

# Learning Factorized Diffusion Policies for Conditional Action Diffusion

Omkar Patil<sup>1</sup>, Eric Rosen, and Nakul Gopalan<sup>2</sup>

<sup>1,2</sup>SCAI, Arizona State University

**Abstract**—Diffusion models have emerged as a promising choice for learning robot skills from demonstrations. However, they face three problems: diffusion models are not sample-efficient, data is expensive to collect in robotics, and the space of tasks is combinatorially large. The established way of learning diffusion policies has little room to accommodate solutions for the aforementioned challenges, in addition to scaling the model size and paired observation-action data. In this work, we propose a novel method for training diffusion models termed ‘Composable Diffusion Guidance’ CoDiG to compositionally learn diffusion policies for robot skills with respect to different observation modalities, such as proprioception, vision and tactile. CoDiG provides more flexibility to deal with such observational modalities, leading to sample-efficiency gains over 20% in applicable tasks. CoDiG opens up more avenues for research on foundation models as it ameliorates the requirement of scaling all the observational modalities together during data collection.

## I. INTRODUCTION

Diffusion models have emerged as a promising choice for learning robot skills from demonstrations [4]. However, diffusion models are not sample efficient and require the collection of large amounts of demonstration data. Unfortunately, data collection in robotics is not easy due to the coupling of various modalities such as robot actions, vision, tactile, text etc. Previous work [4] has used diffusion models to learn the conditional distribution of action with respect to visual, proprioceptive, and other observations. However, this forces the collection of paired action, visual, tactile, and other modality data as demonstrations, which can be restrictive and difficult to scale. Instead, we propose a novel method CoDiG ‘Composable Diffusion Guidance’ utilizing Bayes’ theorem to split the distribution of the robot actions conditioned on different observation modalities. For instance, the existing way to train visuomotor diffusion policies is to learn a conditional distribution of actions with respect to the visual and proprioceptive observations. Instead, CoDiG learns a diffusion ‘motion-model’ on the action distribution of skills and a visual-‘guidance model’, that when composed with the motion-model results in generation of the action conditioned on the visual observations.

opatil3@asu.edu

Critically, CoDiG enables the training of models conditioned on an additional observational modality with respect to the rest, so that data collection efforts can be concentrated on the cheaper modality. CoDiG does not require manual adjustment of compositional weights such as Wang et al. [25]. Our preliminary experiments on learning visuomotor policies for selected tasks show that CoDiG is highly sample efficient and robust to visual changes in the environment, where current diffusion policies drastically fail. Our contributions are as follows.

- We present a novel framework for training diffusion models on robot demonstration data called CoDiG, abbreviated for ‘Composable Diffusion Guidance’. CoDiG decouples the observational modalities, enabling residual learning of an additional modality with respect to the rest.
- Our results show that visuomotor CoDiG performs favorably with respect to the conventional diffusion models on Adroit hand manipulation environments [10] and several RLBench [14] tasks. Further experimentation is required to validate the efficacy and robustness of CoDiG for long horizon tasks and precise manipulation.
- Existing diffusion policies drastically fail with small variations in the visual scene from the training data distribution. Visuomotor CoDiG is relatively robust to visual adaptation and responds stronger to additional demonstrations collected in the new visual scene.

## II. METHODOLOGY

### A. Theoretical Results

Assume that we have robot demonstrations  $D = \{(\mathbf{x}, \mathbf{y})_i\}$  where  $i = 1..N$ , consisting of actions  $\mathbf{x}$  and different observation modalities  $\mathbf{y}^k$  such as camera images, task description and proprioception data. We are interested in learning  $p(\mathbf{x}|\mathbf{y})$  from the data, such that given a task description, current camera images, state of the robot and other observations, we can sample an action  $\mathbf{x}$  with a high likelihood in the demonstration data distribution.

Most treatments of diffusion models have considered distributions of a single entity such as images [13, 16, 20]. This formulation has been directly adopted by the robotics community [4] leading to the popular optimization objective shown in Equation 1.

$$\begin{aligned}\mathcal{L}_t(\theta) &= \mathbb{E}_{q(\mathbf{x}_0, \mathbf{y}) \mathcal{N}(\epsilon_0; 0, \mathcal{I})} [\|\epsilon_0 - \hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)\|_2^2] \\ \mathbf{x}_{t-1} &\sim \mathcal{N}\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)\right), \sqrt{1 - \alpha_t} \mathcal{I}\right)\end{aligned}\quad (1)$$

Here, the network  $\epsilon_\theta$  is also parametrized with  $\mathbf{y}$  for the conditional prediction of the noise added to the action  $\mathbf{x}$ . Once the network  $\hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)$  is trained, ancestral sampling shown in Equation 2 can be used to sample actions given the observations  $\mathbf{y}$ . However, while prior work assumes that  $\epsilon_\theta$  learns the score of the conditional distribution  $p(\mathbf{x}|\mathbf{y})$ , there is no formal proof provided for the same. Hence, we present our first result, where we show that a conditional diffusion process as defined by [5] does indeed result in the loss of Equation 1 [4] being a maximizer of the evidence lower bound (ELBO) of the log-likelihood of the conditional data distribution  $\log q(\mathbf{x}|\mathbf{y})$ .

**Theorem II.1.** *The diffusion loss function  $\mathcal{L}_t(\theta)$  as defined in Equation 1, in expectation over the time-steps  $1 \leq t \leq T$ , is equivalent to maximizing the ELBO of the log-likelihood of the conditional data distribution  $\log q(\mathbf{x}|\mathbf{y})$ , under a conditional Markovian noising process  $\hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1})$  and the reverse transition kernel as  $\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$ .*

The proof for Theorem II.1 can be found in the Appendix B-A. We now proceed to our main result. We would like to decouple different observational modalities  $\mathbf{y}^k$ ,  $1 \leq k \leq M$ , so that we can prioritize their collection based on cost. Assume that we have trained a diffusion policy with  $k$  observational modalities and we collect  $N$  additional demonstrations with the  $(k+1)^{th}$  modality in conjunction with the existing  $k$ . The score of the new conditional distribution  $p(\mathbf{x}|\mathbf{y}^{1:k}, \mathbf{y}^{k+1})$ , where  $\mathbf{y}^{1:k} \equiv \mathbf{y}^1, \dots, \mathbf{y}^k$ , can be written using Bayes's theorem along the lines of Equation 18 at diffusion time-step  $t$  as follows:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}^{1:k}, \mathbf{y}^{k+1}; \theta, \phi) \quad (3)$$

$$= \nabla_{\mathbf{x}_t} \log p(\mathbf{y}^{k+1}|\mathbf{x}_t, \mathbf{y}^{1:k}; \phi) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}^{1:k}; \theta) \quad (4)$$

Here  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}^{1:k}; \theta)$  is the score of the model trained on the original  $k$  observational modalities, while  $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}^{k+1}|\mathbf{x}_t, \mathbf{y}^{1:k}; \phi)$  corresponds to the score of the classifier for the modality  $\mathbf{y}^{k+1}$ . The method proposed by classifier guided-diffusion [5] to explicitly train a classifier  $p(\mathbf{y}^{k+1}|\mathbf{x}_t, \mathbf{y}^{1:k})$  on the noisy samples of  $\mathbf{x}_t$  and  $\mathbf{y}^{1:k}$  does not apply due to continuous and high-dimensional observation modalities such as images, rather than simple classes such as cats or dogs.

The central idea of this work is to instead directly parametrize the gradient of the classifier  $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}^{k+1}|\mathbf{x}_t, \mathbf{y}^{1:k}; \phi)$  using a neural network, which we refer to as the guidance model  $\pi_g$ , rather than to learn a classifier and then obtain its gradients. To learn the guidance model, we minimize the Fisher divergence between the guidance model and the true score of the classifier.

$$\begin{aligned}D_F(p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})||p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) &= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \\ &\left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2\right]\end{aligned}\quad (5)$$

Here, observation modalities  $\mathbf{y}^{1:k+1}$  can be noised with a Gaussian kernel  $\mathcal{N}(\tilde{\mathbf{y}}; \mathbf{y}, \tau^2 \mathcal{I})$  of variance  $\tau^2$  that is small enough such that  $p_\tau(\tilde{\mathbf{y}}^i) \approx p(\mathbf{y}^i)$ . Robot action  $\mathbf{x}$  is noised with the diffusion transition kernel of  $\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}, (1 - \alpha_t) \mathcal{I})$ . However, Equation 5 is difficult to use in its current form as estimation of the true score is difficult for large datasets. Chao et al. [3] derive the denoising likelihood score matching (DLSM) objective for conditional distributions, which forms the basis of our next result.

**Theorem II.2.** *Denoising score matching for the guidance model  $\pi_g$ :  $\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})$ , as expressed in Equation 5 is equal up to a constant to the following loss:*

$$\begin{aligned}L_{DLSM}^t(\phi) &= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\right. \\ &\quad \left. + \nabla_{\mathbf{x}_t} \log p_\theta(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_\alpha(\mathbf{x}_t|\mathbf{x})\|_2^2\right]\end{aligned}\quad (6)$$

For diffusion models, we take the forward transition kernel as  $p_\alpha(\mathbf{x}_t|\mathbf{x}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}, (1 - \alpha_t) \mathcal{I})$ . We noise the observations  $\mathbf{y}^{1:k+1}$  with a Gaussian kernel  $\mathcal{N}(\tilde{\mathbf{y}}; \mathbf{y}, \tau^2 \mathcal{I})$  of variance  $\tau^2$  that is small enough such that  $p_\tau(\tilde{\mathbf{y}}^i) \approx p(\mathbf{y}^i)$ .

The proof for Theorem II.2 is a conditional variant of the one derived by [3], and is presented in Appendix B-B. Equation 6 can now be used in practice to learn the guidance model. Simplifying  $\nabla_{\mathbf{x}_t} \log p_\alpha(\mathbf{x}_t|\mathbf{x})$  to  $-\epsilon_0/\sqrt{1 - \alpha_t}$ , where  $\epsilon_0 \sim \mathcal{N}(0, \mathcal{I})$ , we obtain:

$$\begin{aligned}L_{DLSM}(\phi) &= \mathbb{E}_{p_\tau(\mathbf{x}, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} \mathbb{E}_{\epsilon_0 \sim \mathcal{N}(0, \mathcal{I})} \left[\frac{1}{2} \|\hat{\epsilon}_\phi(\tilde{\mathbf{y}}^{k+1}, \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\right. \\ &\quad \left. + \epsilon_\theta(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) - \epsilon_0\|_2^2\right]\end{aligned}\quad (7)$$

Here,  $\epsilon_\theta$  is frozen and known from prior training with observational modalities  $\mathbf{y}^{1:k}$ , while  $\epsilon_\phi$  is learned in expectation over data collected in conjunction with modality  $\mathbf{y}^{k+1}$ . Equation 7 is very similar to the diffusion loss Equation 1, except the model is expected to learn a residual of the noise added to the action sample  $\mathbf{x}_t$  taken as the input. Once,  $\epsilon_\phi$  is learned, actions can be

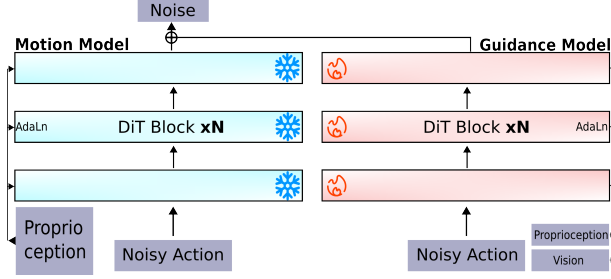


Fig. 1. The simple architecture adds the noise predictions from the two models at each denoising time-step. The motion model is frozen while training the guidance model, which takes an additional observational modality as vision.

sampled from the conditional distribution  $p(\mathbf{x}|\mathbf{y}^{1:k+1})$  using ancestral sampling and Equation 4:

$$\mathbf{x}_{t-1} \sim \mathcal{N}\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon(\mathbf{x}_t, \mathbf{y}^{1:k+1}, t) \right), \sqrt{1 - \alpha_t} \mathbf{I}\right) \quad (8)$$

$$\epsilon(\mathbf{x}_t, \mathbf{y}^{1:k+1}, t) = \hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}^{1:k+1}, t) + \hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}^{1:k}, t) \quad (9)$$

Note that the score of the prior distribution  $\hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}^{1:k}, t)$  or  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}^{1:k}; \theta)$  can further be decomposed with respect to observational modalities as shown above. However, data collection for subsequent modalities  $\mathbf{y}^{k+1}$  must be done in conjunction with existing ones  $\mathbf{y}^{1:k}$ .

### B. Visual CoDiG with Motion-model

We now provide a concrete implementation of the method proposed in Section II-A. Existing diffusion models trained as visuo-motor policies learn the score for the conditional distribution  $p(\mathbf{x}|\mathbf{y}^r, \mathbf{y}^c)$ , where  $\mathbf{y}^r$  corresponds to the proprioceptive observations and  $\mathbf{y}^c$  correspond to the visual observations from cameras. We propose to decouple these observational modalities, and instead learn the scores for a ‘motion-model’  $p(\mathbf{x}|\mathbf{y}^r)$  and a vision ‘guidance-model’  $p(\mathbf{y}^c|\mathbf{x}, \mathbf{y}^r)$ . The equations of diffusion loss for the motion model 1 and the guidance model 7 can be written as follows, and is showcased in Figure 1.

$$\mathcal{L}_{MM}(\theta) = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{y}^r) \sim \mathcal{N}(\epsilon_0, 0, \mathbf{I})} [\|\epsilon_0 - \hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}^r, t)\|_2^2] \quad (10)$$

$$\mathcal{L}_G(\phi) = \mathbb{E}_{p_r(\mathbf{x}, \mathbf{y}^r, \mathbf{y}^c, \tilde{\mathbf{y}}^r, \tilde{\mathbf{y}}^c) \sim \mathbb{E}_{\epsilon_0 \sim \mathcal{N}(0, \mathbf{I})}} \left[ \frac{1}{2} \|\epsilon_0 - \hat{\epsilon}_\phi(\tilde{\mathbf{y}}^c, \mathbf{x}_t, \tilde{\mathbf{y}}^r, t) - \epsilon_\theta(\mathbf{x}_t, \tilde{\mathbf{y}}^r, t)\|_2^2 \right] \quad (11)$$

Instead of having two learned components in Equation 11, we can simplify the learning as follows.

$$\mathcal{L}_G(\phi) = \mathbb{E}_{p_r(\mathbf{x}, \mathbf{y}^r, \mathbf{y}^c, \tilde{\mathbf{y}}^r, \tilde{\mathbf{y}}^c) \sim \mathbb{E}_{\epsilon_0 \sim \mathcal{N}(0, \mathbf{I})}} \left[ \frac{1}{2} \|\epsilon_0 - \hat{\epsilon}_\gamma(\tilde{\mathbf{y}}^c, \mathbf{x}_t, \tilde{\mathbf{y}}^r, t)\|_2^2 \right] \quad (12)$$

$$\hat{\epsilon}_\gamma(\phi, \mathbf{x}_t, \tilde{\mathbf{y}}^r, t) = \epsilon_\theta(\mathbf{x}_t, \tilde{\mathbf{y}}^r, t)$$

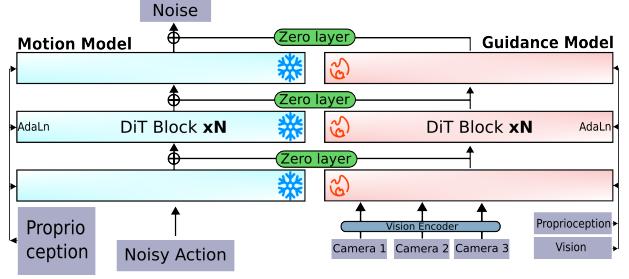


Fig. 2. The improved architecture **CoDiG**, inspired from Controlnet, where the instead of adding the model outputs at the end, we add interrim block outputs from DiT. We pass the visual embeddings into the transformer and condition them on the noisy action using adaLn.

We implement this model using a ControlNet [26] inspired architecture, as shown in Figure 2. Unlike the simple architecture shown in Figure 1, this variant adds the outputs of each DiT block before passing the sum to the consequent block in the motion model. To prevent harmful updates at the start of training, a zero-initialized layer is applied to the guidance model outputs, similar to ControlNet. Moreover, our architecture benefits from attention being applied to the image embeddings, which are in-turn conditioned on the noisy actions in the guidance model.

Note that we expect this approach to have a distinct advantage for visuomotor skills that are heavily dependent on proprioception, as we learn a residual for vision. We hypothesize that the existing way to train visuomotor policies using diffusion models do not learn the right dependencies when presented with all the observational modalities together, leading to lower sample-efficiency.

## III. RESULTS

### A. Environments and Baselines

Since we propose an alternative to the existing way of training diffusion models, our baseline is a model that takes in all the observations and predicts the robot action, learning the conditional distribution. We choose DiT-small ( $\sim 90\text{M}$ ) [17] as our model architecture, which is kept the same for both the baseline and CoDiG. For comparison, we also include UNet [18] implemented by [4] as a part of the baselines. We also include POCO [25], where we compose a motion policy and pre-learned visuomotor policy using a compositional weight of 0.2, as suggested in the paper. Finally, we also showcase the results for classifier-free guidance proposed by Ho and Salimans [11], where we switch out the vision condition with a probability of 0.1. Further details on the baselines, experiments, and the environments used are provided in Appendix C.

TABLE I  
PERFORMANCE RESULTS ACROSS DIFFERENT RL BENCH TASKS FOR 10, 50 AND 100 DEMONSTRATIONS.

Task	Number of Demonstrations	DiT		UNet		DiT-cfg		POCO		CoDiG	
		Mean	Std	Mean	Std	Mean	Std	Mean	Std	Mean	Std
OpenBox	10	21.0	± 2.0	7.5	± 0.7	20.3	± 4.0	27.7	± 4.9	<b>59.0</b>	± 4.0
	50	86.0	± 1.7	86.0	± 1.0	86.0	± 2.6	87.3	± 0.6	<b>90.0</b>	± 2.6
	100	<b>98.3</b>	± 1.5	93.0	± 1.0	92.3	± 2.5	<b>98.3</b>	± 1.5	91.3	± 4.5
CloseBox	10	25.7	± 5.9	30.0	± 6.2	19.7	± 4.7	29.7	± 3.5	<b>71.7</b>	± 6.8
	50	79.7	± 1.2	70.0	± 3.6	75.0	± 3.5	79.7	± 1.2	<b>87.7</b>	± 3.1
	100	85.7	± 1.5	83.7	± 3.8	<b>88.7</b>	± 1.5	<b>88.3</b>	± 1.5	<b>88.3</b>	± 2.9
OpenDoor	10	24.7	± 4.2	7.7	± 1.5	12.7	± 1.5	13.0	± 4.4	<b>42.0</b>	± 5.2
	50	44.0	± 1.7	11.3	± 3.5	22.3	± 9.3	38.0	± 5.0	<b>54.7</b>	± 5.0
	100	44.7	± 4.2	23.3	± 3.5	60.0	± 0.0	42.0	± 2.6	<b>65.0</b>	± 2.6
CloseDoor	10	2.0	± 1.0	3.0	± 1.0	0.0	± 0.0	4.3	± 0.6	<b>9.0</b>	± 1.7
	50	0.0	± 0.0	<b>7.3</b>	± 1.5	1.7	± 0.6	1.7	± 0.6	<b>6.7</b>	± 3.5
	100	2.7	± 1.5	<b>8.0</b>	± 0.0	5.0	± 0.0	5.0	± 2.6	<b>9.0</b>	± 1.0
Basketball in Hoop	10	1.7	± 1.5	2.7	± 1.5	1.7	± 1.2	2.0	± 1.0	<b>10.7</b>	± 2.5
	50	18.0	± 3.5	17.0	± 1.0	13.7	± 5.5	20.7	± 4.2	<b>47.3</b>	± 3.2
	100	38.7	± 4.2	34.3	± 4.2	46.0	± 4.0	42.7	± 4.0	<b>63.3</b>	± 3.8

### B. Learning Visuomotor Policies

The results for RL Bench are shown in Table I. CoDiG generally outperforms all baselines in the presented tasks across different number of demonstrations. CoDiG achieves 20% higher performance on average with 10 demonstrations, and 10% higher performance on average with 100 demonstrations. CoDiG is very sample efficient with low number of demonstrations as the model learns the motion of the task and then learns a residual on the visual observations. On the Adroit dataset, CoDiG does better by 3% on average, and the results can be found in Appendix D.

### C. Robustness to Visual Changes

We introduce modifications in tasks from RL Bench to test how learned models adapt to distribution shifts in the visual observations. We introduce two types of modifications- change in the color of the object being manipulated, and addition of visual distractors. Further, we also collect 5 additional demonstrations in the new modified environment to understand how models benefit from few-shot training on the out-of-distribution data. The results are shown in Table II. CoDiG is more robust to changes in the visual scene of the tasks, while the conventionally trained diffusion models fail miserably. Note that, CoDiG is not limited to the performance of the motion-model in zero-shot setting, which is 17%, 4% and 19% for *OpenBox*, *BasketballInHoop* and *OpenDoor* respectively, with 100 demonstrations. Further, CoDiG also responds better to additional demonstrations in the modified environment. This is expected, as we only update the guidance model with the additional demonstrations, in-effect, only updating the conditional distribution of the visual observation modality  $p(y^c|x, y^r)$ ,

rather than updating the conditional action distribution  $p(x|y^r, y^c)$ .

Task	Number of Demonstrations	DiT	DiT GMod
OpenBox	100 demos original	<b>98.3</b> ± 1.5	91.3 ± 4.5
	zero-shot blue box	43.3 ± 2.5	<b>46.7</b> ± 1.5
	5 demos blue box	34.7 ± 3.5	<b>76.7</b> ± 0.6
	zero-shot distractors	1.7 ± 2.1	<b>16.7</b> ± 2.1
	5 demos distractors	42.3 ± 4.0	<b>53.3</b> ± 2.3
Basketball in Hoop	100 demos original	38.7 ± 4.2	<b>63.3</b> ± 3.8
	zero-shot blue ball	29.0 ± 8.7	<b>63.0</b> ± 1.0
	5 demos blue ball	13.0 ± 0.0	<b>45.0</b> ± 2.6
	zero-shot distractors	0.7 ± 1.2	<b>56.3</b> ± 3.2
	5 demos distractors	2.7 ± 1.2	<b>39.7</b> ± 4.9
Open Door	100 demos original	44.7 ± 4.2	<b>65.0</b> ± 2.6
	zero-shot blue door	0.0 ± 0.0	<b>14.3</b> ± 3.1
	5 demos blue door	17.7 ± 3.1	<b>52.0</b> ± 7.2
	zero-shot red door	0.3 ± 0.6	<b>30.7</b> ± 2.5
	5 demos red door	20.0 ± 5.2	<b>53.7</b> ± 3.8

TABLE II  
PERFORMANCE COMPARISON ACROSS TASKS WITH VISUAL ADAPTATIONS FOR DiT AND CoDiG.

## IV. CONCLUSION

We present a novel compositional method to train diffusion models by decoupling different observational modalities such as proprioception, vision and tactile. Our factorized approach yields a loss function that is simple to use, while also presenting a more intuitive modeling paradigm. We also showcase CoDiG, an architectural implementation that makes the factorized training simpler and more efficient. We find that visuomotor CoDiG performs strongly over its baselines for tasks with a narrow distribution of robot motion. We also show that CoDiG is more robust to distribution shift in its visual observations than current diffusion policies. Several limitations exist, such as the choice to split the observation conditionals must be undertaken based on the task at hand. More experiments are needed to validate the efficacy, robustness and scaling properties of CoDiG.



# REFERENCES

- [1] Brian DO Anderson. Reverse-time diffusion equation models. Stochastic Processes and their Applications, 12(3):313–326, 1982.
- [2] Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I Jordan. An introduction to mcmc for machine learning. Machine learning, 50: 5–43, 2003.
- [3] Chen-Hao Chao, Wei-Fang Sun, Bo-Wun Cheng, Yi-Chen Lo, Chia-Che Chang, Yu-Lun Liu, Yu-Lin Chang, Chia-Ping Chen, and Chun-Yi Lee. Denoising likelihood score matching for conditional score-based data generation, 2022. URL <https://arxiv.org/abs/2203.14206>.
- [4] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. The International Journal of Robotics Research, page 02783649241273668, 2023.
- [5] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021. URL <https://arxiv.org/abs/2105.05233>.
- [6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. Advances in neural information processing systems, 34:8780–8794, 2021.
- [7] Carl Doersch. Tutorial on variational autoencoders. arXiv preprint arXiv:1606.05908, 2016.
- [8] Yilun Du, Conor Durkan, Robin Strudel, Joshua B. Tenenbaum, Sander Dieleman, Rob Fergus, Jascha Sohl-Dickstein, Arnaud Doucet, and Will Grathwohl. Reduce, reuse, recycle: Compositional generation with energy-based diffusion models and mcmc, 2023.
- [9] Bradley Efron. Tweedie’s formula and selection bias. Journal of the American Statistical Association, 106(496):1602–1614, 2011.
- [10] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning, 2020.
- [11] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598, 2022.
- [12] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022. URL <https://arxiv.org/abs/2207.12598>.
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. URL <https://arxiv.org/abs/2006.11239>.
- [14] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J. Davison. Rlbench: The robot learning benchmark & learning environment. CoRR, abs/1909.12271, 2019. URL <http://arxiv.org/abs/1909.12271>.
- [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.
- [16] Calvin Luo. Understanding diffusion models: A unified perspective, 2022. URL <https://arxiv.org/abs/2208.11970>.
- [17] William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. URL <https://arxiv.org/abs/2212.09748>.
- [18] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III 18, pages 234–241. Springer, 2015.
- [19] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics, 2015.
- [20] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022. URL <https://arxiv.org/abs/2010.02502>.
- [21] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution, 2020.
- [22] Yang Song and Diederik P. Kingma. How to train your energy-based models, 2021.
- [23] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2020.
- [24] Pascal Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661–1674, 2011.
- [25] Lirui Wang, Jialiang Zhao, Yilun Du, Edward H Adelson, and Russ Tedrake. Poco: Policy composition from and for heterogeneous robot learning. arXiv preprint arXiv:2402.02511, 2024.
- [26] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. URL <https://arxiv.org/abs/2302.05543>.

## APPENDIX A BACKGROUND

### A. Diffusion Models

Gaussian diffusion models [19] learn the reverse diffusion kernel  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  for a fixed forward kernel that adds Gaussian noise at each step  $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathcal{I})$ , such that  $q(\mathbf{x}_T) \approx \mathcal{N}(0, \mathcal{I})$ . Here,  $t \leq T$  represents the diffusion time-step and  $\alpha_t$  the noise schedule. To generate trajectories from the learned data distribution  $p_\theta(\mathbf{x}_0)$ , we sample at time step  $T$  from  $\mathcal{N}(0, \mathcal{I})$  and apply the reverse diffusion kernel  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  at each time step. For training the model, maximization of the evidence lower bound derived from the log-likelihood of the data distribution  $\log q(\mathbf{x}_0)$  yields the commonly used loss function in Equation 13 [13].

$$\mathcal{L}_t(\theta) = \mathbb{E}_{q(\mathbf{x}_0)\mathcal{N}(\epsilon_0; 0, \mathcal{I})} [\lambda_t \|\epsilon_0 - \hat{\epsilon}_\theta(\mathbf{x}_t, t)\|_2^2] \quad (13)$$

Here,  $\lambda_t$ , a function of  $\alpha_t$  is the weighting parameter for different time-steps, usually taken as 1. We train our model to predict the noise  $\epsilon_0$  added to the data sample  $\mathbf{x}_0$  to generate the noisy sample  $\mathbf{x}_t$  taken as input to the network. The Tweedie formula [9] can be used to show that  $\epsilon_0$ , and consequently  $\epsilon_\theta$  are proportional to the score of the diffused data distribution  $q(\mathbf{x}_t) = \int q(\mathbf{x}_t|\mathbf{x}_0)q(\mathbf{x}_0)d\mathbf{x}_0$  [16].

$$\frac{-1}{\sqrt{1-\alpha_t}}\hat{\epsilon}_\theta(\mathbf{x}_t, t) \approx \frac{-1}{\sqrt{1-\alpha_t}}\epsilon_0 = \nabla_{\mathbf{x}} \log q(\mathbf{x}_t) \quad (14)$$

### B. Energy Based Models (EBMs)

EBMs are a class of probabilistic models of the form  $p_\theta(\mathbf{x}) = \frac{e^{f_\theta(\mathbf{x})}}{Z}$  where  $Z(\theta) = \int e^{f_\theta(\mathbf{x})}d\mathbf{x}$  is the normalizing constant. Denoising score matching (DSM) [24] used to train EBMs [22] minimizes the Fisher divergence between the model  $p_\theta(\mathbf{x})$  and the Gaussian-smoothed data distribution  $q(\tilde{\mathbf{x}}) = \int q(\mathbf{x})\mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma_t^2 \mathcal{I})d\mathbf{x}$  at various noise scales  $\sigma_t$ .

$$\mathcal{J}_{\sigma_t}(\theta) = \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[ \frac{1}{2} \|\nabla_{\tilde{\mathbf{x}}} \log q(\tilde{\mathbf{x}}|\mathbf{x}) - \nabla_{\tilde{\mathbf{x}}} \log p_\theta(\tilde{\mathbf{x}})\|_2^2 \right] \quad (15)$$

This circumvents the normalizing constant by evaluating the gradient of the log-probability of the model  $\nabla_{\mathbf{x}} p_\theta(\tilde{\mathbf{x}}) = \nabla_{\mathbf{x}} f_\theta(\tilde{\mathbf{x}})$  at different noise scales. Equation 15 simplifies to the following when  $q(\mathbf{x}|\tilde{\mathbf{x}}) = \mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma_t^2 \mathcal{I})$  [24]:

$$\mathcal{J}_{\sigma_t}(\theta) = \mathbb{E}_{q(\mathbf{x})\mathcal{N}(\epsilon; 0, \mathcal{I})} \left[ \left\| \frac{\epsilon}{\sigma_t} + \nabla_{\mathbf{x}} f_\theta(\mathbf{x} + \sigma_t \epsilon) \right\|_2^2 \right] \quad (16)$$

Once trained, MCMC methods such as Langevin [2] can be used to sample from EBMs since they only depend on the score of the data distribution. This approach is also known in the literature as score-based modeling [21].

### C. Sampling from the Policy

Song et al. [23] show that score-based and denoising diffusion models can be considered as discretizations of a family of stochastic differential equations (SDE) that slowly add noise to the data distribution. For the generation process, the time-reversal of this SDE was given by Anderson [1]  $d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log q_t(\mathbf{x})] dt + g(t)d\bar{\mathbf{w}}$ , where  $\bar{\mathbf{w}}$  is a standard Wiener process for reverse time,  $g(t)$  is the scalar diffusion coefficient, and  $\mathbf{f}(\cdot, t)$  is the drift coefficient. Some manifestations of the reverse SDE equation are ancestral sampling [23] proposed by Ho et al. [13], shown in Equation 17, and Langevin dynamics [2].

$$\mathbf{x}_{t-1} \sim \mathcal{N} \left( \mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} [\mathbf{x}_t + (1 - \alpha_t) \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)], \sqrt{1 - \alpha_t} \mathcal{I} \right) \quad (17)$$

Equation 15 is then used to obtain an estimate of the score of the perturbed data distribution  $\nabla_{\mathbf{x}} \log q_t(\mathbf{x})$ , where the transition kernel  $q(\tilde{\mathbf{x}}|\mathbf{x})$  varies between approaches. Diffusion models use a forward transition kernel of  $\mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathcal{I})$ , yielding the same loss as Equation 13, while score-based model typically use  $\mathcal{N}(\mathbf{x}_t; \mathbf{x}_{t-1}, \sigma_t^2 \mathcal{I})$ , where  $\alpha_t$  and  $\sigma_t$  are respective noise scales.

### D. Classifier Guided Diffusion

Classifier guided diffusion [5] has received considerable attention due to its ability to reuse existing diffusion models to sample images conforming to specific classes. To sample from a class  $\mathbf{y}$ , Equation 18 as a result of Bayes' theorem allows us to decompose the conditional score at time-step  $t$  into the gradient of the classifier and the unconditional score.

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}; \theta, \phi) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t; \phi) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t; \theta) \quad (18)$$

Classifier guided-diffusion requires a classifier trained on noisy samples to get accurate estimate of the gradients  $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t; \phi)$ . Sampling from the class  $\mathbf{y}$  can then be achieved for diffusion models by substituting a re-weighted version of Equation 18 for  $\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$  in the ancestral sampling Equation 17.

## APPENDIX B METHODOLOGY

### A. Proof for Theorem II.1

**Theorem B.1.** *The diffusion loss function  $\mathcal{L}_t(\theta)$  as defined in Equation 1, in expectation over the time-steps  $1 \leq t \leq T$ , is equivalent to maximizing the ELBO of the log-likelihood of the conditional data distribution  $\log q(\mathbf{x}|\mathbf{y})$ , under a conditional Markovian noising process  $\hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1})$  and the reverse transition kernel as  $\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$ .*

Here, we derive the diffusion loss function for the conditional distribution  $p(\mathbf{x}|\mathbf{y})$  instead of only  $p(\mathbf{x})$ . A parallel derivation for conditional variational auto-encoders can be found in Doersch [7]. Following Dhariwal and Nichol [6], we start with a conditional Markovian noising forward process  $\hat{q}$  similar to  $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1 - \alpha_t)\mathcal{I})$ , and define the following:

$$\hat{q}(\mathbf{x}_0) := q(\mathbf{x}_0) \quad (19)$$

$$\hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{y}) := q(\mathbf{x}_{t+1}|\mathbf{x}_t) \quad (20)$$

$$\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y}) := \prod_{t=1}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y}) \quad (21)$$

We now reproduce some results that will be used later in the derivation of diffusion loss for conditional distributions. Dhariwal and Nichol [6] also show that

$$\hat{q}(\mathbf{y}|\mathbf{x}_t, \mathbf{x}_{t+1}) = \hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{y}) \frac{\hat{q}(\mathbf{y}|\mathbf{x}_t)}{\hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t)} \quad (22)$$

$$= \hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t) \frac{\hat{q}(\mathbf{y}|\mathbf{x}_t)}{\hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t)} \quad (23)$$

$$= \hat{q}(\mathbf{y}|\mathbf{x}_t) \quad (24)$$

Moreover, the unconditional reverse transition kernels can be shown to be equal using Bayes theorem, given Equations 19 and 20:  $\hat{q}(\mathbf{x}_t|\mathbf{x}_{t+1}) = q(\mathbf{x}_t|\mathbf{x}_{t+1})$ . Dhariwal and Nichol [6] use the result from Equation 24 to show the following for conditional reverse transition kernels.

$$\hat{q}(\mathbf{x}_t|\mathbf{x}_{t+1}, \mathbf{y}) = \frac{\hat{q}(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{y})}{\hat{q}(\mathbf{x}_{t+1}, \mathbf{y})} \quad (25)$$

$$= \frac{\hat{q}(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{y})}{\hat{q}(\mathbf{y}|\mathbf{x}_{t+1})\hat{q}(\mathbf{x}_{t+1})} \quad (26)$$

$$= \frac{\hat{q}(\mathbf{x}_t|\mathbf{x}_{t+1})\hat{q}(\mathbf{y}|\mathbf{x}_t, \mathbf{x}_{t+1})\hat{q}(\mathbf{x}_{t+1})}{\hat{q}(\mathbf{y}|\mathbf{x}_{t+1})\hat{q}(\mathbf{x}_{t+1})} \quad (27)$$

$$= \frac{\hat{q}(\mathbf{x}_t|\mathbf{x}_{t+1})\hat{q}(\mathbf{y}|\mathbf{x}_t, \mathbf{x}_{t+1})}{\hat{q}(\mathbf{y}|\mathbf{x}_{t+1})} \quad (28)$$

$$= \frac{q(\mathbf{x}_t|\mathbf{x}_{t+1})\hat{q}(\mathbf{y}|\mathbf{x}_t)}{\hat{q}(\mathbf{y}|\mathbf{x}_{t+1})} \quad (29)$$

Further, we can show the following using Equations 20 and 21 and the Markovian noising process. It states that the joint distribution of the noised samples conditioned on  $\mathbf{y}$  and  $\mathbf{x}_0$  are the same for both  $\hat{q}$  and  $q$ .

$$\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y}) = \prod_{t=1}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y}) \quad (30)$$

$$= \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (31)$$

$$= q(\mathbf{x}_{1:T}|\mathbf{x}_0) \quad (32)$$

We adapt the derivation of diffusion loss from Luo [16] to work with conditional distributions by maximizing the log-likelihood of the conditional data distribution  $\log p(\mathbf{x}|\mathbf{y})$  leading to evidence lower bound (ELBO).

$$\log p(\mathbf{x}|\mathbf{y}) = \log \int p(\mathbf{x}_{0:T}|\mathbf{y}) d\mathbf{x}_{1:T} \quad (33)$$

$$= \log \int \frac{p(\mathbf{x}_{0:T}|\mathbf{y})\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})}{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} d\mathbf{x}_{1:T} \quad (34)$$

$$= \log \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \frac{p(\mathbf{x}_{0:T}|\mathbf{y})}{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \right] \quad (35)$$

$$\geq \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_{0:T}|\mathbf{y})}{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \right] \quad (36)$$

The ELBO can be further simplified as follows

$$\log p(\mathbf{x}|\mathbf{y}) \geq \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_{0:T}|\mathbf{y})}{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \right] \quad (37)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\prod_{t=1}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y})} \right] \quad (38)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y}) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y}) \prod_{t=2}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y})} \right] \quad (39)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y}) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y}) \prod_{t=2}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0, \mathbf{y})} \right] \quad (40)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0, \mathbf{y})} \right] \quad (41)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\frac{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) \hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})}{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_0, \mathbf{y})}} \right] \quad (42)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\frac{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) \hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})}{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_0, \mathbf{y})}} \right] \quad (43)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} + \log \frac{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})}{\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y})} + \log \prod_{t=2}^T \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \right] \quad (44)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})}{\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y})} + \sum_{t=2}^T \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \right] \quad (45)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})] + \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y})}{\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y})} \right] + \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \right] \quad (46)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})] + \mathbb{E}_{\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p(\mathbf{x}_T|\mathbf{y})}{\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y})} \right] + \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t, \mathbf{x}_{t-1}|\mathbf{x}_0, \mathbf{y})} \left[ \log \frac{p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \right] \quad (47)$$

$$= \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y}) \parallel p(\mathbf{x}_T|\mathbf{y}))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}))]}_{\text{denoising matching term}} \quad (48)$$

The prior matching term does not contain any trainable parameters. We further simplify the denoising matching term using Equation 29 further conditioned on  $\mathbf{x}_0$ .

$$- \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}))] \quad (49)$$

$$= - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [\mathbb{E}_{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} [\log \hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) - \log p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})]] \quad (50)$$

$$= - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} \left[ \mathbb{E}_{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \left[ \log \hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) + \log \frac{\hat{q}(\mathbf{y}|\mathbf{x}_{t-1}, \mathbf{x}_0)}{\hat{q}(\mathbf{y}|\mathbf{x}_t, \mathbf{x}_0)} - \log p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}) \right] \right] \quad (51)$$

$$= - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}))] - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} \left[ \mathbb{E}_{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \left[ \log \frac{\hat{q}(\mathbf{y}|\mathbf{x}_{t-1}, \mathbf{x}_0)}{\hat{q}(\mathbf{y}|\mathbf{x}_t, \mathbf{x}_0)} \right] \right] \quad (52)$$

Using the result of Equation 52, Equation 48 can be rewritten as

$$\mathbb{E}_{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})] - D_{\text{KL}}(\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y}) \parallel p(\mathbf{x}_T|\mathbf{y})) - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}))] \quad (53)$$

$$\begin{aligned}
&= \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y}) \parallel p(\mathbf{x}_T|\mathbf{y}))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}))]}_{\text{denoising matching term}} \\
&\quad - \underbrace{\sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} \left[ \mathbb{E}_{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \left[ \log \frac{\hat{q}(\mathbf{y}|\mathbf{x}_{t-1}, \mathbf{x}_0)}{\hat{q}(\mathbf{y}|\mathbf{x}_t, \mathbf{x}_0)} \right] \right]}_{\text{label consistency term}}
\end{aligned} \tag{54}$$

The expression derived from the ELBO for conditional distribution introduces an additional term for label consistency. This minimizes the difference in the likelihood of the labels between consecutive denoising steps. However, since it does not have trainable parameters along with the prior matching term, we will ignore it. Moreover, it is easy to see from Luo [16] that the reconstruction term and the denoising matching term when developed further lead to the diffusion loss of Equation 1, with the additional parametrization of the model with  $\mathbf{y}$ . Note that the expectation is calculated over the same distributions, since  $\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y}) = q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ , as shown in Equation 32.

### B. Proof for theorem II.2

**Theorem B.2.** *Denoising score matching for the guidance model  $\pi_g: \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})$ , as expressed in Equation 5 is equal up to a constant to the following loss:*

$$L_{\text{DLSM}}(\phi) = \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) + \nabla_{\mathbf{x}_t} \log p_\theta(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_\alpha(\mathbf{x}_t|\mathbf{x})\|_2^2 \right] \tag{55}$$

For diffusion models, we take the forward transition kernel as  $p_\alpha(\mathbf{x}_t|\mathbf{x}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}, (1 - \alpha_t)\mathcal{I})$ . We noise the observations  $\mathbf{y}^{1:k+1}$  with a Gaussian kernel  $\mathcal{N}(\tilde{\mathbf{y}}; \mathbf{y}, \tau^2 I)$  of variance  $\tau^2$  that is small enough such that  $p_\tau(\tilde{\mathbf{y}}^i) \approx p(\mathbf{y}^i)$ .

Chao et al. [3] in their insightful work show that the following two losses differ only by a constant.-

$$D_F(p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t) \| p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t)) = \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t)\|_2^2 \right] \tag{56}$$

$$L_{\text{DLSM}}(\phi) = \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}, \mathbf{x}_t, \mathbf{y}^{k+1}, \tilde{\mathbf{y}}^{k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p_\theta(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p_\alpha(\mathbf{x}_t|\mathbf{x})\|_2^2 \right] \tag{57}$$

We extend their proof for multiple conditionals below. The Fisher divergence between the guidance model and the true score of the classifier (Equation 5) can be further expanded as:

$$D_F(p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) \| p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) \tag{58}$$

$$= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] \tag{59}$$

$$\begin{aligned}
&= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] + \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] \\
&\quad - \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) \rangle]
\end{aligned} \tag{60}$$

$$\begin{aligned}
&= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] + \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] \\
&\quad - \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k}, \tilde{\mathbf{y}}^{k+1}) - \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k}) \rangle]
\end{aligned} \tag{61}$$

$$\begin{aligned}
&= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] + \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] \\
&\quad + \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k}) \rangle] \\
&\quad - \underbrace{\mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k}, \tilde{\mathbf{y}}^{k+1}) \rangle]}_{\text{Term 1}}
\end{aligned} \tag{62}$$

Simplifying the Term 1 further:

$$\begin{aligned}
&= -\mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k+1}) \rangle] \\
&= -\int_{\mathbf{x}_t} \int_{\tilde{\mathbf{y}}^{1:k+1}} p_\tau(\tilde{\mathbf{y}}^{1:k+1}) p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k+1}) \langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \frac{\nabla_{\mathbf{x}_t} p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k+1})}{p_{\alpha, \tau}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k+1})} \rangle d\tilde{\mathbf{y}}^{1:k+1} d\mathbf{x}_t \\
&= -\int_{\mathbf{x}_t} \int_{\tilde{\mathbf{y}}^{1:k+1}} p_\tau(\tilde{\mathbf{y}}^{1:k+1}) \langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \int_{\mathbf{x}_0} p_{0, \tau}(\mathbf{x}_0|\tilde{\mathbf{y}}^{1:k+1}) p_{\alpha, \tau}(\mathbf{x}_t|\mathbf{x}_0, \tilde{\mathbf{y}}^{1:k+1}) d\mathbf{x}_0 \rangle d\tilde{\mathbf{y}}^{1:k+1} d\mathbf{x}_t \\
&= -\int_{\mathbf{x}_t} \int_{\tilde{\mathbf{y}}^{1:k+1}} p_\tau(\tilde{\mathbf{y}}^{1:k+1}) \langle \nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \int_{\mathbf{x}_0} \int_{\mathbf{y}^{1:k+1}} p_{0, \tau}(\mathbf{x}_0|\tilde{\mathbf{y}}^{1:k+1}) p_{\alpha, \tau}(\mathbf{x}_t|\mathbf{x}_0, \tilde{\mathbf{y}}^{1:k+1}, \mathbf{y}^{1:k+1}) p(\mathbf{y}^{1:k+1}|\mathbf{x}_0, \tilde{\mathbf{y}}^{1:k+1}) d\mathbf{y}^{1:k+1} d\mathbf{x}_0 \rangle d\tilde{\mathbf{y}}^{1:k+1} d\mathbf{x}_t
\end{aligned}$$

$$\begin{aligned}
&= - \int_{\mathbf{x}_t} \int_{\tilde{\mathbf{y}}^{1:k+1}} \int_{\mathbf{x}_0} \int_{\mathbf{y}^{1:k+1}} p_{\tau}(\mathbf{x}_0, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1}) \langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t | \mathbf{x}_0, \tilde{\mathbf{y}}^{1:k+1}, \mathbf{y}^{1:k+1}) \rangle d\mathbf{y}^{1:k+1} d\mathbf{x}_0 d\tilde{\mathbf{y}}^{1:k+1} d\mathbf{x}_t \\
&= - \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_0, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0) \rangle]
\end{aligned}$$

Plugging this back into Equation 62, we get-

$$\begin{aligned}
&D_F(p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) || p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) \\
&= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} ||\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})||_2^2 \right] + \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} ||\nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})||_2^2 \right] \\
&+ \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) \rangle] \\
&- \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_0, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0) \rangle]
\end{aligned} \tag{63}$$

Here,  $\mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} ||\nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})||_2^2 \right]$  is a constant. Further, adding the constant  $\mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})} \left[ \frac{1}{2} ||\nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0)||_2^2 \right]$  to Equation 63, we get:

$$\begin{aligned}
&D_F(p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) || p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) \\
&= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} ||\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})||_2^2 \right] + \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) \rangle] + C \\
&- \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_0, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0) \rangle] + \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})} \left[ \frac{1}{2} ||\nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0)||_2^2 \right]
\end{aligned} \tag{64}$$

$$= \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} \left[ \frac{1}{2} ||\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) + \nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x})||_2^2 \right] + C \tag{65}$$

## APPENDIX C EXPERIMENTAL SETUP

### A. Environments

We use 5 tasks from RL Bench [14] and all the tasks from the dexterous hand-manipulation environment Adroit [10] to evaluate our performance. RL Bench has a wide variety of tasks and an inbuilt planner that enables collection of demonstrations. We train visuomotor policies on RL Bench with a 5 camera setup recording 96x96 RGB images that use joint position as the action modality. The tasks considered in this paper are shown in Figure 3. Further we also modify the visual scene of these tasks to evaluate the performance of models under visual distribution shift, as shown in Figure 4.

The Adroit dataset focuses on high-dimensional dexterous manipulation tasks using a 24-degree-of-freedom anthropomorphic hand. The dataset encompasses a suite of challenging manipulation tasks designed to test fine motor control and complex hand-object interactions. The environment state is specified as a low-dimensional vector tailored to each task.



Fig. 3. RL Bench Tasks: *OpenBox*; *CloseBox*; *BasketballInHoop*; *OpenDoor*; *CloseDoor*

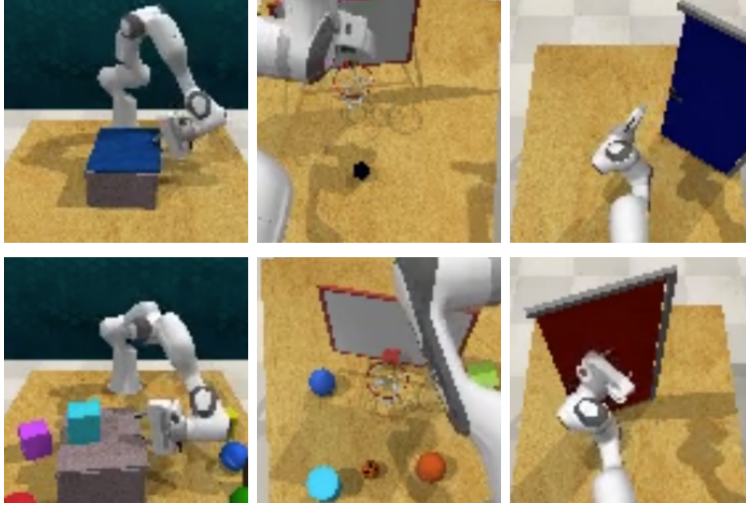


Fig. 4. **Top row from left to right:** *OpenBox* task with a blue lid; *BasketballInHoop* with a blue ball; *OpenDoor* with a blue door. **Bottom row from left to right:** *OpenBox* task with distractors; *BasketballInHoop* with distractors; *OpenDoor* with a red door.

### B. Baselines

POCO [25] can be formulated for composing a motion and a vision guidance model as follows. We train the motion and the vision guidance models separately prior to the compositional sampling proposed by Du et al. [8]. The motion model is only trained on proprioceptive observations, while the vision guidance model also uses the vision observations. We set  $\lambda = 0.2$  in Equation 66.

$$\begin{aligned} \mathbf{x}_{t-1} &\sim \mathcal{N}\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \boldsymbol{\epsilon}(\mathbf{x}_t, \mathbf{y}^r, \mathbf{y}^c, t) \right), \sqrt{1-\alpha_t} \mathcal{I} \right) \\ \boldsymbol{\epsilon}(\mathbf{x}_t, \mathbf{y}^r, \mathbf{y}^c, t) &= \hat{\boldsymbol{\epsilon}}_\phi(\mathbf{x}_t, \mathbf{y}^r, \mathbf{y}^c, t) + \lambda * \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, \mathbf{y}^r, t) \end{aligned} \quad (66)$$

For Classifier-free guidance [12], we train a single model and switch out the vision modality with a probability of 0.1. We then sample according to the Equation 67, where we set  $\lambda_1 = 1.1$  and  $\lambda_2 = 0.1$ .

$$\begin{aligned} \mathbf{x}_{t-1} &\sim \mathcal{N}\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \boldsymbol{\epsilon}(\mathbf{x}_t, \mathbf{y}^r, \mathbf{y}^c, t) \right), \sqrt{1-\alpha_t} \mathcal{I} \right) \\ \boldsymbol{\epsilon}(\mathbf{x}_t, \mathbf{y}^r, \mathbf{y}^c, t) &= \lambda_1 * \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, \mathbf{y}^r, \mathbf{y}^c, t) + \lambda_2 * \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, \mathbf{y}^r, \phi, t) \end{aligned} \quad (67)$$

### C. Training and Rollout

All models are trained for 2000 epochs for RLBench tasks and 3000 epochs for Adroit. We use Adam optimizer [15] and a learning rate of  $1e-4$ , in accordance to Peebles and Xie [17]. All our results are reported as the mean and standard deviation of the success rates over 100 rollouts for 3 different seeds (total of 300 rollouts) for 10, 50 and 100 demonstrations for each task. We rollout models for 300 and 475 timesteps for RLBench and Adroit tasks, respectively. We report results for the number of rollouts with cumulative success signals greater than 1 and 25 for all Adroit tasks.

## APPENDIX D RESULTS: ADROIT

We show the results in the Adroit environments in Table III. Note that we modify the *Hammer* environment to consider the distance between the nail and the board as greater than 0.2 to be successful, instead of 0.1. Here **DiT CMod** refers to an architectural ablation, which resembles ControlNet closely. Here, the noisy actions are passed into the transformer and then conditioned upon the environment state using adaLn.

Task	Criteria	DiT	UNet	DiT-cfg	DiT CMod	CoDiG
Door Human	Success achieved > 1	58.3 $\pm$ 6.7	67.3 $\pm$ 4.7	42.7 $\pm$ 5.5	60.3 $\pm$ 6.7	<b>74.3</b> $\pm$ 3.1
	Success achieved > 25	34.3 $\pm$ 2.5	30.7 $\pm$ 0.6	22.3 $\pm$ 7.5	17.7 $\pm$ 3.1	<b>45.7</b> $\pm$ 5.8
Pen Human	Success achieved > 1	58.3 $\pm$ 1.5	63.3 $\pm$ 3.5	<b>67.7</b> $\pm$ 3.5	62.3 $\pm$ 5.9	62.3 $\pm$ 4.7
	Success achieved > 25	53.3 $\pm$ 1.2	57.3 $\pm$ 2.5	<b>62.3</b> $\pm$ 3.2	56.0 $\pm$ 4.4	56.0 $\pm$ 3.5
Hammer Human	Success achieved > 1	<b>54.0</b> $\pm$ 1.7	47.3 $\pm$ 6.0	40.7 $\pm$ 2.3	<b>56.7</b> $\pm$ 4.6	51.7 $\pm$ 6.4
	Success achieved > 25	<b>52.7</b> $\pm$ 2.3	47.0 $\pm$ 6.0	39.0 $\pm$ 2.0	<b>56.0</b> $\pm$ 5.2	<b>50.7</b> $\pm$ 6.4
Relocate Human	Success achieved > 1	<b>63.0</b> $\pm$ 4.6	56.0 $\pm$ 7.2	1.3 $\pm$ 0.6	43.3 $\pm$ 4.0	63.0 $\pm$ 6.9
	Success achieved > 25	<b>60.0</b> $\pm$ 3.6	52.0 $\pm$ 7.0	0.3 $\pm$ 0.6	38.7 $\pm$ 4.6	<b>58.7</b> $\pm$ 4.0

TABLE III  
PERFORMANCE COMPARISON ACROSS ADROIT TASKS FOR DIFFERENT MODELS.