

LEARNING COMPOSABLE DIFFUSION GUIDANCE FOR MOTION PRIORS

Omkar Patil

Arizona State University

opatil3@asu.edu

Eric Rosen

eric.andrew.rosen@gmail.com

Nakul Gopalan

Arizona State University

ng@asu.edu

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 method to train diffusion models on skill demonstrations borrow from the literature on image generation, and results in a conditional distribution of robot actions given the visual, proprioceptive and other observations. However, they have 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. CoDiG decouples the observation modalities allowing the residual learning of one modality with respect to the others. While presenting a more intuitive modeling paradigm, CoDiG also enables the scaling of modalities such as robot motions independently. Our preliminary results show that visual CoDiG with motion-priors outperforms the conventional way of learning visuomotor policies using diffusion models on skills with relatively low-diversity of robot motion. Further experimentation is needed to evaluate the performance and robustness of CoDiG for different observation modalities, and on different classes of skills, such as long-horizon and precise manipulation.

1 INTRODUCTION

Diffusion models have emerged as a promising choice for learning robot skills from demonstrations Chi et al. (2023). 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. Prior work Chi et al. (2023) 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 visuo-motor 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.

Recent efforts to exploit compositionality through diffusion models have yielded promising results Du et al. (2023); Wang et al. (2024). For instance, POCO Wang et al. (2024) composes multi-task conditional and unconditional policies to show better generalization over using a single multi-task policy for tool use. Wang et al. (2024) emphasize the heterogeneous nature of data collection in robotics and showcase composition of policies trained on different data modalities such as vision-action and tactile-action. However, the candidate policies are learned independently and require manual tuning of their compositional weights. Moreover, these methods are limited to the combination of pre-learned policies, requiring a careful consideration of the candidate tasks and data collection.

Critically, CoDiG enables the training of models conditioned on one observational modality with respect to the rest so that data collection efforts can be concentrated on the cheaper modality and consequent models learn on the residual rather than in an isolated manner. Consequently, CoDiG does not require manual tuning of the compositional weights. Specifically, we showcase this by training a ‘motion-model’ on the action distribution of a task, and learning a vision ‘guidance-model’ relative to it on vision-action coupled demonstration data. Through our experiments, we show that not all tasks are created equal, and policy learning using our method can benefit from further collection of only the robot motion data corresponding to the task, while reducing the requirement of collecting coupled vision-action data. Ultimately, CoDiG is a general policy learning method that can be flexibly used with different modalities of data

and shows better scaling properties for decoupled task-specific data collection. Further experiments are needed to validate the efficacy of CoDiG, especially for long-horizon tasks and for its scaling properties. Our contributions are as follows-

- We present a novel method to train diffusion models on robot demonstration data called CoDiG, abbreviated for ‘Composable Diffusion Guidance’. Conventional diffusion policies Chi et al. (2023) learn a conditional distribution of actions with respect to the all observational modalities. In contrast, CoDiG decouples the observational modalities allowing the residual learning of one modality with respect to the others.
- CoDiG allows data collection to focus on the cheapest modality, removing the necessity of collecting paired modality data together. Our experiments show that not all skills are created equal, and some skills are less reliant on visual observations for their motions. We use this insight to show that CoDiG benefits from motion-only data collection for such skills, while limiting the collection of paired visual-action data.
- Our results show that visual CoDiG with a motion-model performs favorably with respect to the conventional diffusion models on several RL Bench tasks. Further experimentation is required to validate the efficacy and robustness of CoDiG for long horizon tasks and precise manipulation. Our initial results also show that CoDiG scales better with additional motion data.

2 BACKGROUND

2.1 DIFFUSION MODELS

Gaussian diffusion models Sohl-Dickstein et al. (2015) 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 α_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 1 Ho et al. (2020).

$$\mathcal{L}_t(\theta) = \mathbb{E}_{q(\mathbf{x}_0)\mathcal{N}(\epsilon_0; 0, \mathcal{I})} [\lambda_t \|\epsilon_0 - \hat{\epsilon}_{\theta}(\mathbf{x}_t, t)\|_2^2] \quad (1)$$

Here, λ_t , a function of α_t is the weighting parameter for different time-steps, usually taken as 1. We train our model to predict the noise ϵ_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 Efron (2011) can be used to show that ϵ_0 , and consequently ϵ_{θ} are proportional to the score of the diffused data distribution $q(\mathbf{x}_t) = \int q(\mathbf{x}_t|\mathbf{x}_0)q(\mathbf{x}_0)d\mathbf{x}_0$ Luo (2022).

$$\frac{-1}{\sqrt{1 - \alpha_t}} \hat{\epsilon}_{\theta}(\mathbf{x}_t, t) \approx \frac{-1}{\sqrt{1 - \alpha_t}} \epsilon_0 = \nabla_{\mathbf{x}} \log q(\mathbf{x}_t) \quad (2)$$

2.2 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) Vincent (2011) used to train EBMs Song and Kingma (2021) 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 σ_t .

$$\mathcal{J}_{\sigma_t}(\theta) = \mathbb{E}_{q(\mathbf{x}, \tilde{\mathbf{x}})} \left[\frac{1}{2} \|\nabla_{\tilde{\mathbf{x}}} \log q(\tilde{\mathbf{x}}|\mathbf{x}) - \nabla_{\tilde{\mathbf{x}}} \log p_{\theta}(\tilde{\mathbf{x}})\|_2^2 \right] \quad (3)$$

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 3 simplifies to the following when $q(\mathbf{x}|\tilde{\mathbf{x}}) = \mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma_t^2\mathcal{I})$ Vincent (2011):

$$\mathcal{J}_{\sigma_t}(\theta) = \mathbb{E}_{q(\mathbf{x})\mathcal{N}(\epsilon; 0, \mathcal{I})} \left[\left\| \frac{\epsilon}{\sigma_t} + \nabla_{\mathbf{x}} f_{\theta}(\mathbf{x} + \sigma_t\epsilon) \right\|_2^2 \right] \quad (4)$$

Once trained, MCMC methods such as Langevin Andrieu et al. (2003) 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 Song and Ermon (2020).

2.3 SAMPLING FROM THE POLICY

Song et al. (2020) 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 (1982) $d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - g(t)^2 \nabla_{\mathbf{x}} \log q_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}}$, where $\bar{\mathbf{w}}$ is a standard Wiener process for reverse time, $g(t)$ is the scalar diffusion coefficient, and $f(\cdot, t)$ is the drift coefficient. Some manifestations of the reverse SDE equation are ancestral sampling Song et al. (2020) proposed by Ho et al. (2020), shown in Equation 5, and Langevin dynamics Andrieu et al. (2003).

$$\mathbf{x}_{t-1} \sim \mathcal{N}\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} [\mathbf{x}_t + (1 - \alpha_t) \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)], \sqrt{1 - \alpha_t} \mathcal{I}\right) \quad (5)$$

Equation 3 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 $\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_{t-1}, (1 - \bar{\alpha}_t) \mathcal{I})$, yielding the same loss as Equation 1, while score-based model typically use $\mathcal{N}(\mathbf{x}_t; \mathbf{x}_{t-1}, \sigma_t^2 \mathcal{I})$, where α_t and σ_t are respective noise scales.

2.4 CLASSIFIER GUIDED DIFFUSION

Classifier guided diffusion Dhariwal and Nichol (2021a) 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 6 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.

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}; \boldsymbol{\theta}, \phi) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t; \phi) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t; \boldsymbol{\theta}) \quad (6)$$

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 6 for $\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$ in the ancestral sampling Equation 5.

3 METHODOLOGY

3.1 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 Ho et al. (2020); Luo (2022); Song et al. (2022). This formulation has been directly adopted by the robotics community Chi et al. (2023) leading to the popular optimization objective shown in Equation 7.

$$\mathcal{L}_t(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{y}) \mathcal{N}(\epsilon_0; 0, \mathcal{I})} [\|\epsilon_0 - \hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, \mathbf{y}, t)\|_2^2] \quad (7)$$

$$\mathbf{x}_{t-1} \sim \mathcal{N}\left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, \mathbf{y}, t)\right), \sqrt{1 - \alpha_t} \mathcal{I}\right) \quad (8)$$

Here, the network $\epsilon_{\boldsymbol{\theta}}$ is also parametrized with \mathbf{y} for the conditional prediction of the noise added to the action \mathbf{x} . Once the network $\hat{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, \mathbf{y}, t)$ is trained, ancestral sampling shown in Equation 8 can be used to sample actions given the observations \mathbf{y} . However, while prior work assumes that $\epsilon_{\boldsymbol{\theta}}$ 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 Dhariwal and Nichol (2021a) does indeed result in the loss of Equation 7 Chi et al. (2023) 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 3.1. *The diffusion loss function $\mathcal{L}_t(\boldsymbol{\theta})$ as defined in Equation 7, 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 3.1 can be found in the Appendix A.1. 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 6 at diffusion time-step t as follows:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}^{1:k}, \mathbf{y}^{k+1}; \boldsymbol{\theta}, \phi) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y}^{k+1}|\mathbf{x}_t, \mathbf{y}^{1:k}; \phi) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}^{1:k}; \boldsymbol{\theta}) \quad (9)$$

Here $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}^{1:k}; \boldsymbol{\theta})$ is the score of the model trained on the original k observational modalities, while $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}^{k+1} | \mathbf{x}_t, \mathbf{y}^{1:k}; \boldsymbol{\phi})$ corresponds to the score of the classifier for the modality \mathbf{y}^{k+1} . The method proposed by classifier guided-diffusion Dhariwal and Nichol (2021a) 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 $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}^{k+1} | \mathbf{x}_t, \mathbf{y}^{1:k}; \boldsymbol{\phi})$ using a neural network, which we refer to as the guidance model π_g , rather than to learn a classifier and then obtain its gradients. To learn the guidance model, following the suit of Equation 3, we minimize the Fisher divergence between the guidance model and the true score of the classifier.

$$D_F(p_\phi(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) || p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) = \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha, \tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] \quad (10)$$

Here, observation modalities $\mathbf{y}^{1:k+1}$ can be noised with a Gaussian kernel $\mathcal{N}(\tilde{\mathbf{y}}; \mathbf{y}, \tau^2 I)$ of variance τ^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 10 is difficult to use in its current form as estimation of the true score is difficult for large datasets. Chao et al. (2022) derive the denoising likelihood score matching (DLSM) objective from the denoising score matching of Equation 3 for conditional distributions, which forms the basis of our next result.

Theorem 3.2. *Denoising score matching for the guidance model π_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 10 is equal up to a constant to the following loss:*

$$L_{DLSM}^t(\phi) = \mathbb{E}_{p_{\alpha, \tau}(\mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_\phi(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) + \nabla_{\mathbf{x}_t} \log p_\theta(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_\alpha(\mathbf{x}_t | \mathbf{x})\|_2^2 \right] \quad (11)$$

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 τ^2 that is small enough such that $p_\tau(\tilde{\mathbf{y}}^i) \approx p(\mathbf{y}^i)$.

The proof for Theorem 3.2 is a conditional variant of the one derived by Chao et al. (2022), and is presented in Appendix A.2. Equation 53 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 $\epsilon_0 \sim \mathcal{N}(0, \mathcal{I})$, and using Equation 2, we obtain:

$$L_{DLSM}(\phi) = \mathbb{E}_{p_\tau(\mathbf{x}, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} \mathbb{E}_{\epsilon_0 \sim \mathcal{N}(0, \mathcal{I})} \left[\frac{1}{2} \|\hat{\epsilon}_\phi(\tilde{\mathbf{y}}^{k+1}, \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) + \epsilon_\theta(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) - \epsilon_0\|_2^2 \right] \quad (12)$$

Here, ϵ_θ is frozen and known from prior training with observational modalities $\mathbf{y}^{1:k}$, while ϵ_ϕ is learned in expectation over the new data collected in conjugation with modality \mathbf{y}^{k+1} . Equation 12 is very similar to the diffusion loss Equation 7, 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, ϵ_ϕ is learned, actions can be sampled from the conditional distribution $p(\mathbf{x} | \mathbf{y}^{1:k+1})$ using ancestral sampling and Equation 9:

$$\mathbf{x}_{t-1} \sim \mathcal{N} \left(\mathbf{x}_t; \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon(\mathbf{x}_t, \mathbf{y}^{1:k+1}, t) \right), \sqrt{1 - \alpha_t} \mathcal{I} \right) \quad (13)$$

$$\epsilon(\mathbf{x}_t, \mathbf{y}^{1:k+1}, t) = \hat{\epsilon}_\phi(\mathbf{x}_t, \mathbf{y}^{1:k+1}, t) + \hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}^{1:k}, t) \quad (14)$$

Note that the score of the prior distribution $\hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}^{1:k}, t)$ or $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}^{1:k}; \boldsymbol{\theta})$ 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 $\mathbf{y}^{1:k}$.

3.2 VISUAL CODIG WITH MOTION-MODEL

We now provide a concrete implementation of the method proposed in Section 3.1. 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’ $p(\mathbf{y}^c | \mathbf{x}, \mathbf{y}^r)$. This implies that for certain tasks, we can collect a larger set of pure motions of the robot performing the specific skill, and fewer visual demonstrations of the robot actually performing the task. We learn the motion-model using diffusion loss of Equation 7 and the guidance model using DLSM loss of Equation 12.

$$L_{MM}(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{x}_0, \mathbf{y}^r) \mathcal{N}(\epsilon_0; 0, \mathcal{I})} \left[\|\epsilon_0 - \hat{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}^r, t)\|_2^2 \right] \quad (15)$$

$$L_G(\phi) = \mathbb{E}_{p_\tau(\mathbf{x}, \mathbf{y}^r, \mathbf{y}^c, \tilde{\mathbf{y}}^r, \tilde{\mathbf{y}}^c)} \mathbb{E}_{\epsilon_0 \sim \mathcal{N}(0, \mathcal{I})} \left[\frac{1}{2} \|\epsilon_0 - \hat{\epsilon}_\phi(\tilde{\mathbf{y}}^c, \mathbf{x}_t, \tilde{\mathbf{y}}^r, t) - \epsilon_\theta(\mathbf{x}_t, \tilde{\mathbf{y}}^r, t)\|_2^2 \right] \quad (16)$$

In Section 4, we showcase that this approach has a distinct advantage for skills that are heavily dependent on proprioception, while performing comparably on others. 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. Moreover, CoDiG allows us to tailor the collection of data to suit the task and costs.

4 RESULTS

4.1 ENVIRONMENTS AND BASELINES

We use 9 randomly-picked tasks from RL Bench James et al. (2019) and the low-dimensional Franka-Kitchen environment Gupta et al. (2019) 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. Franka-Kitchen uses joint velocity as the action modality with a 30-dimensional observational state that includes the position of the robot joints and the objects in the kitchen, along with the goal states encoding the target objects to manipulate. There are three versions of the dataset-complete, partial and mixed. The complete version is a smaller part of the dataset where all the demonstrations correspond to a fixed set of objects (microwave, kettle, light and slider) that are manipulated in the same order.

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 ($\sim 90\text{M}$) Peebles and Xie (2023) as our model architecture, which is kept the same for both the baseline and CoDiG. Since CoDiG uses two separate models for motion and guidance, we also include DiT-base ($\sim 190\text{M}$) in the baseline. All models are trained for 2000 epochs for RL Bench tasks and 3000 epochs for Franka-Kitchen. We use Adam optimizer Kingma and Ba (2017) and a learning rate of $1e - 4$, in accordance to Peebles and Xie (2023).

In this work, we present two types of results. In the first, we train both the baseline and our method using only 10 visual demonstrations from RL Bench to compare how CoDiG fares relative to the existing way of learning visuomotor policies. For Franka-Kitchen we train all models, including the motion-model only on the ‘complete’ version of the dataset. In the second, we collect an additional ‘motions’ of the robot performing the skill and record the evolution of the action and the proprioceptive states. In simulation, the motions are collected using the same planner used to collect visual demonstration in the environment. Our objective for the second set of results is to understand how collection of additional motion data can benefit CoDiG in comparison to directly learning the conditional distribution. In addition to these, we also show results for the motion-model that does not use any visual information. All our results for RL Bench are reported as the mean of the success rates over 50 rollouts for 3 different seeds (total of 150 rollouts). For Franka-Kitchen we report the mean task-completion rate out of 4 tasks. We leave real-world experiments and evaluation on play-motion data for future work.

4.2 LEARNING VISUOMOTOR POLICIES

For RL Bench, we train the DiT-small and DiT-base baselines on 10 visual demonstrations using the loss defined in Equation 7. For CoDiG, we train the motion model π_{mm} on only the proprioception-action data using the loss defined in Equation 15 and the vision guidance-model relative to the motion model using Equation 16. We also report the mean success rate for the motion-model π_{mm} . The results for RL Bench and Franka-Kitchen (complete) are shown in Table 1.

CoDiG performs comparably with the baselines, and shows strong improvements when the motion of the robot is heavily dependent on the robot proprioception. An indicator of this is the performance of the motion model (MM), since it is only trained on the proprioception of the robot and does not use any visual information. CoDiG outperforms the baselines on the tasks of *OpenBox* and *CloseBox* by large margins, which can be corroborated by the better performance of the motion-model on them, as compared to the other tasks.

4.3 PERFORMANCE WITH ADDITIONAL MOTION DEMONSTRATIONS

To understand how the performance of the model scales with independent modalities, we collect 40 and 90 additional motions of the robot performing the task in simulation. The implication on training CoDiG is that the guidance-models learn (with only 10 visual demonstrations) relative to the motion-models trained on different number of motion demonstrations (40 and 90). The results are shown in Table 2.

Although pre-training DiT with additional motion data does help improve the performance for tasks such as *OpenBox* and *CloseBox*, it still falls short of CoDiG even without any additional motion demonstrations. We see performance gains for CoDiG with increasing the available task-specific robot motions, for skills that are not heavily dependent on

Tasks	DiT	DiT x2	CoDiG	MM
Demonstrations	10	10	10	(50)
Open Box	18	20	39.33	17.33
Open Door	24	27.33	20.67	7.33
Open Microwave	2	1.33	1.33	0
Open Fridge	0.66	1.33	2	0
Close Box	29.33	18	52	33.33
Close Door	2	1.33	0.67	0.67
Close Microwave	49.33	49.33	36	24
Close Fridge	53.33	52	54	46
Basketball in Hoop	0.66	2	10	2
Kitchen (complete)	1.2	1.12	1.22	1.06

Table 1: Performance comparison of different models across various tasks on RLBench and Franka-Kitchen (complete). All the models on RLBench are trained using 10 visual demonstrations, except the motion-model which is trained on 50 motions of the robot performing the task. DiT x2 denotes a DiT-base model that is twice the size of the smaller variant, and comparable to CoDiG.

Task	CoDiG	CoDiG	CoDiG	MM	DiT x2 Pt.
Visual(+Motion) Demos	10(10)	10(50)	10(100)	0(50)	10(50)
Open Box	39.33	31.33	44.6	17.33	22
Open Door	20.67	18	18.67	7.33	3.33
Open Microwave	1.33	0.67	2	0	6
Open Fridge	2	0	0	0	2
Close Box	52	50.67	56.67	33.33	38.67
Close Door	0.67	2	2.67	0.67	2.67
Close Microwave	36	36.67	33.34	24	38.67
Close Fridge	54	58.67	44	46	54
Basketball in Hoop	10	9.33	9.33	2	3.33

Table 2: Performance comparison on scaling the collection of robot motions specific to the task. All the models except MM use only 10 visual demonstrations, and further robot motions as indicated in brackets. DiT x2 Pt. refers to DiT-base that has been pre-trained using the motion data and then fine-tuned on the visual demonstrations.

visual observations. However, the magnitude of improvement is small, indicating that there is scope to improve the interaction between the motion and guidance models.

5 CONCLUSION

We present CoDiG, a novel compositional method to train diffusion models by decoupling different observational modalities such as proprioception, vision and tactile. CoDiG yields a loss function that is simple to use, while also presenting a more intuitive modeling paradigm. Our experiments show that not all tasks are equal, and some tasks depend heavily on certain observational modalities than others. CoDiG allows us to tailor the data collection for tasks based on their dependence on observational modalities and costs. We find that visual CoDiG with a motion-model performs strongly over its baselines for visuomotor tasks that rely more on the proprioception. More experiments are needed to validate the efficacy, robustness and scaling properties of CoDiG.

REFERENCES

- Brian DO Anderson. Reverse-time diffusion equation models. Stochastic Processes and their Applications, 12(3): 313–326, 1982.
- Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I Jordan. An introduction to mcmc for machine learning. Machine learning, 50:5–43, 2003.
- 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>.
- 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.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021a. URL <https://arxiv.org/abs/2105.05233>.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. Advances in neural information processing systems, 34:8780–8794, 2021b.
- Carl Doersch. Tutorial on variational autoencoders. arXiv preprint arXiv:1606.05908, 2016.
- 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.
- Bradley Efron. Tweedie’s formula and selection bias. Journal of the American Statistical Association, 106(496): 1602–1614, 2011.
- Abhishek Gupta, Vikash Kumar, Corey Lynch, Sergey Levine, and Karol Hausman. Relay policy learning: Solving long-horizon tasks via imitation and reinforcement learning, 2019. URL <https://arxiv.org/abs/1910.11956>.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. URL <https://arxiv.org/abs/2006.11239>.
- 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>.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL <https://arxiv.org/abs/1412.6980>.
- Calvin Luo. Understanding diffusion models: A unified perspective, 2022. URL <https://arxiv.org/abs/2208.11970>.
- William Peebles and Saining Xie. Scalable diffusion models with transformers, 2023. URL <https://arxiv.org/abs/2212.09748>.
- Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics, 2015.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022. URL <https://arxiv.org/abs/2010.02502>.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution, 2020.
- Yang Song and Diederik P. Kingma. How to train your energy-based models, 2021.
- 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.
- Pascal Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661–1674, 2011.
- 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.

A METHODOLOGY

A.1 PROOF FOR THEOREM 3.1

Theorem A.1. *The diffusion loss function $\mathcal{L}_t(\boldsymbol{\theta})$ as defined in Equation 7, 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 (2016). Following Dhariwal and Nichol (2021b), 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:

$$\hat{q}(\mathbf{x}_0) := q(\mathbf{x}_0) \quad (17)$$

$$\hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t, \mathbf{y}) := q(\mathbf{x}_{t+1}|\mathbf{x}_t) \quad (18)$$

$$\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y}) := \prod_{t=1}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y}) \quad (19)$$

We now reproduce some results that will be used later in the derivation of diffusion loss for conditional distributions. Dhariwal and Nichol (2021b) also show that

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

$$= \hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t) \frac{\hat{q}(\mathbf{y}|\mathbf{x}_t)}{\hat{q}(\mathbf{x}_{t+1}|\mathbf{x}_t)} \quad (21)$$

$$= \hat{q}(\mathbf{y}|\mathbf{x}_t) \quad (22)$$

Moreover, the unconditional reverse transition kernels can be shown to be equal using Bayes theorem, given Equations 17 and 18: $\hat{q}(\mathbf{x}_t|\mathbf{x}_{t+1}) = q(\mathbf{x}_t|\mathbf{x}_{t+1})$. Dhariwal and Nichol (2021b) use the result from Equation 22 to show the following for conditional reverse transition kernels.

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

$$= \frac{\hat{q}(\mathbf{x}_t, \mathbf{x}_{t+1}, \mathbf{y})}{\hat{q}(\mathbf{y}|\mathbf{x}_{t+1})\hat{q}(\mathbf{x}_{t+1})} \quad (24)$$

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

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

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

Further, we can show the following using Equations 18 and 19 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 .

$$\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y}) = \prod_{t=1}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y}) \quad (28)$$

$$= \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (29)$$

$$= q(\mathbf{x}_{1:T}|\mathbf{x}_0) \quad (30)$$

We adapt the derivation of diffusion loss from Luo (2022) 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).

$$\log p(\mathbf{x}|\mathbf{y}) = \log \int p(\mathbf{x}_{0:T}|\mathbf{y}) d\mathbf{x}_{1:T} \quad (31)$$

$$= \log \int \frac{p(\mathbf{x}_{0:T}|\mathbf{y})\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})}{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} d\mathbf{x}_{1:T} \quad (32)$$

$$= \log \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[\frac{p(\mathbf{x}_{0:T}|\mathbf{y})}{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \right] \quad (33)$$

$$\geq \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[\log \frac{p(\mathbf{x}_{0:T}|\mathbf{y})}{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \right] \quad (34)$$

The ELBO can be further simplified as follows

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

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[\log \frac{p(\mathbf{x}_T|\mathbf{y}) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\prod_{t=1}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y})} \right] \quad (36)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[\log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y}) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y}) \prod_{t=2}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{y})} \right] \quad (37)$$

$$= \mathbb{E}_{\hat{q}(\mathbf{x}_{1:T}|\mathbf{x}_0, \mathbf{y})} \left[\log \frac{p(\mathbf{x}_T|\mathbf{y}) p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y}) \prod_{t=2}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})}{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y}) \prod_{t=2}^T \hat{q}(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0, \mathbf{y})} \right] \quad (38)$$

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

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

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

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

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

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

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

$$= \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_1|\mathbf{x}_0, \mathbf{y})} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1, \mathbf{y})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(\hat{q}(\mathbf{x}_T|\mathbf{x}_0, \mathbf{y}) \parallel p(\mathbf{x}_T|\mathbf{y}))}_{\text{prior matching term}} - \underbrace{\sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}))]}_{\text{denoising matching term}} \quad (46)$$

The prior matching term does not contain any trainable parameters. We further simplify the denoising matching term using Equation 27 further conditioned on \mathbf{x}_0 .

$$- \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}))] \quad (47)$$

$$= - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} [\mathbb{E}_{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} [\log \hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y}) - \log p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})]] \quad (48)$$

$$= - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0, \mathbf{y})} \left[\mathbb{E}_{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{y})} \left[\log \hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) + \log \frac{\hat{q}(\mathbf{y}|\mathbf{x}_{t-1}, \mathbf{x}_0)}{\hat{q}(\mathbf{y}|\mathbf{x}_t, \mathbf{x}_0)} - \log p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}) \right] \right] \quad (49)$$

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

Using the result of Equation 50, Equation 46 can be rewritten as

$$\mathbb{E}_{\hat{q}(\mathbf{x}_1|\mathbf{x}_0,\mathbf{y})} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1,\mathbf{y})] - D_{\text{KL}}(\hat{q}(\mathbf{x}_T|\mathbf{x}_0,\mathbf{y}) \parallel p(\mathbf{x}_T|\mathbf{y})) - \sum_{t=2}^T \mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0,\mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{x}_0,\mathbf{y}) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{y}))] \quad (51)$$

$$= \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_1|\mathbf{x}_0,\mathbf{y})} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1,\mathbf{y})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(\hat{q}(\mathbf{x}_T|\mathbf{x}_0,\mathbf{y}) \parallel p(\mathbf{x}_T|\mathbf{y}))}_{\text{prior matching term}} - \sum_{t=2}^T \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0,\mathbf{y})} [D_{\text{KL}}(\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{y}))]}_{\text{denoising matching term}} \quad (52)$$

$$- \sum_{t=2}^T \underbrace{\mathbb{E}_{\hat{q}(\mathbf{x}_t|\mathbf{x}_0,\mathbf{y})} \left[\mathbb{E}_{\hat{q}(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{x}_0,\mathbf{y})} \left[\log \frac{\hat{q}(\mathbf{y}|\mathbf{x}_{t-1},\mathbf{x}_0)}{\hat{q}(\mathbf{y}|\mathbf{x}_t,\mathbf{x}_0)} \right] \right]}_{\text{label consistency term}}$$

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 (2022) that the reconstruction term and the denoising matching term when developed further lead to the diffusion loss of Equation 7, 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 30.

A.2 PROOF FOR THEOREM 3.2

Theorem A.2. *Denoising score matching for the guidance model π_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 10 is equal up to a constant to the following loss:*

$$L_{\text{DLMS}}(\phi) = \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x},\mathbf{x}_t,\mathbf{y}^{1:k+1},\tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) + \nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t|\tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t|\mathbf{x})\|_2^2 \right] \quad (53)$$

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 τ^2 that is small enough such that $p_{\tau}(\tilde{\mathbf{y}}^i) \approx p(\mathbf{y}^i)$.

Chao et al. (2022) in their insightful work show that the following two losses differ only by a constant.-

$$D_F(p_{\phi}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t) \parallel p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t)) = \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t,\tilde{\mathbf{y}}^{k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t)\|_2^2 \right] \quad (54)$$

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

We extend their proof for multiple conditionals below. The Fisher divergence between the guidance model and the true score of the classifier (Equation 10) can be further expanded as:

$$D_F(p_{\phi}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) \parallel p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) \quad (56)$$

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

$$= \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t,\tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] + \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t,\tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] - \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t,\tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1}|\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) \rangle] \quad (58)$$

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

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

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 \\
& = - \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 60, we get-

$$\begin{aligned}
& D_F(p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) || p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) \\
& = \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] + \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] \\
& + \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha,\tau}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) \rangle] \\
& - \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_0, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0) \rangle] \tag{61}
\end{aligned}$$

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 61, we get:

$$\begin{aligned}
& D_F(p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) || p_{\alpha,\tau}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})) \\
& = \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})\|_2^2 \right] + \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha,\tau}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) \rangle] + C \\
& - \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_0, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} [\langle \nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}), \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0) \rangle] + \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}_t, \tilde{\mathbf{y}}^{1:k})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\alpha,\tau}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x}_0)\|_2^2 \right] \tag{62}
\end{aligned}$$

$$= \mathbb{E}_{p_{\alpha,\tau}(\mathbf{x}, \mathbf{x}_t, \mathbf{y}^{1:k+1}, \tilde{\mathbf{y}}^{1:k+1})} \left[\frac{1}{2} \|\nabla_{\mathbf{x}_t} \log p_{\phi}(\tilde{\mathbf{y}}^{k+1} | \mathbf{x}_t, \tilde{\mathbf{y}}^{1:k}) + \nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t | \tilde{\mathbf{y}}^{1:k}) - \nabla_{\mathbf{x}_t} \log p_{\alpha}(\mathbf{x}_t | \mathbf{x})\|_2^2 \right] + C \tag{63}$$