

While this general approach has seen success in prior robotic manipulation work [4]–[9], the interface between the high-level planner generating subgoals and the low-level policy that must reach these subgoals can be brittle. State-of-the-art (SOTA) image or video prediction models are effective at generating likely subgoal images given a language prompt describing the task, but these subgoal generations may not be functionally useful for control. First, generative models may

Correspondence to: kyle.hatch@tri.global

¹Toyota Research Institute

²UC Berkeley

³Princeton University

occasionally sample subgoals that do not progress towards completing a given language instruction. If one such “off-task” subgoal is followed, it can have a compounding errors effect, leading to subsequent subgoals being increasingly “off-task.” Second, even if the generated subgoals lead to task progress, they can contain subtle visual artifacts that degrade the performance of a naively trained low-level policy.

We propose Generative Hierarchical Imitation Learning-Glue (GHIL-Glue) (Fig. 1), a method to *robustly* “Glue” together image or video generative models to a low-level robotic control policy. Our method is based on two components. **First**, we filter out “off-task” subgoals that are physically inconsistent with the commanded language instruction. We do this by training a subgoal classifier to predict the likelihood of the transition between the current state and a given subgoal resulting in progress towards completing the provided language instruction. We then sample a number of candidate subgoals from the generative model and choose the subgoal with the highest classifier ranking. **Second**, we identify a simple yet non-obvious data augmentation practice to robustify both the low-level policy and our subgoal classifier to visual artifacts in the generated subgoals. While image augmentations are ubiquitous in robot learning methods, our key finding is that the standard way of applying image augmentations does not make low-level policies robust to visual artifacts in generated subgoal images.

Experiments on the CALVIN [10] simulation benchmark and four language-conditioned tasks on the Bridge V2 physical robot platform [11] suggest that GHIL-Glue improves upon prior SOTA methods for zero-shot generalization while adding minimal additional algorithmic complexity.

II. RELATED WORK

Generative Models for Robotic Control: Prior works have explored diverse ways to leverage generative models, such as diffusion models [12], [13] and Transformers [14], for robotic control. They have employed highly expressive generative models, potentially pre-trained on Internet-scale data, for low-level control [15]–[20], data augmentation [21]–[23], object detection [24], [25], semantic planning [26]–[30], and visual planning [4]–[9]. Among them, our work is most related to prior works that employ image or video prediction models to generate intermediate subgoal images for the given language task [4]–[9]. These works use diffusion models to convert language instructions into visual subgoal plans, which are then fed into low-level subgoal-conditioned policies to produce actions. While sensible, this configuration leads to failures due to the misalignment of the generative models and the low-level policies that control the robot behavior, as shown in our experiments (Section V).

Rejection Sampling: One of our key ideas in this paper is based on rejection sampling, where we sample multiple subgoal proposals from an image or video prediction model and pick the best one based on a learned subgoal classifier. The idea of test-time rejection sampling has been widely used in diverse areas of machine learning, such as filtering-based action selection in offline reinforcement

learning (RL) [31]–[34], response verification in natural language processing [35]–[37], and planning and exploration in robotics [28], [29], [38]–[40]. Previous works in robotics have proposed several ways to filter out infeasible plans generated by pre-trained foundation models [28], [29], [38], [39], [41]. Unlike these works, we focus on filtering visual subgoals instead of language plans [28], [39], [41], and do not involve any planning procedures [29] or structural knowledge [38]. While the subgoal classifier we train resembles the classifier from [42], our classifier differs in two key ways. First, we use our classifier to filter out “off-task” subgoals, whereas the classifier in [42] is used as a reward function for training downstream policies. Second, the classifier from [42] is conditioned on the initial state s_0 and the current state s , whereas our classifier is conditioned on the current state s and a generated subgoal g .

Goal-Conditioned Policy Learning: Our method is broadly related to goal-conditioned policy learning [43]–[45], language-conditioned policy learning [46]–[50], and hierarchical control [4], [5], [51]–[54]. Most prior works in hierarchical policy learning either train a high-level policy from scratch that produces subgoals or latent skills [52], [55]–[68] or employ subgoal planning [65], [69]–[81]. Unlike these works, we do not train a high-level subgoal prediction model from scratch nor involve a potentially complex planning procedure. Instead, we sample multiple potential subgoals from a pre-trained (or potentially fine-tuned) image or video prediction model and pick the best one based on a trained subgoal classifier. Among hierarchical policy methods, perhaps the closest work to ours is IRIS [51], which trains a conditional variational autoencoder to generate subgoal proposals and selects the best subgoal that maximizes the task value function. While conceptually similar, our method differs from IRIS in that we do not assume access to a reward function in order to train a value function. Our classifier is trained on trajectories consisting only of images and language descriptions.

Diffusion Model Guidance: The generative models we consider in our paper [82], [83] are diffusion-based models trained using classifier-free guidance (CfG) [84]. Although we use a large value for the language-prompt guidance parameter at inference in our experiments, we find that producing “off-task” subgoals is still a common failure mode that is not solved by increasing this parameter alone.

Classifier guidance [12], [85], [86] is also a plausible alternative to rejection sampling, but there are some practical challenges in training a subgoal classifier for this purpose. First, the diffusion models we consider use latent diffusion [87], and therefore would require training the subgoal classifier to operate in the latent space of the diffusion model. Second, the subgoal classifier would need to be trained on noised data in order to guide the diffusion denoising process of the generative model. Nevertheless, classifier guidance is a potentially appealing direction for future work.

III. PRELIMINARIES

We consider the same problem setting as [4], where the goal is for a robot to perform a task described by some previously unseen language command l . To do this, we consider the same three dataset categories as in [4]: (1) language-labeled video clips \mathcal{D}_l which contain no robot actions; (2) language-labeled robot data $\mathcal{D}_{l,a}$ that includes both language labels and robot actions; (3) unlabeled robot data that only includes actions \mathcal{D}_a . The dataset $\mathcal{D}_{l,a}$ consists of a set of trajectory and task language pairs, $\{(\tau^n, l^n)\}_{n=1}^N$, and a trajectory contains a sequence of state, $s_t^n \in \mathcal{S}$, and action, $a_t^n \in \mathcal{A}$, pairs, $\tau^n = (s_0^n, a_0^n, s_1^n, a_1^n, \dots)$. Given these datasets, we assume access to two learned modules:

- 1) a **subgoal generation module** from which we can sample multiple possible future subgoals. This can be trained on \mathcal{D}_l and $\mathcal{D}_{l,a}$.
- 2) a **low-level goal-reaching policy** that chooses actions to reach generated subgoals. This can be trained on \mathcal{D}_a and/or $\mathcal{D}_{l,a}$.

Our contribution is a set of approaches to robustify the interface between these two modules.

While GHIL-Glue can be applied to any hierarchical imitation learning method consisting of the two components mentioned above, in this work we apply GHIL-Glue to two specific algorithms: (1) UniPi [5], in which a high-level model generates a subgoal video, and a low-level inverse-dynamics model predicts the actions needed to “connect” the images in the video, and (2) SuSIE [4], in which a high-level model generates a subgoal image by “editing” the current image observation, and a goal-conditioned policy predicts actions to achieve the subgoal image. We define subgoals, $g \in \mathcal{G}$, as video or image samples from the high-level models used in these algorithms.

IV. GHIL-GLUE

Many modern hierarchical policy methods aim to improve generalization by using language-conditioned image or video models to generate intermediate subgoal images for a given task. The interface between these image or video models and the low-level policies that choose actions to reach generated subgoals is a major performance bottleneck for these hierarchical policy methods. GHIL-Glue improves the robustness of this interface (see Fig. 1). In Section IV-A, we propose a simple method to filter subgoals that do not progress towards completing the task specified by language instruction l . Then, in Section IV-B, we describe a simple yet non-obvious data augmentation practice to robustify the low-level policy and our subgoal classifier to harmful visual artifacts in the generated subgoals. We note that the two components of GHIL-Glue work together synergistically: when applied together, the resulting performance improvement is larger than the sum of improvements that results from applying each component individually (see Section V).

A. Subgoal Filtering

The image and video generative models we consider are first pre-trained on general Internet-scale image and video

data, and then fine-tuned on a modest amount of robot data. Despite being fine-tuned on robot data, a common failure mode we observe across different models is that, over the course of executing a task, the model begins to go “off-task,” meaning that it starts generating subgoals that are consistent with the current image observation but that do not progress towards completing the language instruction l . We hypothesize that this is due to the distribution shift between the Internet data these image or video prediction models are pre-trained on and the robot data they are fine-tuned on.

To address this challenge, we train a subgoal classifier $f_\theta(s, g, l)$ on \mathcal{D}_l and/or $\mathcal{D}_{l,a}$ that predicts the probability that the transition between the current image observation s and the next subgoal g makes progress towards completing language instruction l . Note that although we train the subgoal classifier on robot data in our experiments, action labels are not used in the training of the classifier, and the subgoal classifier can be trained on action-free data, including large, non-robotics Internet video datasets. During training, we sample positive examples of state-goal transitions for l from the set of trajectories that successfully complete the instruction. We construct negative examples in the following three ways:

- 1) **Wrong Instruction:** (s, g, l') where l' is sampled from a different transition than s and g .
- 2) **Wrong Goal Image:** (s, g', l) where g' is sampled from a different transition than s and l .
- 3) **Reverse Direction:** (g, s, l) , where the order of the current image observation and the subgoal image have been switched. This is important for learning whether a candidate goal image is making temporal progress towards completing the language instruction.

We refer to this dataset of negative examples constructed from $\mathcal{D}_{l,a}$ as $\mathcal{D}_{l,a}^-$. We then train the subgoal classifier by minimizing the binary cross entropy loss between the positive examples and the constructed negative examples (see Appendix A for additional training details):

$$\mathcal{J}(\theta) = \mathbb{E}_{(s,g,l) \sim \mathcal{D}_{l,a}} [\log(f_\theta(s, g, l))] + \mathbb{E}_{(s^-,g^-,l^-) \sim \mathcal{D}_{l,a}^-} [\log(1 - f_\theta(s^-, g^-, l^-))]. \quad (1)$$

Given a set of K subgoals predicted by the image or video model, GHIL-Glue uses the classifier to select the subgoal with the highest progress probability and passes that subgoal to the low-level policy for conditioning.

B. Image Augmentation De-Synchronization

While the method proposed in Section IV-A increases robustness to predicted subgoals that do not make task progress, generated subgoals can also contain visual artifacts that degrade the performance of both the low-level control policy and the subgoal classifier. This performance degradation results from the distribution shift between the subgoal images seen by the policy during training, which come from the robot dataset, and the subgoal images seen

during inference, which come from the generative model. Ideally, the low-level policy and subgoal classifier would be trained on the same distribution of *generated* subgoal images that they will see at inference time. However, due to the high degree of variance in sampling images from a generative model, there is not a clear way to obtain generated subgoal images that match the actual future states reached in trajectories in the training data. To address this issue, we identify a simple yet non-obvious data augmentation practice to train the low-level policy and subgoal classifier on goals from the robot dataset while also robustifying them to visual artifacts in generated subgoals.

Applying image augmentation procedures such as random cropping or color jitter during training is a standard approach in image-based robot learning methods [88] to improve the robustness of learned models to distribution shifts between their training and evaluation domains. More formally, let ϕ be the set of image augmentation parameters to be randomly sampled from space Φ , $p_{\Phi}(\cdot)$ be some probability distribution over Φ , and let $\hat{\phi} \sim p_{\Phi}(\cdot)$ be some realization of augmentations sampled from $p_{\Phi}(\cdot)$. Typically, for each training sample, a different value $\hat{\phi}$ is applied during training to make a model robust to any augmentation in the space Φ .

For both the low-level goal-conditioned policy and the subgoal classifier, each training sample includes two images: the current state s and the corresponding goal g . Standard practice is to sample augmentation parameters $\hat{\phi}$ and apply them to all images in a given training sample [4], [89], which corresponds to applying the same $\hat{\phi}$ to both s and g . In a non-hierarchical policy setting, this makes sense, because at inference time s and g will both be sampled from the camera observations of the current environment instantiation. However, when using an image or video prediction model for subgoal generation, at inference time the low-level policy and subgoal classifier will see states from the camera observations, but the goals will be generated by the image or video prediction model. There will often be differences in the visual artifacts between a camera observation s and the corresponding generated subgoal image g , such as differences in color, contrast, blurriness, and the shapes of objects, which can degrade the performance of low-level policies and subgoal classifiers.

To encourage robustness to this distribution shift, we sample separate augmentation parameters for s and g , denoted by $\hat{\phi}_s$ and $\hat{\phi}_g$ (i.e., we de-synchronize the image augmentations applied to s and g). Random cropping, brightness shifts, contrast shifts, saturation shifts, and hue shifts comprise our space of augmentations (see Appendix B for details). Concretely, for each s and g pair sampled during training, a different random crop, brightness, contrast, saturation, and hue shift are applied to s than are applied to g . This forces the low-level policy and the subgoal classifier to learn to make accurate predictions on (s, g) pairs that have differences in visual artifacts.

While image augmentations are ubiquitous in robot learning methods, our experiments show that the standard way of applying image augmentations for goal-conditioned policies

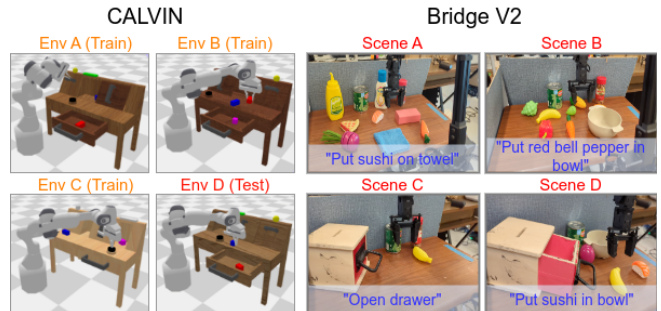


Fig. 2: **Experimental Domains.** Simulation Environments (Left): Train/test environments in the CALVIN simulation benchmark. The environments each have different table textures, furniture positions, and initial configurations of the colored blocks. Each environment contains 34 tasks, each with an associated 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.

and classifiers is deficient for the hierarchical policy methods that we consider. We also note that augmentation desynchronization is applied not only to the policy, but also to the subgoal classifier (Section IV-A), which has a significant impact on overall performance (Section V).

V. EXPERIMENTS

We study the degree to which GHIL-Glue improves existing hierarchical imitation learning algorithms across a number of tasks in simulation and physical experiments that assess zero-shot generalization. We analyze the influence of each component of GHIL-Glue on task performance and also perform extensive qualitative analysis in Appendix C.

A. Experimental Domains

We evaluate our method on the CALVIN [10] simulation benchmark and the Bridge V2 [11] physical experiment setup with a WidowX250 robot.

Simulation Experiment Setup: Simulation experiments are performed in the CALVIN [10] benchmark, which focuses on long-horizon language-conditioned robot manipulation. We follow the same protocol as in [4], and train on data from three environments (A, B, and C) and test policies on a fully unseen environment (D). Each environment contains a Franka Emika Panda robot arm that is placed in front of a desk with a variety of objects and is associated with 34 possible tasks (Fig. 2). The held-out environment (D) contains unseen desk and object colors, object and furniture positions, and object shapes. The corresponding language instructions are similarly held out.

Physical Experiment Setup: For physical experiments, we use the same datasets as in [4] for training both the high-level image prediction model and the low-level goal-conditioned policy. The Bridge V2 dataset contains 45K language-annotated trajectories, which are used for the language-labeled robot dataset $\mathcal{D}_{l,a}$. The remaining 15K trajectories are used for the action-only dataset \mathcal{D}_a . As

"Put the red bell pepper in the bowl."

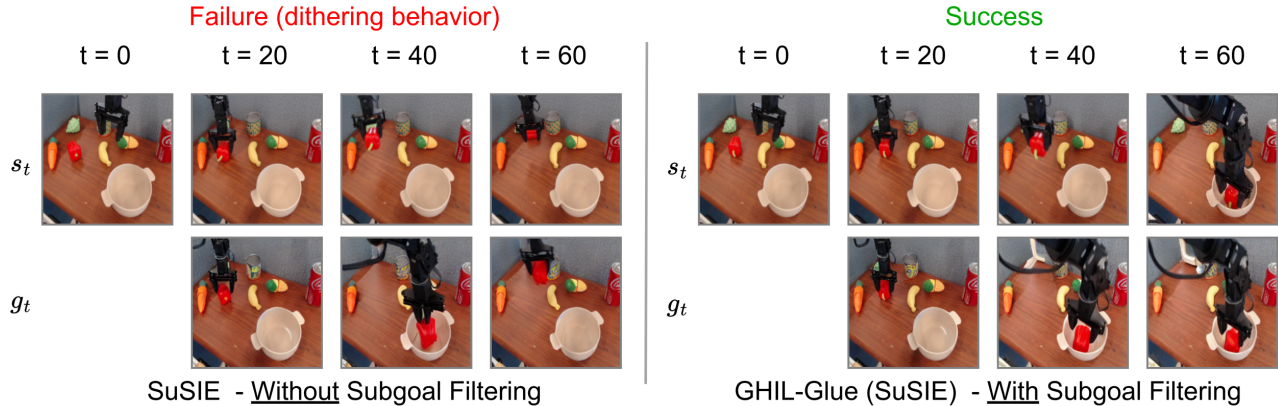


Fig. 3: **GHIL-Glue Subgoal Filtering.** We visualize policy rollouts of SuSIE without subgoal filtering vs. GHIL-Glue SuSIE with subgoal filtering. We show the states reached every 20 timesteps (top row) and the corresponding predicted subgoals (bottom row). Without subgoal filtering, the subgoal at $t = 60$ 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.

in [4], we use a filtered version of the Something-Something V2 dataset [90] with the same filtering scheme as in [4] (resulting in 75K video clips) as our video-only dataset \mathcal{D}_v .

We test our policies on four tasks on four different cluttered table top scenes (Fig. 2) on the Bridge V2 physical robot platform. These environments require generalizing to novel scenes, with novel objects, and with novel language commands that are not seen in the Bridge V2 dataset.

B. Comparison Algorithms

To evaluate GHIL-Glue’s performance, we study the impact of applying it to two SOTA hierarchical imitation learning algorithms: SuSIE [4] and UniPi [5]. To evaluate the importance of hierarchy more generally, we also compare GHIL-Glue to a flat language-conditioned diffusion policy (LCBC Diffusion Policy). Finally, we consider ablations where we separately study the impact of each of our proposed contributions: subgoal filtering (Section IV-A) and de-synchronizing augmentations (Section IV-B). For physical experiments, we additionally consider a comparison to OpenVLA [91], which is trained on the Open X-Embodiment dataset [2] (which includes the Bridge V2 dataset).

- 1) **LCBC Diffusion Policy:** Low-level language-conditioned behavior cloning diffusion policy [16] trained only on robot trajectories with language annotations. We use the same implementation as in [4].
- 2) **OpenVLA [17]:** A SOTA language-conditioned vision-language-action model (VLA) trained on the Open X-Embodiment dataset [2] (which includes the entirety of the Bridge V2 dataset).
- 3) **SuSIE [4]:** A method which fine-tunes Instruct-Pix2Pix [82], an image-editing diffusion model, to generate subgoal images given the current image observation. Low-level control is performed using a goal-conditioned policy. For SuSIE and all methods that build on it, we predict subgoals 20 steps in the future as in the original paper.

- 4) **UniPi [5]:** A method which fine-tunes a language-conditioned video prediction model on robot data and then uses an inverse dynamics model for low-level goal reaching. For UniPi and all methods that build on it, we predict video sequences of 16 frames. As the original UniPi model is not publicly available, we re-implement UniPi by fine-tuning the video model from [83].
- 5) **GHIL-Glue (SuSIE / UniPi):** GHIL-Glue applied on top of either SuSIE or UniPi. For all experiments we implement the subgoal filtering step by sampling four to eight subgoals from the high-level video prediction model and selecting amongst them (see Appendix D for details). We directly filter the subgoal images generated by the SuSIE model. We filter the video sequences generated by the UniPi model based on the final frame of each sequence.
- 6) **GHIL-Glue (SuSIE / UniPi) - Subgoal Filtering Only:** GHIL-Glue applied to SuSIE or UniPi using subgoal filtering but without augmentation de-synchronization.
- 7) **GHIL-Glue (SuSIE / UniPi) - Aug De-sync Only:** GHIL-Glue applied to SuSIE or UniPi using augmentation de-synchronization but without subgoal filtering.

C. Experimental Results

Simulation Experiments: We present results on the CALVIN benchmark in Table III. Applying GHIL-Glue yields significant performance increases for SuSIE and UniPi, increasing the average successful task sequence length from **2.94** to **3.69** for SuSIE and from **1.02** to **1.56** for UniPi. **GHIL-Glue (SuSIE) achieves a new SOTA on CALVIN** for policies that use observations from a single RGB camera. The two components of GHIL-Glue (subgoal filtering and image augmentation de-synchronization) improve performance when applied individually, but, when applied together, these components build on each other, leading to a performance increase greater than the sum of the individual benefits. Specifically, for SuSIE, image augmentation

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 I: **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 [91]	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 II: **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.

de-synchronization and subgoal filtering individually yield increases in sequence length of 0.56 and 0.02 respectively, whereas when applied together they yield an increase of 0.75. Similarly, for UniPi, the individual improvements yield increases in sequence length of 0.08 and 0.34 respectively, compared to an increase of 0.54 when applied together.

When applied alone, image augmentation de-synchronization increases the average successful task sequence length from 2.94 to 3.51 for SuSIE and from 1.02 to 1.1 for UniPi. We hypothesize that augmentation de-synchronization improves performance a large amount with SuSIE because its low-level policy is conditioned on a camera observation image s from the environment and a subgoal image g generated by the image model. When generalizing to the held-out test environment D, the SuSIE image model generates subgoal images with visual discrepancies from the camera observation images. In contrast, the UniPi video model predicts a sequence of frames as opposed to a single subgoal image. The UniPi low-level policy functions as an inverse dynamics model, choosing actions to link between the frames of the generated subgoal video, and is therefore conditioned on an s and g that both come from the predicted subgoal video.

When applied alone, subgoal filtering has a small effect on SuSIE, while on UniPi it increases the average successful task sequence length from 1.02 to 1.36. This suggests that unless the SuSIE low-level policy is made robust to visual artifacts in generated subgoals, simply selecting the most task relevant subgoal is insufficient to improve performance. As discussed previously, the SuSIE low-level policy is more sensitive to visual artifacts in generated subgoals than is the UniPi inverse dynamics model.

Physical Experiments: We present results (Table II) comparing GHIL-Glue (SuSIE) to OpenVLA and SuSIE across

four environments on the Bridge V2 robot platform that require interacting with a number of objects on a cluttered table (Fig. 2). These environments require generalizing to novel scenes, with novel objects, and with novel language commands that are not seen in the Bridge V2 dataset. GHIL-Glue applied to SuSIE outperforms SuSIE across all tasks and outperforms OpenVLA, a 7-billion parameter SOTA VLA, on 3 out of 4 tasks. Significantly, the baseline SuSIE implementation does not outperform OpenVLA on a single task, whereas GHIL-Glue (SuSIE) outperforms OpenVLA on 3 out of 4 tasks, demonstrating that hierarchical goal conditioned architectures with well-tuned interfaces between the high and low-level policies can outperform SOTA VLA methods on zero-shot generalization tasks. See Appendix C for qualitative examples of success and failure cases of GHIL-Glue in physical experiments, and for examples of generated subgoals for a subset of the tasks in addition to their scores under our subgoal filtering method.

VI. CONCLUSION

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

One exciting avenue for future work would be to explore training image or video prediction models for subgoal generation on a broader distribution of robot data, such as the data

available in the Open-X embodiment dataset [2]. Another interesting direction would be to filter subgoals based on the capability of the low-level policy to actually achieve them, for example, by training a goal-conditioned value function for the low-level policy and using it to evaluate subgoal feasibility. Finally, while we trained the subgoal classifiers on robot datasets, in principle these could be trained in the same way on much larger, non-robotics video datasets in order to improve generalization.

REFERENCES

- [1] S. Dasari, F. Ebert, S. Tian, S. Nair, B. Bucher, K. Schmeckpeper, S. Singh, S. Levine, and C. Finn, “Robonet: Large-scale multi-robot learning,” in *Conference on Robot Learning (CoRL)*, 2019.
- [2] O. X.-E. Collaboration, A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Gupta, A. Wang, A. Kolobov, A. Singh, A. Garg, A. Kembhavi, A. Xie, A. Brohan, A. Raffin, A. Sharma, A. Yavary, A. Jain, A. Balakrishna, A. Wahid, B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Wulfe, B. Ichter, C. Lu, C. Xu, C. Le, C. Finn, C. Wang, C. Xu, C. Chi, C. Huang, C. Chan, C. Agia, C. Pan, C. Fu, C. Devin, D. Xu, D. Morton, D. Driess, D. Chen, D. Pathak, D. Shah, D. Büchler, D. Jayaraman, D. Kalashnikov, D. Sadigh, E. Johns, E. Foster, F. Liu, F. Ceola, F. Xia, F. Zhao, F. V. Frujeri, F. Stulp, G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Feng, G. Schiavi, G. Berseth, G. Kahn, G. Yang, G. Wang, H. Su, H.-S. Fang, H. Shi, H. Bao, H. B. Amor, H. I. Christensen, H. Furuta, H. Walke, H. Fang, H. Ha, I. Mordatch, I. Radosavovic, I. Leal, J. Liang, J. Abou-Chakra, J. Kim, J. Drake, J. Peters, J. Schneider, J. Hsu, J. Bohg, J. Bingham, J. Wu, J. Gao, J. Hu, J. Wu, J. Wu, J. Sun, J. Luo, J. Gu, J. Tan, J. Oh, J. Wu, J. Lu, J. Yang, J. Malik, J. Silvério, J. Hejna, J. Booher, J. Tompson, J. Yang, J. Salvador, J. J. Lim, J. Han, K. Wang, K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund, K. Kawaharazuka, K. Black, K. Lin, K. Zhang, K. Ehsani, K. Lekkala, K. Ellis, K. Rana, K. Srinivasan, K. Fang, K. P. Singh, K.-H. Zeng, K. Hatch, K. Hsu, L. Itti, L. Y. Chen, L. Pinto, L. Fei-Fei, L. Tan, L. J. Fan, L. Ott, L. Lee, L. Weihs, M. Chen, M. Lepert, M. Memmel, M. Tomizuka, M. Itkina, M. G. Castro, M. Spero, M. Du, M. Ahn, M. C. Yip, M. Zhang, M. Ding, M. Heo, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf, N. Liu, N. D. Palo, N. M. M. Shafiqullah, O. Mees, O. Kroemer, O. Bastani, P. R. Sanketi, P. T. Miller, P. Yin, P. Wohlhart, P. Xu, P. D. Fagan, P. Mitran, P. Sermanet, P. Abbeel, P. Sundaresan, Q. Chen, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Mart’ in-Mart’ in, R. Bajjal, R. Scalise, R. Hendrix, R. Lin, R. Qian, R. Zhang, R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Lin, S. Moore, S. Bahl, S. Dass, S. Sonawani, S. Song, S. Xu, S. Haldar, S. Karamcheti, S. Adebola, S. Guist, S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Ramamoorthy, S. Dasari, S. Belkale, S. Park, S. Nair, S. Mirchandani, T. Osa, T. Gupta, T. Harada, T. Matsushima, T. Xiao, T. Kollar, T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, T. Chung, V. Jain, V. Vanhoucke, W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Chen, X. Wang, X. Zhu, X. Geng, X. Liu, X. Liangwei, X. Li, Y. Pang, Y. Lu, Y. J. Ma, Y. Kim, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Wu, Y. Xu, Y. Wang, Y. Bisk, Y. Dou, Y. Cho, Y. Lee, Y. Cui, Y. Cao, Y.-H. Wu, Y. Tang, Y. Zhu, Y. Zhang, Y. Jiang, Y. Li, Y. Li, Y. Iwasawa, Y. Matsuo, Z. Ma, Z. Xu, Z. J. Cui, Z. Zhang, Z. Fu, and Z. Lin, “Open X-Embodiment: Robotic learning datasets and RT-X models,” 2024.
- [3] A. Khazatsky, K. Pertsch, S. Nair, A. Balakrishna, S. Dasari, S. Karamcheti, S. Nasiriany, M. K. Srirama, L. Y. Chen, K. Ellis, P. D. Fagan, J. Hejna, M. Itkina, M. Lepert, Y. J. Ma, P. T. Miller, J. Wu, S. Belkale, S. Dass, H. Ha, A. Jain, A. Lee, Y. Lee, M. Memmel, S. Park, I. Radosavovic, K. Wang, A. Zhan, K. Black, C. Chi, K. B. Hatch, S. Lin, J. Lu, J. Mercat, A. Rehman, P. R. Sanketi, A. Sharma, C. Simpson, Q. Vuong, H. R. Walke, B. Wulfe, T. Xiao, J. H. Yang, A. Yavary, T. Z. Zhao, C. Agia, R. Bajjal, M. G. Castro, D. Chen, Q. Chen, T. Chung, J. Drake, E. P. Foster, J. Gao, D. A. Herrera, M. Heo, K. Hsu, J. Hu, D. Jackson, C. Le, Y. Li, K. Lin, R. Lin, Z. Ma, A. Maddukuri, S. Mirchandani, D. Morton, T. Nguyen, A. O’Neill, R. Scalise, D. Seale, V. Son, S. Tian, E. Tran, A. E. Wang, Y. Wu, A. Xie, J. Yang, P. Yin, Y. Zhang, O. Bastani, G. Berseth, J. Bohg, K. Goldberg, A. Gupta, A. Gupta, D. Jayaraman, J. J. Lim, J. Malik, R. Martín-Martín, S. Ramamoorthy, D. Sadigh, S. Song, J. Wu, M. C. Yip, Y. Zhu, T. Kollar, S. Levine, and C. Finn, “Droid: A large-scale in-the-wild robot manipulation dataset,” 2024.
- [4] K. Black, M. Nakamoto, P. Atreya, H. Walke, C. Finn, A. Kumar, and S. Levine, “Zero-shot robotic manipulation with pretrained image-editing diffusion models,” *arXiv preprint arXiv:2310.10639*, 2023.
- [5] Y. Du, S. Yang, B. Dai, H. Dai, O. Nachum, J. Tenenbaum, D. Schurman, and P. Abbeel, “Learning universal policies via text-guided video generation,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [6] I. Kapelyukh, V. Vosylius, and E. Johns, “Dall-e-bot: Introducing web-scale diffusion models to robotics,” *IEEE Robotics and Automation Letters*, 2023.
- [7] Y. Du, M. Yang, P. Florence, F. Xia, A. Wahid, B. Ichter, P. Sermanet, T. Yu, P. Abbeel, J. B. Tenenbaum *et al.*, “Video language planning,” *arXiv preprint arXiv:2310.10625*, 2023.
- [8] A. Ajay, S. Han, Y. Du, S. Li, A. Gupta, T. Jaakkola, J. Tenenbaum, L. Kaelbling, A. Srivastava, and P. Agrawal, “Compositional foundation models for hierarchical planning,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [9] J. Gao, K. Hu, G. Xu, and H. Xu, “Can pre-trained text-to-image models generate visual goals for reinforcement learning?” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [10] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard, “Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks,” in *IEEE Robotics and Automation Letters (RAL)*, 2021.
- [11] H. Walke, K. Black, A. Lee, M. J. Kim, M. Du, C. Zheng, T. Zhao, P. Hansen-Estruch, Q. Vuong, A. He, V. Myers, K. Fang, C. Finn, and S. Levine, “Bridgedata v2: A dataset for robot learning at scale,” in *Conference on Robot Learning (CoRL)*, 2023.
- [12] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics,” in *International conference on machine learning*. PMLR, 2015, pp. 2256–2265.
- [13] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [15] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu *et al.*, “Rt-1: Robotics transformer for real-world control at scale,” *arXiv preprint arXiv:2212.06817*, 2022.
- [16] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song, “Diffusion policy: Visuomotor policy learning via action diffusion,” *arXiv preprint arXiv:2303.04137*, 2023.
- [17] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choro-manski, T. Ding, D. Driess, A. Dubey, C. Finn *et al.*, “Rt-2: Vision-language-action models transfer web knowledge to robotic control,” *arXiv preprint arXiv:2307.15818*, 2023.
- [18] Octo Model Team, D. Ghosh, H. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, C. Xu, J. Luo, T. Kreiman, Y. Tan, L. Y. Chen, P. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine, “Octo: An open-source generalist robot policy,” in *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024.
- [19] R. Doshi, H. Walke, O. Mees, S. Dasari, and S. Levine, “Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation,” in *Conference on Robot Learning*, 2024.
- [20] M. Zawalski, W. Chen, K. Pertsch, O. Mees, C. Finn, and S. Levine, “Robotic control via embodied chain-of-thought reasoning,” in *Conference on Robot Learning*, 2024.
- [21] Z. Mandi, H. Bharadhwaj, V. Moens, S. Song, A. Rajeswaran, and V. Kumar, “Cacti: A framework for scalable multi-task multi-scene visual imitation learning,” *arXiv preprint arXiv:2212.05711*, 2022.
- [22] Z. Chen, S. Kiani, A. Gupta, and V. Kumar, “Genaug: Retargeting behaviors to unseen situations via generative augmentation,” *arXiv preprint arXiv:2302.06671*, 2023.
- [23] T. Yu, T. Xiao, A. Stone, J. Tompson, A. Brohan, S. Wang, J. Singh, C. Tan, J. Peralta, B. Ichter *et al.*, “Scaling robot learning with

- semantically imagined experience,” *arXiv preprint arXiv:2302.11550*, 2023.
- [24] A. Stone, T. Xiao, Y. Lu, K. Gopalakrishnan, K.-H. Lee, Q. Vuong, P. Wohlhart, S. Kirmani, B. Zitkovich, F. Xia *et al.*, “Open-world object manipulation using pre-trained vision-language models,” *arXiv preprint arXiv:2303.00905*, 2023.
- [25] A. Peng, I. Sucholutsky, B. Z. Li, T. R. Sumers, T. L. Griffiths, J. Andreas, and J. A. Shah, “Learning with language-guided state abstractions,” *arXiv preprint arXiv:2402.18759*, 2024.
- [26] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch, “Language models as zero-shot planners: Extracting actionable knowledge for embodied agents,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 9118–9147.
- [27] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, A. Zeng, J. Tompson, I. Mordatch, Y. Chebotar *et al.*, “Inner monologue: Embodied reasoning through planning with language models,” *arXiv preprint arXiv:2207.05608*, 2022.
- [28] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian *et al.*, “Do as i can, not as i say: Grounding language in robotic affordances,” in *Conference on robot learning*. PMLR, 2023, pp. 287–318.
- [29] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg, “Text2motion: From natural language instructions to feasible plans,” *Autonomous Robots*, vol. 47, no. 8, pp. 1345–1365, 2023.
- [30] Z. Wang, S. Cai, G. Chen, A. Liu, X. Ma, and Y. Liang, “Describe, explain, plan and select: Interactive planning with large language models enables open-world multi-task agents,” *arXiv preprint arXiv:2302.01560*, 2023.
- [31] S. Fujimoto, D. Meger, and D. Precup, “Off-policy deep reinforcement learning without exploration,” in *International conference on machine learning*. PMLR, 2019, pp. 2052–2062.
- [32] S. K. S. Ghasemipour, D. Schuurmans, and S. S. Gu, “Emaq: Expected-max q-learning operator for simple yet effective offline and online rl,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 3682–3691.
- [33] H. Chen, C. Lu, C. Ying, H. Su, and J. Zhu, “Offline reinforcement learning via high-fidelity generative behavior modeling,” *arXiv preprint arXiv:2209.14548*, 2022.
- [34] P. Hansen-Estruch, I. Kostrikov, M. Janner, J. G. Kuba, and S. Levine, “Idql: Implicit q-learning as an actor-critic method with diffusion policies,” *arXiv preprint arXiv:2304.10573*, 2023.
- [35] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano *et al.*, “Training verifiers to solve math word problems,” *arXiv preprint arXiv:2110.14168*, 2021.
- [36] H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe, “Let’s verify step by step,” *arXiv preprint arXiv:2305.20050*, 2023.
- [37] A. Hosseini, X. Yuan, N. Malkin, A. Courville, A. Sordoni, and R. Agarwal, “V-star: Training verifiers for self-taught reasoners,” *arXiv preprint arXiv:2402.06457*, 2024.
- [38] W. Liu, Y. Du, T. Hermans, S. Chernova, and C. Paxton, “Structdiffusion: Language-guided creation of physically-valid structures using unseen objects,” *arXiv preprint arXiv:2211.04604*, 2022.
- [39] W. Huang, F. Xia, D. Shah, D. Driess, A. Zeng, Y. Lu, P. Florence, I. Mordatch, S. Levine, K. Hausman *et al.*, “Grounded decoding: Guiding text generation with grounded models for robot control,” *arXiv preprint arXiv:2303.00855*, 2023.
- [40] A. Z. Ren, J. Clark, A. Dixit, M. Itkina, A. Majumdar, and D. Sadigh, “Explore until confident: Efficient exploration for embodied question answering,” in *Robotics Science and Systems (RSS)*, 2024.
- [41] “Robots that ask for help: Uncertainty alignment for large language model planners,” *arXiv preprint arXiv:2307.01928*, 2023.
- [42] S. Nair, E. Mitchell, K. Chen, B. Ichter, S. Savarese, and C. Finn, “Learning language-conditioned robot behavior from offline data and crowd-sourced annotation,” *Conference on Robot Learning (CoRL)*, 2021.
- [43] L. P. Kaelbling, “Learning to achieve goals,” in *IJCAI*, vol. 2. Citeseer, 1993, pp. 1094–8.
- [44] T. Schaul, D. Horgan, K. Gregor, and D. Silver, “Universal value function approximators,” in *International conference on machine learning*. PMLR, 2015, pp. 1312–1320.
- [45] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, O. Pieter Abbeel, and W. Zaremba, “Hindsight experience replay,” *Advances in neural information processing systems*, vol. 30, 2017.
- [46] S. Tellex, N. Gopalan, H. Kress-Gazit, and C. Matuszek, “Robots that use language,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 3, pp. 25–55, 2020.
- [47] S. Stepputtis, J. Campbell, M. Phielipp, S. Lee, C. Baral, and H. Ben Amor, “Language-conditioned imitation learning for robot manipulation tasks,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 13 139–13 150, 2020.
- [48] O. Mees, L. Hermann, and W. Burgard, “What matters in language conditioned robotic imitation learning over unstructured data,” *IEEE Robotics and Automation Letters (RA-L)*, vol. 7, no. 4, pp. 11 205–11 212, 2022.
- [49] O. Mees, J. Borja-Diaz, and W. Burgard, “Grounding language with visual affordances over unstructured data,” in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, London, UK, 2023.
- [50] C. Lynch and P. Sermanet, “Language conditioned imitation learning over unstructured data,” *arXiv preprint arXiv:2005.07648*, 2020.
- [51] A. Mandlekar, F. Ramos, B. Boots, S. Savarese, L. Fei-Fei, A. Garg, and D. Fox, “Iris: Implicit reinforcement without interaction at scale for learning control from offline robot manipulation data,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 4414–4420.
- [52] S. Park, D. Ghosh, B. Eysenbach, and S. Levine, “Hlq: Offline goal-conditioned rl with latent states as actions,” *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [53] R. S. Sutton, D. Precup, and S. Singh, “Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning,” *Artificial intelligence*, vol. 112, no. 1-2, pp. 181–211, 1999.
- [54] P.-L. Bacon, J. Harb, and D. Precup, “The option-critic architecture,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 31, no. 1, 2017.
- [55] J. Schmidhuber, “Learning to generate sub-goals for action sequences,” in *Artificial neural networks*, 1991, pp. 967–972.
- [56] P. Dayan and G. E. Hinton, “Feudal reinforcement learning,” *Advances in neural information processing systems*, vol. 5, 1992.
- [57] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum, “Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation,” *Advances in neural information processing systems*, vol. 29, 2016.
- [58] A. S. Vechnets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu, “Feudal networks for hierarchical reinforcement learning,” in *International conference on machine learning*. PMLR, 2017, pp. 3540–3549.
- [59] A. Levy, G. Konidaris, R. Platt, and K. Saenko, “Learning multi-level hierarchies with hindsight,” *arXiv preprint arXiv:1712.00948*, 2017.
- [60] O. Nachum, S. S. Gu, H. Lee, and S. Levine, “Data-efficient hierarchical reinforcement learning,” *Advances in neural information processing systems*, vol. 31, 2018.
- [61] O. Nachum, S. Gu, H. Lee, and S. Levine, “Near-optimal representation learning for hierarchical reinforcement learning,” *arXiv preprint arXiv:1810.01257*, 2018.
- [62] A. Gupta, V. Kumar, C. Lynch, S. Levine, and K. Hausman, “Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning,” *arXiv preprint arXiv:1910.11956*, 2019.
- [63] A. Ajay, A. Kumar, P. Agrawal, S. Levine, and O. Nachum, “Opal: Offline primitive discovery for accelerating offline reinforcement learning,” *arXiv preprint arXiv:2010.13611*, 2020.
- [64] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet, “Learning latent plans from play,” in *Conference on Robot Learning (CoRL)*. PMLR, 2020, pp. 1113–1132.
- [65] E. Rosete-Beas, O. Mees, G. Kalweit, J. Boedecker, and W. Burgard, “Latent plans for task-agnostic offline reinforcement learning,” in *Conference on Robot Learning*. PMLR, 2023, pp. 1838–1849.
- [66] T. Zhang, S. Guo, T. Tan, X. Hu, and F. Chen, “Generating adjacency-constrained subgoals in hierarchical reinforcement learning,” *Advances in neural information processing systems*, vol. 33, pp. 21 579–21 590, 2020.
- [67] K. Pertsch, Y. Lee, and J. Lim, “Accelerating reinforcement learning with learned skill priors,” in *Conference on robot learning*. PMLR, 2021, pp. 188–204.
- [68] E. Chane-Sane, C. Schmid, and I. Laptev, “Goal-conditioned reinforcement learning with imagined subgoals,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 1430–1440.
- [69] N. Savinov, A. Dosovitskiy, and V. Koltun, “Semi-parametric topological memory for navigation,” *arXiv preprint arXiv:1803.00653*, 2018.

- [70] B. Eysenbach, R. R. Salakhutdinov, and S. Levine, "Search on the replay buffer: Bridging planning and reinforcement learning," *Advances in neural information processing systems*, vol. 32, 2019.
- [71] S. Nair and C. Finn, "Hierarchical foresight: Self-supervised learning of long-horizon tasks via visual subgoal generation," *arXiv preprint arXiv:1909.05829*, 2019.
- [72] S. Nasiriany, V. Pong, S. Lin, and S. Levine, "Planning with goal-conditioned policies," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [73] Z. Huang, F. Liu, and H. Su, "Mapping state space using landmarks for universal goal reaching," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [74] C. Hoang, S. Sohn, J. Choi, W. Carvalho, and H. Lee, "Successor feature landmarks for long-horizon goal-conditioned reinforcement learning," *Advances in neural information processing systems*, vol. 34, pp. 26 963–26 975, 2021.
- [75] J. Kim, Y. Seo, and J. Shin, "Landmark-guided subgoal generation in hierarchical reinforcement learning," *Advances in neural information processing systems*, vol. 34, pp. 28 336–28 349, 2021.
- [76] L. Zhang, G. Yang, and B. C. Stadie, "World model as a graph: Learning latent landmarks for planning," in *International conference on machine learning*. PMLR, 2021, pp. 12 611–12 620.
- [77] D. Shah, B. Eysenbach, G. Kahn, N. Rhinehart, and S. Levine, "Rapid exploration for open-world navigation with latent goal models," *arXiv preprint arXiv:2104.05859*, 2021.
- [78] K. Fang, P. Yin, A. Nair, and S. Levine, "Planning to practice: Efficient online fine-tuning by composing goals in latent space," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 4076–4083.
- [79] J. Li, C. Tang, M. Tomizuka, and W. Zhan, "Hierarchical planning through goal-conditioned offline reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 10 216–10 223, 2022.
- [80] J. Kim, Y. Seo, S. Ahn, K. Son, and J. Shin, "Imitating graph-based planning with goal-conditioned policies," *arXiv preprint arXiv:2303.11166*, 2023.
- [81] K. Fang, P. Yin, A. Nair, H. R. Walke, G. Yan, and S. Levine, "Generalization with lossy affordances: Leveraging broad offline data for learning visuomotor tasks," in *Conference on Robot Learning*. PMLR, 2023, pp. 106–117.
- [82] T. Brooks, A. Holynski, and A. A. Efros, "Instructpix2pix: Learning to follow image editing instructions," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [83] J. Xing, M. Xia, Y. Zhang, H. Chen, W. Yu, H. Liu, X. Wang, T.-T. Wong, and Y. Shan, "Dynamicrafter: Animating open-domain images with video diffusion priors," *arXiv preprint arXiv:2310.12190*, 2023.
- [84] J. Ho and T. Salimans, "Classifier-free diffusion guidance," *arXiv preprint arXiv:2207.12598*, 2022.
- [85] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," *arXiv preprint arXiv:2011.13456*, 2020.
- [86] P. Dhariwal and A. Nichol, "Diffusion models beat gans on image synthesis," *Advances in neural information processing systems*, vol. 34, pp. 8780–8794, 2021.
- [87] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [88] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel, "Domain randomization for transferring deep neural networks from simulation to the real world," *International Conference on Intelligent Robots and Systems*, 2017.
- [89] C. Zheng, B. Eysenbach, H. Walke, P. Yin, K. Fang, R. Salakhutdinov, and S. Levine, "Stabilizing contrastive rl: Techniques for offline goal reaching," *arXiv preprint arXiv:2306.03346*, 2023.
- [90] R. Goyal, S. E. Kahou, V. Michalski, J. Materzynska, S. Westphal, H. Kim, V. Haenel, I. Fruend, P. Yianilos, M. Mueller-Freitag, and et al., "The" something something" video database for learning and evaluating visual common sense," in *IEEE international conference on computer vision (ICCV)*, 2017.
- [91] M. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine, P. Liang, and C. Finn, "Openvla: An open-source vision-language-action model," *arXiv preprint arXiv:2406.09246*, 2024.
- [92] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "Film: Visual reasoning with a general conditioning layer," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, no. 1, 2018.

A. Classifier Training

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

$$\mathcal{J}(\theta) = \mathbb{E}_{(s,g,l) \sim \mathcal{D}_{l,a}} [\log(f_{\theta}(s,g,l))] + \mathbb{E}_{(s^-,g^-,l^-) \sim \mathcal{D}_{l,a}^-} [\log(1 - f_{\theta}(s^-,g^-,l^-))]. \quad (2)$$

where $\mathcal{D}_{l,a}$ is the language-annotated dataset that consists of trajectory and language task pairs, and N is a function for generating negative examples from the dataset. Given a dataset $\mathcal{D}_{l,a}$, N generates negatives from $\mathcal{D}_{l,a}$ in the following ways:

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

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

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

Network architecture and training hyperparameters: The classifier network architecture consists of a ResNet-34 encoder from [11], followed by a two-layer MLP with layers of dimension 256. Separate encoders are used to encode the image observations and the goal images (parameters are not shared between the two). Both of these encoders use FiLM conditioning [92] after each residual block to condition on the language instruction. Classifier networks are trained using a learning rate of 3×10^{-4} and a batch size of 256 for 100,000 gradient steps. A dropout rate of 0.1 is used.

B. Image Augmentations

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

- 1) Random Resized Crop:
 - scale: (0.8, 1.0)
 - ratio: (0.9, 1.1)
- 2) Random Brightness Shift:
 - shift ratio: 0.2
- 3) Random Contrast:
 - Contrast range: (0.8, 1.2)
- 4) Random Saturation:
 - Saturation range: (0.8, 1.2)
- 5) Random Hue:
 - shift ratio: 0.1

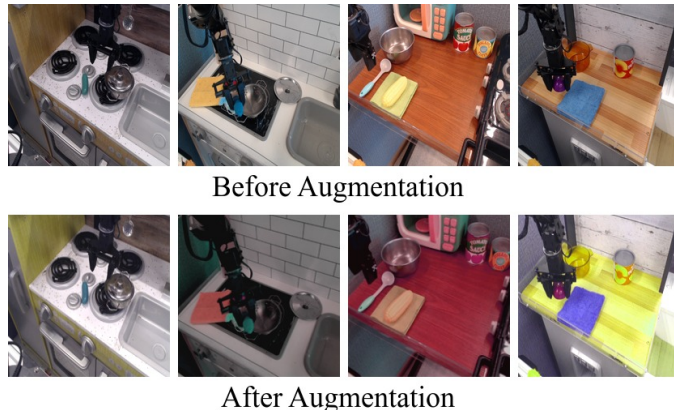


Fig. 4: **Image augmentation examples** Examples of images from the Bridge dataset before and after having the image augmentations applied to them that are used during policy and classifier training.

Figure 4 visualizes examples from the Bridge dataset before and after augmentations are applied.

C. Qualitative Analysis

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

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

Fig. 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.

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

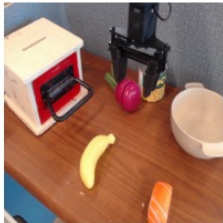
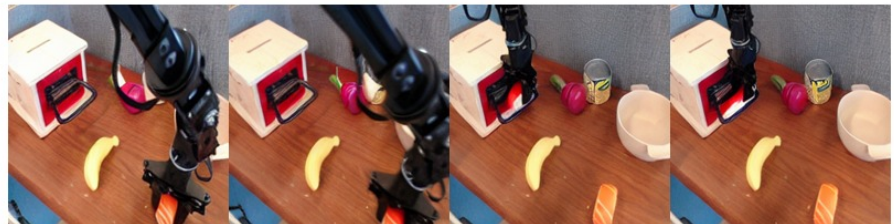


Image Observation



1st 2nd 3rd 4th



5th 6th 7th 8th

Generated goal images ranked by the classifier

(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."

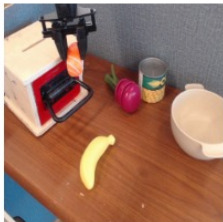


Image Observation



1st 2nd 3rd 4th



5th 6th 7th 8th

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



1st

2nd

3rd

4th



5th

6th

7th

8th

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."

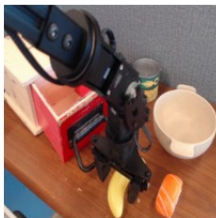


Image Observation



1st

2nd

3rd

4th



5th

6th

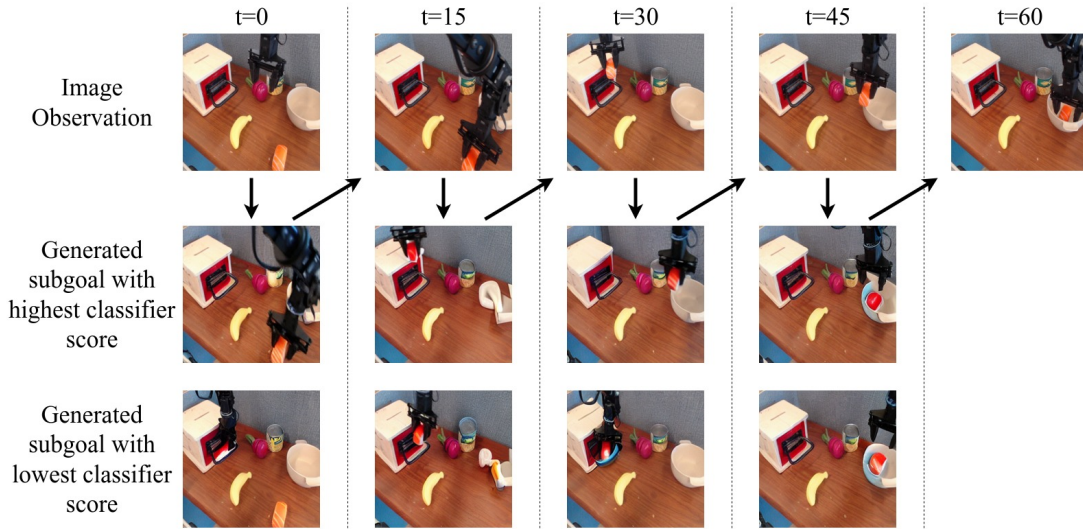
7th

8th

Generated goal images ranked by the classifier

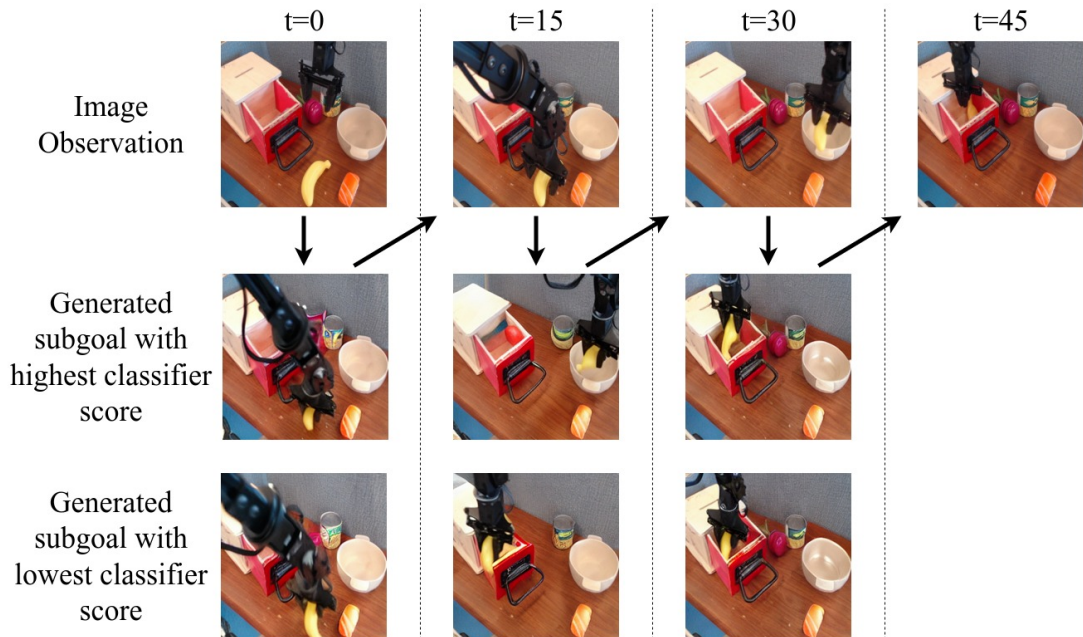
(d) The classifier incorrectly ranks the subgoal images higher where the robot is placing the banana into the bowl than it ranks the subgoal images where the robot is placing the banana into the drawer. This could be due to there being a strong bias for placing objects in bowls in the Bridge V2 training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

Fig. 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.

3) *Qualitative Analysis of Augmentation De-synchronization*: We see that when applying augmentation de-synchronization, the number of failures due to low-level policy errors (missed grasps, dropping held objects, etc.) decreases, indicating that augmentation de-synchronization is important for the low-level policy to be able to correctly interpret and follow the subgoal images generated by the video prediction model. This is particularly important in domains where there is a large visual generalization gap between the training data and the evaluation tasks. For example, in the CALVIN benchmark, the colors and shapes of objects differ between the training and evaluation scenes. This difference causes the subgoals generated by the video prediction model to often contain objects with incorrect shapes and colors (Figure 7). Augmentation de-synchronization seems to be critical to allowing the low-level policy to be robust to these hallucinations and artifacts.

D. Additional Ablation Experiments

1) *Effect of Augmentation Desynchronization*: We ablate the different components of GHIL-Glue when applied to SuSIE in the CALVIN benchmark (Table III). Removing the augmentation desynchronization from only the low-level policy results in similar performance to base-SuSIE and GHIL-Glue (SuSIE) with the augmentation desynchronization removed from both the low-level policy and the classifier. This suggests that the low-level policy performance of SuSIE without augmentation desynchronization is a significant bottleneck for SuSIE—even when selecting better subgoals via the use of filtering, performance cannot increase if the low-level policy cannot reliably reach those goals. Conversely, removing the augmentation desynchronization from only the classifier results in similar performance to GHIL-Glue (SuSIE) without subgoal filtering. This suggests that, like the low-level policy, augmentation desynchronization is important for the classifier network to correctly perform its function in GHIL-Glue (SuSIE).

Method	Tasks completed in a row					Avg. Len.
	1	2	3	4	5	
SuSIE	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) - w/o Aug De-sync on policy	91.5%	74.2%	56.0%	41.2%	29.8%	2.93
GHIL-Glue (SuSIE) - w/o Aug De-sync on classifier	95.0%	86.2%	70.0%	57.8%	47.0%	3.56
GHIL-Glue (SuSIE)	95.2%	88.5%	73.2%	62.5%	49.8%	3.69

TABLE III: **Effect of Augmentation Desynchronization in GHIL-Glue (SuSIE)** Success rates on the validation tasks from the D environment of the CALVIN Challenge averaged across 4 random seeds. Results are shown comparing the performance of SuSIE, GHIL-Glue (SuSIE), and ablations of GHIL-Glue (SuSIE). *GHIL-Glue (SuSIE) - Aug De-sync Only* is GHIL-Glue without applying subgoal filtering, *GHIL-Glue (SuSIE) - Subgoal Filtering Only* is GHIL-Glue without applying augmentation de-synchronization to either the low-level policy or the subgoal classifier, *GHIL-Glue (SuSIE) - w/o Aug De-sync on policy* is GHIL-Glue without applying augmentation de-synchronization on the low-level policy, and *GHIL-Glue (SuSIE) - w/o Aug De-sync on classifier* is GHIL-Glue without applying augmentation de-synchronization on the subgoal classifier.

2) *Number of Candidate Subgoals*: We conduct an ablation over the number of candidate subgoals used for subgoal filtering in GHIL-Glue (SuSIE) in the CALVIN benchmark (Table IV). We find that GHIL-Glue (SuSIE) achieves similar performance whether 4, 8, or 16 candidate subgoals are used. In our main results (Section V-C), we report the performance of GHIL-Glue (SuSIE) on the CALVIN benchmark when using 8 candidate subgoals for filtering. For GHIL-Glue (UniPi) on the CALVIN benchmark, we use 4 candidate subgoals for filtering, due to the increased computation burden of generating video subgoals with the UniPi video model vs. generating image subgoals with the SuSIE image model. In our physical experiments, we run GHIL-Glue (SuSIE) using 4 candidate subgoals for filtering.

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

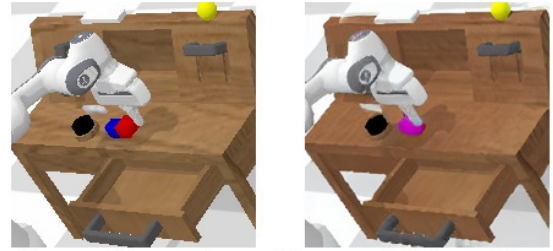


Image Observation Generated Subgoal Image

Fig. 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.

Method	Tasks completed in a row					Avg. Len.
	1	2	3	4	5	
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 IV: **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.