
Understanding Diffusion-based Representation Learning via Low-Dimensional Modeling

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This work addresses the critical question of why and when diffusion models, despite
2 their generative design, are capable of learning high-quality representations in a self-
3 supervised manner. We hypothesize that diffusion models excel in representation
4 learning due to their ability to learn the low-dimensional distributions of image
5 datasets via optimizing a noise-controlled denoising objective. Our empirical
6 results support this hypothesis, indicating that variations in the representation
7 learning performance of diffusion models across noise levels are closely linked to
8 the quality of the corresponding posterior estimation. Grounded on this observation,
9 we offer theoretical insights into the unimodal representation dynamics of diffusion
10 models as noise scales vary, demonstrating how they effectively learn meaningful
11 representations through the denoising process. We also highlight the impact of
12 the inherent parameter-sharing mechanism in diffusion models, which accounts
13 for their advantages over traditional denoising auto-encoders in representation
14 learning.

15 1 Introduction

16 Diffusion models, a new family of likelihood-based generative models, have demonstrated superior
17 performance among many generative tasks, including image generation [Alkhouri et al., 2024, Ho
18 et al., 2020, Rombach et al., 2022, Zhang et al., 2024], video generation [Bar-Tal et al., 2024, Ho
19 et al., 2022], speech and audio synthesis [Kong et al., 2020, 2021], semantic editing [Roich et al.,
20 2022, Ruiz et al., 2023, Chen et al., 2024a] and solving inverse problem [Chung et al., 2022, Song
21 et al., 2024, Li et al., 2024, Alkhouri et al., 2023]. At its core, diffusion models are learning a data
22 distribution from training samples by imitating the non-equilibrium thermodynamic diffusion process
23 [Sohl-Dickstein et al., 2015, Ho et al., 2020, Song et al., 2021]. In the forward process, training
24 samples are gradually combined with increasing Gaussian noise until the data structure is completely
25 destroyed while in the backward process, a model is trained to restore the structure from the noised
26 data [Hyvärinen and Dayan, 2005, Song et al., 2021].

27 In addition to their impressive generative capabilities, recent studies [Baranchuk et al., 2021, Xiang
28 et al., 2023, Mukhopadhyay et al., 2023, Chen et al., 2024b, Tang et al., 2023] have highlighted the
29 exceptional representation power of diffusion models, suggesting that they could serve as a unified
30 foundation model for both generative and discriminative vision tasks. Specifically, recent evaluations
31 across various applications, including classification [Xiang et al., 2023, Mukhopadhyay et al., 2023],
32 semantic segmentation [Baranchuk et al., 2021], and image alignment [Tang et al., 2023], show
33 that diffusion models are capable of learning high-quality representations, often matching or even
34 surpassing the performance of previous state-of-the-art methods. However, it remains unclear whether
35 the representation capabilities of diffusion models stem from the diffusion process or the denoising

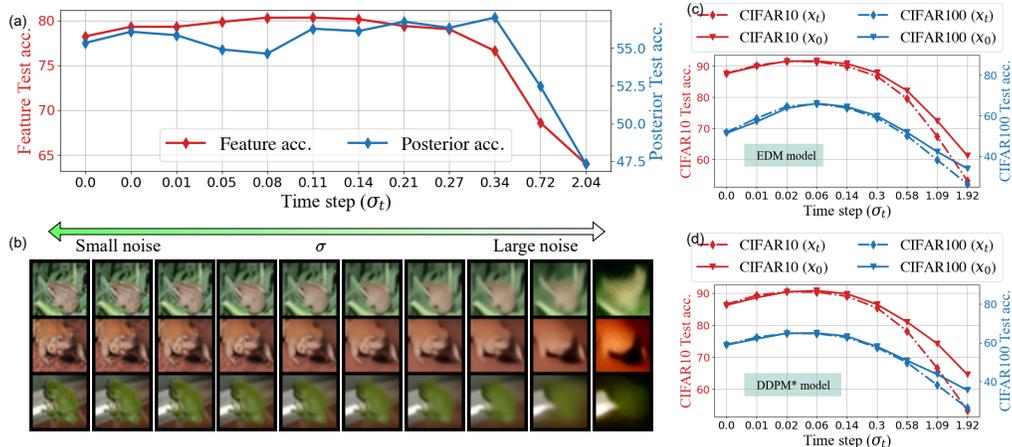


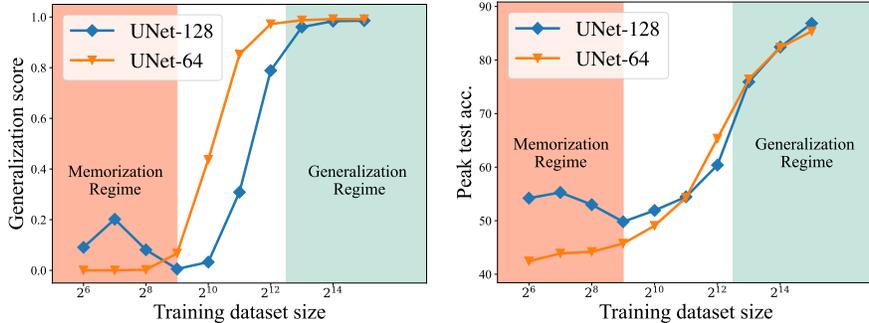
Figure 1: **Representation learning ability of a diffusion model at different time steps reflects the granularity in posterior estimation.** (a) Intermediate feature accuracy and posterior accuracy of the diffusion model exhibit a similar unimodal trend as noise level increases. (b) Posterior estimation for clean image inputs shows a transition from fine to coarse granularity with increasing noise levels. (c)-(d) Using clean image input x_0 for feature extraction achieves comparable or superior representation learning performance compared to using noisy input x_t .

36 mechanism [Fuest et al., 2024]. More fundamentally, given their generative design, *when and why*
 37 *diffusion models can learn high-quality representations in a self-supervised manner?*

38 This work aims to address this question through a comprehensive investigation, both empirically and
 39 theoretically, grounded in the formulation of denoising auto-encoders (DAEs) for learning diffusion
 40 models [Vincent et al., 2008, 2010, Vincent, 2011]. We hypothesize that diffusion models can learn
 41 high-quality representations without supervision due to their superior ability to approximate the
 42 low-dimensional distributions of image datasets, as supported by recent findings [Wang et al., 2024].
 43 Although image dataset can be very high-dimensional, recent results [Pope et al., 2021, Stanczuk
 44 et al., 2022, Wang et al., 2024] demonstrate that the intrinsic dimension of these datasets are much
 45 lower than the ambient dimension, and it has shown that the number of samples to learn the underlying
 46 distribution using diffusion models scales with the intrinsic low-dimensionality. Therefore, by being
 47 trained to capture the underlying structure of data through a controlled process of noise injection and
 48 denoising, diffusion models effectively learn meaningful and compact features.

49 On the empirical side, we support our claim by reconciling several intriguing phenomena related
 50 to the quality of learned representations in diffusion models. Recent studies Zhang et al. [2023]
 51 reveal that diffusion models operate in two regimes: memorization and generalization, depending on
 52 training data size. In the memorization regime with limited samples, the model captures only the
 53 empirical distribution of training data without the ability to generate new samples. In contrast, in the
 54 generalization regime, diffusion models are able to learn the underlying distribution. Our experiments
 55 in Figure 2 confirm that high-quality representations are *only* learned in the generalization regime with
 56 sufficient samples due to its ability of learning the underlying distribution. More importantly, in the
 57 generalization regime, we show that the quality of hidden representations in diffusion models/DAEs
 58 follows a uni-modal curve (see Figure 1 and Figure 7): high-quality representations are learned
 59 at an intermediate step close to the clean image, whereas the representation quality degrades as it
 60 approaches either pure noise or the clean image.

61 Building on these empirical observations, we provide theoretical insights using a noisy mixture of
 62 low-rank Gaussian distributions. Our assumption captures the inherent low-dimensionality of the
 63 image data distribution [Pope et al., 2021, Gong et al., 2019, Stanczuk et al., 2022], where the data
 64 lies on a union of low-dimensional subspaces. We analyze the unimodal trend in representation
 65 performance by relating it to the Class-specific Signal-to-Noise Ratio (CSNR). Specifically, we
 66 consider the optimal posterior estimation function under our data assumption and show that the
 67 CSNR is determined by the interplay between data “denoising” and class confidence rate as the
 68 noise scale increases. Additionally, our study reveals an implicit weight-sharing mechanism inherent
 69 in diffusion models, which helps explain their strengths compared to traditional one-step DAEs,
 70 particularly in the small noise regions.



(a) Phase transition in generalization score (b) Phase transition in representation learning

Figure 2: **Better representations are learned in the generalization regime.** We train EDM-based [Karras et al., 2022] diffusion models on the CIFAR-10 dataset using different training dataset sizes, ranging from 2^6 to 2^{15} . (a) The change in the generalization score [Zhang et al., 2023] as the dataset size increases, where regions with a generalization score close to 0 are labeled as the memorization regime, and those close to 1 are labeled as the generalization regime. (b) The peak representation learning accuracy achieved as a function of dataset size.

71 **Contribution of this work.** In summary, our findings can be highlighted as follows:

- 72 • **Linking posterior estimation ability of diffusion models to representation learning.** Our
73 empirical results reveal that, much like the dynamics of diffusion representation learning, posterior
74 estimation quality across noise levels follows a similar unimodal curve. This indicates that changes
75 in representation quality are a direct reflection of changes in posterior estimation quality, prompting
76 us to explore representation learning through the more fundamental lens of posterior recovery.
- 77 • **Theoretical analysis of the unimodal curve in the denoising process.** Building on the connection
78 between posterior estimation and representation learning, we present the first theoretical framework
79 for analyzing the unimodal evolution of representation quality. Using a mixture of low-rank
80 Gaussian data model, we demonstrate that the unimodal curve arises from the interplay between
81 denoising strength and class confidence as the noise level varies.
- 82 • **Weight sharing in the diffusion process.** Furthermore, we reveal that the diffusion process, by
83 minimizing losses across all noise levels simultaneously, fosters an implicit parameter sharing
84 mechanism within a diffusion model. This mechanism plays a crucial role for diffusion models
85 to achieve superior and more consistent representation learning performances compared with
86 traditional DAEs.

87 2 Representation Learning via diffusion models

88 In this section, we first review the fundamentals of diffusion models and outline the feature extraction
89 method used in this work. Following this, we illustrate the connection between diffusion posterior
90 estimation and representation learning, which serves as the foundation for the subsequent analysis in
91 Section 3.

92 2.1 Preliminaries on denoising diffusion models

93 Diffusion models are a class of probabilistic generative models that aim to reverse a progressive
94 noising process by mapping an underlying data distribution, p_{data} , to a Gaussian distribution.

95 **The forward process.** Starting from clean data \mathbf{x}_0 , noise is gradually introduced according to a
96 noise schedule determined by the time step t until the data becomes indistinguishable from pure
97 Gaussian noise. Specifically, at any time step t , the noised data can be expressed as: $\mathbf{x}_t = s_t \mathbf{x}_0 + s_t \sigma_t \epsilon$
98 where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ represents noise sampled from a Gaussian distribution, s_t and $s_t \sigma_t$ represent the
99 scaling of the signal and noise, respectively.

100 **The reverse process.** Noise is gradually removed from \mathbf{x}_1 following the reverse-time SDE:

$$d\mathbf{x}_t = (f(t)\mathbf{x}_t - g^2(t)\nabla \log p_t(\mathbf{x}_t)) dt + g(t)d\bar{\mathbf{w}}_t, \quad (1)$$

101 where $\{\bar{\mathbf{w}}_t\}_{t \in [0,1]}$ is the standard Wiener process running backward in time from $t = 1$ to $t = 0$ and
102 the functions $f(t), g(t) : \mathbb{R} \rightarrow \mathbb{R}$ respectively denote the drift and diffusion coefficients. Notably, if
103 both \mathbf{x}_1 and $\nabla \log p_t$ are known, the reverse process mirrors the forward process at each time step
104 $t \geq 0$ [Anderson, 1982].

105 **Score approximation and denoising auto-encoders (DAEs).** However, the score function $\nabla \log p_t$
 106 is typically unknown, as it depends on the underlying data distribution p_{data} . To address this, a neural
 107 network s_θ is trained to estimate the score at various time steps [Ho et al., 2020, Song et al., 2021].
 108 Given the relationship between the score function and the posterior mean $\mathbb{E}[\hat{x}_0|\mathbf{x}_t]$ [Vincent, 2011,
 109 Wang et al., 2024]:

$$s_t \mathbb{E}[\hat{x}_0|\mathbf{x}_t] = \mathbf{x}_t + s_t^2 \sigma_t^2 \nabla \log p_t(\mathbf{x}_t) \approx \mathbf{x}_t + s_t^2 \sigma_t^2 s_\theta(\mathbf{x}_t), \quad (2)$$

110 prior works [Chen et al., 2024b, Xiang et al., 2023, Kadkhodaie et al., 2023] have also proposed an
 111 alternative DAE-based training objective that directly estimates the posterior mean $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$:

$$\min_{\theta} \ell(\theta) := \frac{1}{2N} \sum_{i=1}^N \int_0^1 \lambda_t \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)} \left[\left\| \mathbf{x}_\theta(s_t \mathbf{x}_0^{(i)} + s_t \sigma_t \epsilon, t) - \mathbf{x}_0^{(i)} \right\|^2 \right] dt, \quad (3)$$

112 where $\mathbf{x}_\theta(\mathbf{x}_0, t)$ denotes the posterior estimating network, N represents the size of the training
 113 dataset, and λ_t denotes the weighting for each noise level. To simplify the analysis, we assume
 114 throughout the paper that $s_t = 1$ and λ_t remain constant across all noise levels, with the noise level
 115 denoted as σ_t .

116 We note that if we remove the integration in (3) and fix t , the loss simplifies to the traditional
 117 single-level DAE loss [Vincent et al., 2008], where the DAE is trained at a single noise level.
 118 Previous work [Chen et al., 2024b] has decomposed the training objective of diffusion models into the
 119 denoising process (through the denoising loss) and the diffusion process (integrating the loss across
 120 all noise levels in (3)). To comprehensively investigate the distinct roles of these two processes in
 121 representation learning, we consider both diffusion models and individual DAEs in our experiments
 122 where the individual DAEs serve as a control group, allowing us to isolate and analyze the effects of
 123 the denoising process alone.

124 2.2 Extracting representations from diffusion model

125 In this work, we adopt the following feature extraction setups to leverage diffusion models for
 126 representation learning:

127 **Use clean images as network inputs.** First, we use the clean image \mathbf{x}_0 as input to the network
 128 in contrast to conventional approaches that use the noisy image \mathbf{x}_t [Xiang et al., 2023, Baranchuk
 129 et al., 2021, Tang et al., 2023]. This setup aligns with the goal of representation learning, where
 130 additive noise is not necessary (e.g., similar to training a classifier with data augmentations while
 131 using non-augmented data during inference). As demonstrated in Figure 1(c)-(d), this approach
 132 preserves the overall unimodal representation dynamic while achieving better performance at higher
 133 noise levels. As such, throughout the remainder of this paper, we use the clean data \mathbf{x}_0 as input to the
 134 diffusion model, i.e., we always consider $\mathbf{x}_\theta(\mathbf{x}_0, t)$ where t serves solely as an indicator of the noise
 135 level for diffusion model to adopt during feature extraction.

136 **Layer selection for representations.** Second, we extract features only from the bottleneck layer
 137 of the U-Net architecture [Ronneberger et al., 2015],¹ following the protocols used in [Kwon et al.,
 138 2022, Park et al., 2023].² Unlike prior methods [Xiang et al., 2023, Baranchuk et al., 2021], we do
 139 not conduct a grid search for the optimal layer, as our focus is on understanding the process rather
 140 than achieving state-of-the-art results.

141 2.3 Relationship Between Learned Representations & Posterior Estimation

142 Relationship among posterior estimation, distribution recovery, and representation learning.

143 Since directly studying representation ability is challenging, in Section 3 we approach the problem
 144 through its strong correlation with posterior mean estimation, $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$. As we will argue, diffusion
 145 representation quality is closely linked with the semantic information encoded in the posterior
 146 estimation. Additionally, empirical validations can be found in Figure 1.

- 147 • *Posterior estimation and distribution recovery.* Diffusion models are trained to learn the underlying
 148 data distribution by reconstructing the posterior mean $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$ for a given input \mathbf{x}_t at the specified
 149 noise level. Therefore, the quality of posterior estimation $\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$ reflects the degree to which the
 150 underlying distribution is captured [Choi et al., 2022].

¹In other words, the layer with the smallest feature resolution.

²After feature extraction, we apply a global average pooling to the features. For instance, given a feature map
 of dimension $256 \times 4 \times 4$, we pool the last two dimensions, resulting in a 256-dimensional vector.

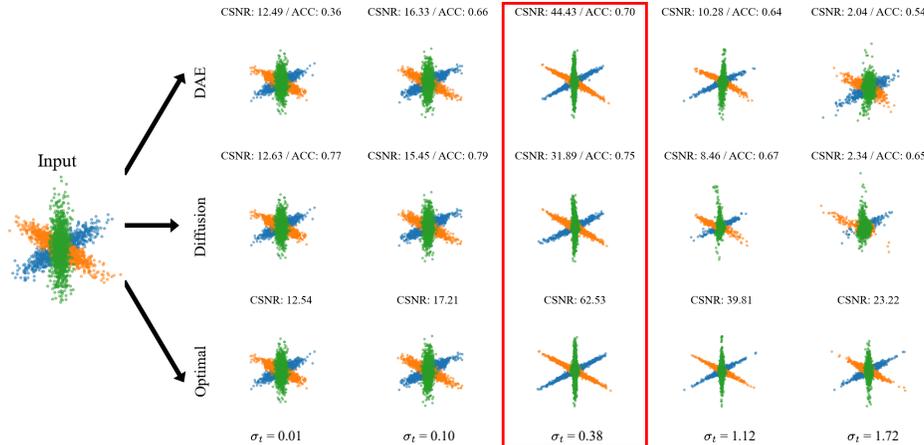


Figure 3: **Visualization of posterior estimation for a clean input.** The same MoLRG data is fed into the models; each row represents a different denoising model, and each column corresponds to a different time step with noise scale (σ_t). The red box indicates the best posterior estimation and feature probing accuracy.

151 • *Representation learning through distribution approximation.* On the other hand, achieving high-
 152 quality distribution approximation results in more meaningful and informative representations in
 153 unsupervised learning. This is supported by Figure 2, where the findings, inspired by recent works
 154 [Zhang et al., 2023], demonstrate that diffusion models transition from memorizing the training
 155 data distribution to accurately approximating the underlying data distribution as the amount of
 156 training data increases. Consequently, better approximation of the underlying data distribution
 157 improves the quality of representation learning.

158 Given this relationship, we use posterior estimation as a proxy for representation quality throughout
 159 our analysis. Additionally, since diffusion models tend to memorize the training data instead of
 160 learning underlying data distribution when the training dataset is small [Zhang et al., 2023], we focus
 161 on the case where sufficient training data is available throughout our analysis in Section 3.

162 **Unimodal curve of representation quality.** Previous studies [Xiang et al., 2023, Baranchuk et al.,
 163 2021, Tang et al., 2023] have empirically shown that the representation dynamics of diffusion models
 164 follow a unimodal curve as the noise scale increases, across various tasks such as classification,
 165 segmentation, and image correspondence. Our findings corroborate this observation, as demonstrated
 166 in Figure 1(a), where the representation quality consistently exhibits a unimodal trend, regardless
 167 of the specific network architecture or dataset used (see Figure 1(c)-(d)). In the following analysis,
 168 we argue that this unimodal behavior arises from subtle differences between the requirements of
 169 representation learning and the generative nature of diffusion models.

170 High-fidelity image generation demands that diffusion models capture every aspect of the data
 171 distribution—from coarse structures to fine details. In contrast, representation learning, particularly
 172 for high-level tasks such as classification [Allen-Zhu and Li, 2022], prefers an abstract representation,
 173 where finer image details may even act as ‘noise’ that hinders performance. As shown in Figure 1(b),
 174 as the noise level increases, the predicted posteriors for clean input x_0 transition from ‘fine’ to ‘coarse’
 175 [Wang and Vastola, 2023], gradually removing fine-grained details. For the classification task in the
 176 plot, the best performance is achieved when the posterior estimation retains the essential information
 177 while discarding some class-irrelevant details. These findings indicate a trade-off between generative
 178 quality and representation performance [Chen et al., 2024b], prompting us to attribute variations in
 179 feature quality across noise levels to differences in posterior prediction.

180 3 Theoretical Understanding Through Low-Dimensional Models

181 In this section, we theoretically examine the representation learning capabilities of diffusion models
 182 across varying noise levels by evaluating the quality of posterior estimation, $\mathbb{E}[x_0|x_t]$ for low-
 183 dimensional distributions.

184 **3.1 Assumptions of Low-Dimensional Data Distribution**

185 Although real-world image datasets are high-dimensional in terms of pixel count and data volume,
 186 extensive empirical studies Gong et al. [2019], Pope et al. [2021], Stanczuk et al. [2022] suggest that
 187 their intrinsic dimensionality is considerably lower. Moreover, state-of-the-art large-scale diffusion
 188 models [Peebles and Xie, 2023, Podell et al., 2023] commonly employ auto-encoders [Kingma, 2013]
 189 to map images to a low-dimensional latent space [Rombach et al., 2022] for better training efficiency.
 190 Consequently, image datasets often reside on a union of low-dimensional manifolds.

191 In light of this, many recent studies of diffusion models have been focused on approximating low-
 192 dimensional distributions [Wang et al., 2024]. Moreover, as union of low-dimensional manifolds can
 193 be locally approximated by a union of linear subspaces, it motivates us to model the underlying data
 194 distribution as a mixture of low-rank Gaussians (MoLRG). The data points generated by MoLRG lie on
 195 a union of subspaces. Within each subspace, the data follows a Gaussian distribution with a low-rank
 196 covariance matrix that represents the subspace basis. Formally, we introduce a noisy version of the
 197 MoLRG distribution as follows:

198 **Assumption 1** (*K*-Subspace Noisy MoLRG Distribution). *For any sample \mathbf{x}_0 drawn from the noisy*
 199 *MoLRG distribution with K subspaces, the following holds:*

$$\mathbf{x}_0 = \mathbf{U}_k \mathbf{a} + \delta \mathbf{U}_k^\perp \mathbf{e}, \text{ with probability } \pi_k \geq 0, k \in [K]. \quad (4)$$

200 Here, $\sum_{k=1}^K \pi_k = 1$, $\mathbf{U}_k \in \mathcal{O}^{n \times d_k}$ denotes an orthonormal basis for the k -th subspace, d_k is the
 201 subspace dimension with $d_k \ll n$, and the coefficient $\mathbf{a} \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_{d_k})$ is drawn from a standard
 202 normal distribution. For the noise, we assume $\mathbf{e} \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_{n-d_k})$ with magnitude controlled by the
 203 scalar $\delta < 1$. Additionally, $\mathbf{U}_k^\perp \in \mathcal{O}^{n \times (n-d_k)}$ is the orthogonal compliment of \mathbf{U}_k .

204 For simplicity of analysis, we let $d_1 = \dots = d_K = d$, and we assume that the basis $\{\mathbf{U}_k\}$ are
 205 orthogonal to each other with $\mathbf{U}_k^T \mathbf{U}_l = \mathbf{0}$ for all $k \neq l$. Additionally, we assume all mixing weights
 206 $\{\pi_k\}$ are equal with $\pi_1 = \dots = \pi_K = 1/K$, and we define $\mathbf{U}_\perp = \bigcap_{k=1}^K \mathbf{U}_k^\perp \in \mathcal{O}^{n \times (n-Kd)}$ to be
 207 the noise space that is the orthogonal complement to all basis $\{\mathbf{U}_k\}_{k=1}^K$.

208 We note that the noise term $\delta \mathbf{U}_k^\perp \mathbf{e}_i$ captures perturbations unrelated to the k -th subspace via the
 209 orthogonal complement \mathbf{U}_k^\perp , thereby aligning the model more closely with real-world scenarios.
 210 These perturbations can be interpreted as attributes irrelevant to the subspace, such as the background
 211 in an image of a bird or the color/texture of a car. The extra noise term may not be relevant for
 212 representation learning, but it plays an importance role for diffusion model to generate high-fidelity
 213 samples. Additionally, for the noisy MoLRG distribution, ground truth posterior mean $\mathbb{E}[\hat{\mathbf{x}}_0 | \mathbf{x}_t]$ is:

214 **Proposition 1.** *For a K -class MoLRG data distribution, for each time $t > 0$, it holds that*

$$\mathbf{x}_\theta^*(\mathbf{x}_t, t) := \mathbb{E}[\hat{\mathbf{x}}_0 | \mathbf{x}_t] = \sum_{k=1}^K w_k^*(\mathbf{x}_t) \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \right) \mathbf{x}_t \quad (5)$$

$$\text{where } w_k^*(\mathbf{x}_t) := \frac{\exp(g_k(\mathbf{x}_t, t))}{\sum_{k=1}^K \exp(g_k(\mathbf{x}_t, t))}, \quad (6)$$

$$\text{and } g_k(\mathbf{x}) = \frac{1}{2\sigma_t^2(1 + \sigma_t^2)} \|\mathbf{U}_k^T \mathbf{x}\|^2 + \frac{\delta^2}{2\sigma_t^2(\delta^2 + \sigma_t^2)} \|\mathbf{U}_k^{\perp T} \mathbf{x}\|^2. \quad (7)$$

215 **Remark.** In the above proposition, we present the ground truth posterior estimation function that
 216 a diffusion model can achieve by minimizing the training objective defined in (3). We denote this
 217 optimal model \mathbf{x}_θ^* . Given the established relationship between posterior estimation and representation
 218 learning on clean inputs \mathbf{x}_0 , we can now analyze the representation learning dynamics under this
 219 optimal setting by evaluating $\mathbf{x}_\theta^*(\mathbf{x}_0, t)$ at different time step t .

220 **3.2 Main Theoretical Results**

221 As we discussed in Section 2.3, based upon the strong correlation between representation quality and
 222 the posterior mean estimation, we analyze $\mathbf{x}_\theta^*(\mathbf{x}_0, t)$ across different time step $t \in [0, 1]$. Here, we
 223 use \mathbf{x}_0 as the input instead of \mathbf{x}_t according to our discussion in Section 2.2.

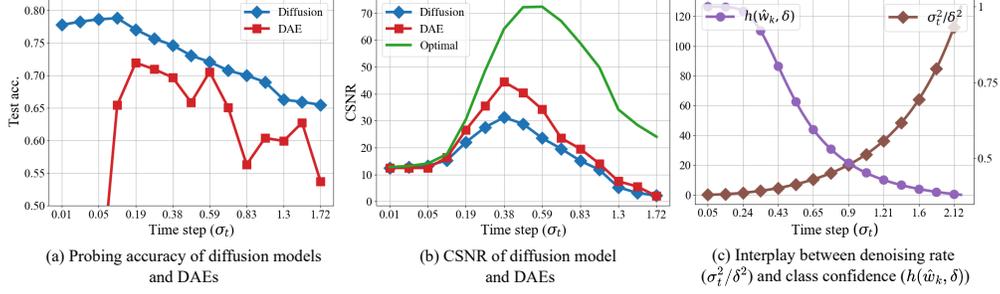


Figure 4: **Dynamics of feature probing accuracy, CSNR, and denoising/class confidence rate with increasing noise levels.** Panels (a) and (b) show the feature probing accuracy and CSNR trends using the same MoLRG data as in Figure 3, both exhibiting a unimodal pattern. The interplay between the “denoising rate” and the class confidence rate for the approximate optimal solution f^* is illustrated in panel (c).

224 Given $\mathbf{x}_0 \sim \text{MoLRG}$ and without loss of generality, let k represent the true class to which \mathbf{x}_0 belongs.
 225 We quantify the accuracy of posterior mean estimation by introducing a measure of Class-specific
 226 Signal-to-Noise Ratio (CSNR) as follows:

$$\text{CSNR}(t, \mathbf{x}_\theta^*) := \frac{\mathbb{E}_{\mathbf{x}_0}[\|\mathbf{U}_k \mathbf{U}_k^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|^2]}{\mathbb{E}_{\mathbf{x}_0}[\sum_{l \neq k} \|\mathbf{U}_l \mathbf{U}_l^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|^2]} \quad (8)$$

227 We know that successful prediction of the class for \mathbf{x}_0 occurs when the class-specific signal
 228 $\|\mathbf{U}_k \mathbf{U}_k^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|$ dominates over the noise term $\|\mathbf{U}_\perp \mathbf{U}_\perp^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|$. On the other hand, be-
 229 cause

$$\|\mathbf{U}_\perp \mathbf{U}_\perp^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|^2 = \sum_{l \neq k} \|\mathbf{U}_l \mathbf{U}_l^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|^2 + \|\mathbf{U}_\perp \mathbf{U}_\perp^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|^2$$

230 and \mathbf{U}_\perp does not affect classification due to its presence in every data point, it leads to our definition
 231 of CSNR in (8) which measures the ratio between the true class signal and irrelevant noise from
 232 other classes.

233 Therefore, intuitively, a higher CSNR indicates a better recovery of the underlying low-dimensional
 234 data subspace, and thus the predicted posterior is more likely to be assigned to the correct class. This
 235 is supported by Figure 4(a)-(b) which shows that both $\text{CSNR}(t)$ and classification accuracy using the
 236 learned representation follow similar unimodal curves.

237 To simplify the calculation of (8), which involves the expectation over the softmax term w_k^* , we
 238 approximate \mathbf{x}_θ^* as follows:

$$f^*(\mathbf{x}, t) = \sum_{k=1}^K \hat{w}_k \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \right) \mathbf{x}, \quad (9)$$

$$\text{where } \hat{w}_k := \frac{\exp(\mathbb{E}_{\mathbf{x}_0}[g_k(\mathbf{x}_0, t)])}{\sum_{k=1}^K \exp(\mathbb{E}_{\mathbf{x}_0}[g_k(\mathbf{x}_0, t)])}$$

239 In other words, we use \hat{w}_k in (9) to approximate $w_k^*(\mathbf{x}_0)$ in (6) by taking expectation inside the
 240 softmax with respect to \mathbf{x}_0 . This allows us to treat \hat{w}_k as a constant when calculating CSNR, making
 241 the analysis more tractable while maintaining $\mathbb{E}[\|\mathbf{U}_l \mathbf{U}_l^T \mathbf{x}_\theta^*(\mathbf{x}_0, t)\|^2] \approx \mathbb{E}[\|\mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}, t)(\mathbf{x}_0, t)\|^2]$
 242 for all $l \in [K]$. We verify the tightness of this approximation at Appendix A.3 (Figure 9). Now, we
 243 are ready to state our main theorem as follows.

244 **Theorem 1.** Let data \mathbf{x}_0 be any arbitrary data point drawn from the MoLRG distribution defined in
 245 Assumption 1 and let k denote the true class \mathbf{x}_0 belongs to. Then CSNR introduced in (8) depends
 246 on the noise level σ_t in the following form:

$$\text{CSNR}(t, f^*) = \frac{1}{(K-1)\delta^2} \cdot \left(\frac{1 + \frac{\sigma_t^2}{\delta^2} h(\hat{w}_k, \delta)}{1 + \frac{\sigma_t^2}{\delta^2} h(\hat{w}_l, \delta)} \right)^2 \quad (10)$$

247 where $h(w, \delta) := (1 - \delta^2)w + \delta^2$. Since δ is fixed, $h(w, \delta)$ is a monotonically increasing function
 248 with respect to w . Note that here δ represents the magnitude of the fixed intrinsic noise in the data
 249 where σ_t denotes the level of additive Gaussian noise introduced during the diffusion training process.

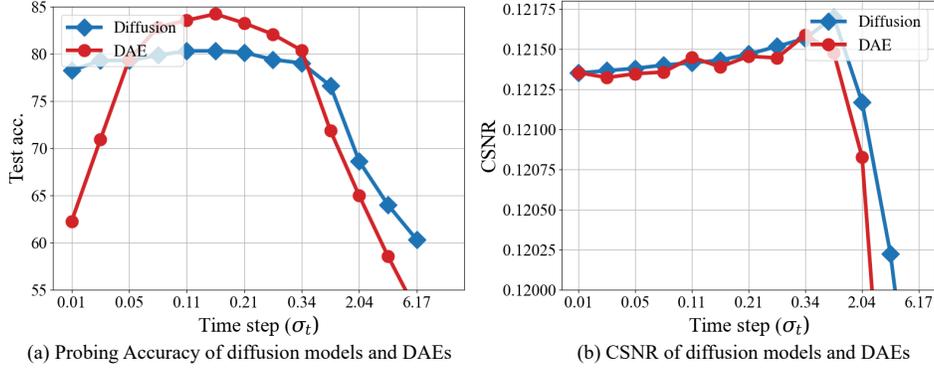


Figure 5: **Dynamics of feature probing accuracy and CSNR on CIFAR10.** Panels (a) and (b) show the feature probing accuracy and CSNR trends computed using the CIFAR10 test dataset, both exhibiting a unimodal pattern.

250 **Remark.** Intuitively, the unimodal curve of CSNR reflects how the additive noise level σ_t in the
 251 diffusion process helps counteract the intrinsic data noise δ . The noise ratio (σ_t/δ) can be interpreted
 252 as the “denoising” rate, where a larger ratio indicates more data noise being canceled out and vice
 253 versa. Meanwhile, $h(\hat{w}_k, \delta)$ represents the class confidence rate, with lower values meaning less class-
 254 specific information is captured by the model. With σ_t increases from 0 to ∞ , the “denoising rate”
 255 rises accordingly, while the class confidence rate decreases monotonically. Thus, from Theorem 1,
 256 we can derive the rationale behind the unimodal behavior of CSNR.

257 • **The unimodal curve of CSNR.** The unimodal curve is decided by the interplay between the
 258 “denoising rate” and the class confidence rate as noise increases. As observed in Figure 4(c), the
 259 “denoising rate” (σ_t^2/δ^2) increases monotonically with σ_t while the class confidence rate $h(\hat{w}_k, \delta)$
 260 monotonically declines. Initially, as σ_t increases, the class confidence rate remains relatively
 261 stable due to its flat slope (as seen in Figure 4(c)), and an increasing “denoising rate” enhances the
 262 CSNR, resulting in improved posterior estimation. However, as indicated by (7), when σ_t becomes
 263 too large, $h(\hat{w}_k, \delta)$ approaches $h(\hat{w}_l, \delta)$, leading to a drop in CSNR, which limits the model’s
 264 ability to project \mathbf{x}_0 onto the correct signal space and ultimately impairs posterior estimation. This
 265 interpretation is validated by the visualization in Figure 3. In the plot, each class is represented by
 266 a colored straight line, while deviations from these lines correspond to the δ -related noise term.
 267 Initially, increasing the noise scale effectively cancels out the δ -related data noise, resulting in a
 268 cleaner posterior estimation and improved probing accuracy. However, as the noise continues to
 269 increase, the class confidence rate drops, leading to an overlap between classes, which ultimately
 270 degrades the feature quality and probing performance.

271 Back to our real-world analogy, the proportion of data associated with δ represents class-irrelevant
 272 attributes or finer image details. The unimodal representation learning dynamic thus captures a “fine-
 273 to-coarse” shift [Choi et al., 2022, Wang and Vastola, 2023], where these details are progressively
 274 stripped away. During this process, peak representation performance is achieved at a balance point
 275 where class-irrelevant attributes are eliminated, while class-essential information is preserved.

276 3.3 Empirical Validation

277 In this subsection, we conduct experiments on both synthetic and real datasets to validate our theory
 278 on the representation learning dynamics.

279 We use two datasets: a 3-class MoLRG dataset, where each subspace has dimension $d = 1$ and ambient
 280 dimension $n = 10$, with noise scale $\delta = 0.2$, and the standard CIFAR10 dataset [Krizhevsky et al.,
 281 2009]. We consider two training settings: (a) a DDPM-based diffusion training configuration and
 282 (b) a vanilla DAE training configuration, where separate DAEs are trained for different noise levels.
 283 Here, the separate DAEs serve as a control group, enabling us to isolate the effects of the denoising
 284 process, as discussed in Section 2.1. We leave further training details in Appendix A.2.

285 After training, we extract intermediate features and posterior predictions from both diffusion models
 286 and DAEs, followed by linear probing on the features and computation of empirical CSNR for
 287 the posterior estimations. The results for the two datasets are presented in Figure 4 and Figure 5,
 288 respectively. As shown in the plots, both feature probing accuracy and the empirical CSNR exhibit a

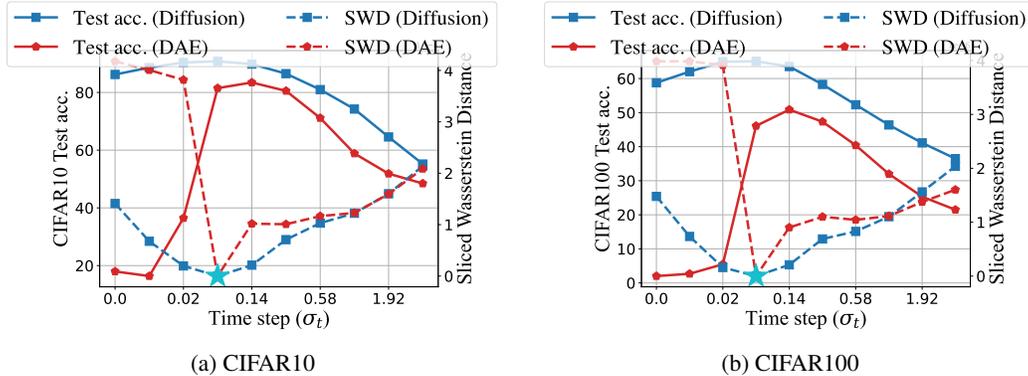


Figure 6: **Comparison of representation learning performance and feature similarity between diffusion model and individual DAEs.** We train DDPM-based diffusion models and individual DAEs on the CIFAR10 and CIFAR100 datasets. After training, we plotted their representation learning performance and feature similarity against the best features (indicated by \star) as the noise level increases.

289 matching unimodal curve, consistent across training configurations and datasets, thus supporting our
 290 theoretical results.

291 4 Additional Experiments

292 In Section 3, we analyzed diffusion representation dynamics with a focus on the denoising process,
 293 assuming sufficient training data for learning the underlying distribution. In this section, we explore
 294 the impact of the diffusion process (Section 4.1) and data complexity (Section 4.2) in shaping
 295 diffusion models’ representation learning dynamics.

296 4.1 Weight sharing in diffusion models helps representation learning

297 In this subsection, we demonstrate how the inherent weight-sharing mechanism in diffusion models,
 298 stemming from their loss design, enhances representation learning performances compared with
 299 traditional DAEs.

300 Previously, in Section 3, we analyzed the optimal posterior function by treating each noise level
 301 independently. However, the training objective for diffusion models in (3) involves minimizing
 302 the loss across all noise levels simultaneously, which results in interactions and parameter sharing
 303 among denoising subcomponents at different noise levels. We hypothesize that these interactions
 304 and parameter sharing create greater feature similarity across noise scales, effectively functioning as
 305 an implicit “ensemble” mechanism that enhances the performance of diffusion models compared to
 306 individual DAEs [Chen et al., 2024b], which accounts for the significant performance gap between
 307 DAEs and diffusion models, as shown in Figure 4(a) and Figure 5(a).

308 To test this hypothesis, we trained 10 individual DAEs, each at a different noise level, as well as a
 309 single DDPM-based diffusion model on CIFAR10 and CIFAR100 datasets. We then conducted linear
 310 probing on the features extracted from both setups. To evaluate feature similarity, we calculated
 311 the sliced Wasserstein distance (SWD) [Doan et al., 2024] between features for both diffusion and
 312 DAE models at various noise levels and their corresponding features at $\sigma_t = 0.06$, which achieves
 313 near-optimal accuracy for all scenarios.

314 As shown in Figure 6, diffusion models consistently outperform individual DAEs, particularly at
 315 lower noise levels, where the performance gap is most pronounced. In these low-noise regions, due
 316 to the almost negligible additive noise, individual DAEs are more likely to be trained as identity
 317 functions, leading to trivial representations. In contrast, the parameter sharing in diffusion models
 318 alleviates this issue significantly. The SWD curve demonstrates an inverse correlation with the test
 319 accuracy curve, indicating that features closer to their optimal state possess stronger representational
 320 capacity. Furthermore, the plot shows that diffusion model features across different noise levels
 321 remain significantly closer to their optimal features at $\sigma_t = 0.06$, while DAE features show less
 322 similarity. These results strongly support our hypothesis.

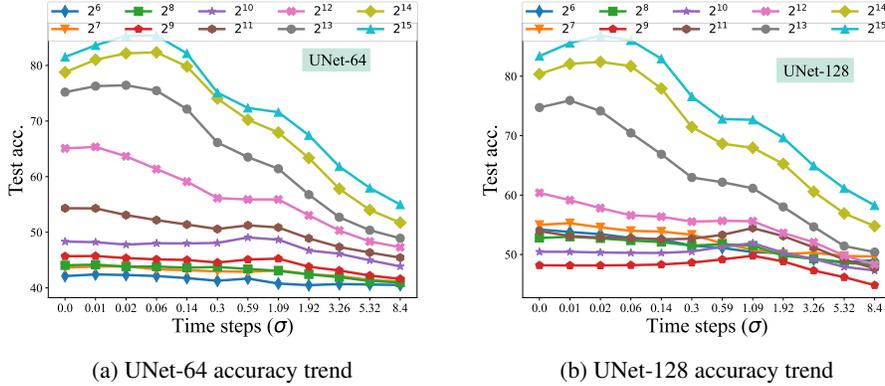


Figure 7: **The influence of data complexity in diffusion-based representation learning.** With the same model trained in Figure 2, we plot the representation learning dynamics for each trained model as a function of changing noise levels.

323 The concept of this “sharing mechanism” is also supported by previous empirical studies on DAEs,
 324 which have shown that sequential training over multiple noise scales enhances representation quality
 325 [Chandra and Sharma, 2014, Geras and Sutton, 2014, Zhang and Zhang, 2018]. In this work, we
 326 conduct an ablation study to explore methods for improving DAE performance at lower noise levels,
 327 finding that training with multiple noise scales provided the most promising results. Further details
 328 can be found in Appendix A.3 (Table 1).

329 4.2 The influence of data complexity in diffusion representation learning

330 So far, our analyses are based on the assumption that the training dataset contains sufficient samples
 331 for the diffusion model to learn the underlying distribution. Interestingly, if this assumption is violated
 332 by training the model on insufficient data, the unimodal representation learning dynamic disappears
 333 and the probing accuracy also drops severely.

334 As illustrated in Figure 7, we train 2 different UNets following the EDM [Karras et al., 2022]
 335 configuration with training dataset size ranging from 2^5 to 2^{15} . The unimodal curve emerges only
 336 when the dataset size exceeds 2^{12} , whereas smaller datasets produce flat curves.

337 The underlying reason for this observation is that, when training data is limited, diffusion models
 338 memorize all individual data points rather than learn the true underlying data structure [Wang et al.,
 339 2024]. In this scenario, the model memorizes an empirical distribution that lacks meaningful low-
 340 dimensional structures and thus deviates from the setting in our theory, leading to the loss of the
 341 unimodal representation dynamic. To confirm this, we calculated the generalization score, which
 342 measures the percentage of generated data that does not belong to the training dataset, as defined in
 343 [Zhang et al., 2023]. As shown in Figure 2, representation learning only achieves strong accuracy
 344 and displays the unimodal dynamic when the generalization score approaches 1, aligning with our
 345 theoretical assumptions.

346 5 Conclusion

347 In this work, we establish a link between distribution recovery, posterior estimation, and representation
 348 learning, providing the first theoretical study of diffusion-based representation learning dynamics
 349 across varying noise scales. Using a low-dimensional mixture of low-rank Gaussians, we show that
 350 the unimodal representation learning dynamic arises from the interplay between data denoising and
 351 class specification. Additionally, our analysis highlights the inherent weight-sharing mechanism
 352 in diffusion models, demonstrating its benefits for peak representation performance as well as its
 353 limitations in optimizing high-noise regions due to increased complexity. Experiments on both
 354 synthetic and real datasets validate our findings.

355 References

356 K. Abstreiter, S. Mittal, S. Bauer, B. Schölkopf, and A. Mehrjou. Diffusion-based representation
 357 learning. *arXiv preprint arXiv:2105.14257*, 2021.

- 358 I. Alkhouri, S. Liang, R. Wang, Q. Qu, and S. Ravishankar. Diffusion-based adversarial purification
359 for robust deep mri reconstruction. *arXiv preprint arXiv:2309.05794*, 2023.
- 360 I. Alkhouri, S. Liang, R. Wang, Q. Qu, and S. Ravishankar. Diffusion-based adversarial purification
361 for robust deep mri reconstruction. In *ICASSP 2024-2024 IEEE International Conference on*
362 *Acoustics, Speech and Signal Processing (ICASSP)*, pages 12841–12845. IEEE, 2024.
- 363 Z. Allen-Zhu and Y. Li. Feature purification: How adversarial training performs robust deep learning.
364 In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages
365 977–988. IEEE, 2022.
- 366 B. D. Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*,
367 12(3):313–326, 1982.
- 368 O. Bar-Tal, H. Chefer, O. Tov, C. Herrmann, R. Paiss, S. Zada, A. Ephrat, J. Hur, Y. Li,
369 T. Michaeli, et al. Lumiere: A space-time diffusion model for video generation. *arXiv preprint*
370 *arXiv:2401.12945*, 2024.
- 371 D. Baranchuk, I. Rubachev, A. Voynov, V. Khruikov, and A. Babenko. Label-efficient semantic
372 segmentation with diffusion models. *arXiv preprint arXiv:2112.03126*, 2021.
- 373 BohaoZou. Denoisingdiffusionprobabilitymodel-ddpm-. [https://github.com/zoubohao/](https://github.com/zoubohao/DenoisingDiffusionProbabilityModel-ddpm-)
374 [DenoisingDiffusionProbabilityModel-ddpm-](https://github.com/zoubohao/DenoisingDiffusionProbabilityModel-ddpm-), 2022.
- 375 B. Chandra and R. K. Sharma. Adaptive noise schedule for denoising autoencoder. In *International*
376 *Conference on Neural Information Processing*, 2014.
- 377 S. Chen, H. Zhang, M. Guo, Y. Lu, P. Wang, and Q. Qu. Exploring low-dimensional subspaces in
378 diffusion models for controllable image editing. *arXiv preprint arXiv:2409.02374*, 2024a.
- 379 X. Chen, Z. Liu, S. Xie, and K. He. Deconstructing denoising diffusion models for self-supervised
380 learning. *arXiv preprint arXiv:2401.14404*, 2024b.
- 381 J. Choi, J. Lee, C. Shin, S. Kim, H. Kim, and S. Yoon. Perception prioritized training of diffusion
382 models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
383 pages 11472–11481, 2022.
- 384 H. Chung, B. Sim, D. Ryu, and J. C. Ye. Improving diffusion models for inverse problems using
385 manifold constraints. *Advances in Neural Information Processing Systems*, 35:25683–25696, 2022.
- 386 A.-D. Doan, B. L. Nguyen, S. Gupta, I. Reid, M. Wagner, and T.-J. Chin. Assessing domain gap for
387 continual domain adaptation in object detection. *Computer Vision and Image Understanding*, 238:
388 103885, 2024.
- 389 M. Fuest, P. Ma, M. Gui, J. S. Fischer, V. T. Hu, and B. Ommer. Diffusion models and representation
390 learning: A survey. *arXiv preprint arXiv:2407.00783*, 2024.
- 391 K. J. Geras and C. Sutton. Scheduled denoising autoencoders. *arXiv preprint arXiv:1406.3269*, 2014.
- 392 S. Gong, V. N. Boddeti, and A. K. Jain. On the intrinsic dimensionality of image representations. In
393 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
394 3987–3996, 2019.
- 395 J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in Neural*
396 *Information Processing Systems*, 33:6840–6851, 2020.
- 397 J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi,
398 D. J. Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv*
399 *preprint arXiv:2210.02303*, 2022.
- 400 D. A. Hudson, D. Zoran, M. Malinowski, A. K. Lampinen, A. Jaegle, J. L. McClelland, L. Matthey,
401 F. Hill, and A. Lerchner. Soda: Bottleneck diffusion models for representation learning. In
402 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
403 23115–23127, 2024.

- 404 A. Hyvärinen and P. Dayan. Estimation of non-normalized statistical models by score matching.
405 *Journal of Machine Learning Research*, 6(4), 2005.
- 406 Z. Kadkhodaie, F. Guth, E. P. Simoncelli, and S. Mallat. Generalization in diffusion models arises
407 from geometry-adaptive harmonic representation. *arXiv preprint arXiv:2310.02557*, 2023.
- 408 T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based generative
409 models. In *Proc. NeurIPS*, 2022.
- 410 D. P. Kingma. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 411 D. P. Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- 412 J. Kong, J. Kim, and J. Bae. Hifi-gan: Generative adversarial networks for efficient and high fidelity
413 speech synthesis. *Advances in neural information processing systems*, 33:17022–17033, 2020.
- 414 Z. Kong, W. Ping, J. Huang, K. Zhao, and B. Catanzaro. DIFFWAVE: A versatile diffusion model for
415 audio synthesis. In *International Conference on Learning Representations*, 2021.
- 416 A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 417 D. Kunin, J. Bloom, A. Goeva, and C. Seed. Loss landscapes of regularized linear autoencoders. In
418 *International conference on machine learning*, pages 3560–3569. PMLR, 2019.
- 419 M. Kwon, J. Jeong, and Y. Uh. Diffusion models already have a semantic latent space. *arXiv preprint*
420 *arXiv:2210.10960*, 2022.
- 421 M. Lewis. Bart: Denoising sequence-to-sequence pre-training for natural language generation,
422 translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- 423 D. Li, H. Ling, A. Kar, D. Acuna, S. W. Kim, K. Kreis, A. Torralba, and S. Fidler. Dreamteacher:
424 Pretraining image backbones with deep generative models. In *Proceedings of the IEEE/CVF*
425 *International Conference on Computer Vision*, pages 16698–16708, 2023.
- 426 X. Li, S. M. Kwon, I. R. Alkhouri, S. Ravishanka, and Q. Qu. Decoupled data consistency with
427 diffusion purification for image restoration. *arXiv preprint arXiv:2403.06054*, 2024.
- 428 S. Mukhopadhyay, M. Gwilliam, V. Agarwal, N. Padmanabhan, A. Swaminathan, S. Hegde,
429 T. Zhou, and A. Shrivastava. Diffusion models beat gans on image classification. *arXiv preprint*
430 *arXiv:2307.08702*, 2023.
- 431 M.-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes.
432 *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729,
433 2008.
- 434 Y.-H. Park, M. Kwon, J. Jo, and Y. Uh. Unsupervised discovery of semantic latent directions in
435 diffusion models. *arXiv preprint arXiv:2302.12469*, 2023.
- 436 W. Peebles and S. Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF*
437 *International Conference on Computer Vision*, pages 4195–4205, 2023.
- 438 D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach.
439 Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint*
440 *arXiv:2307.01952*, 2023.
- 441 P. Pope, C. Zhu, A. Abdelkader, M. Goldblum, and T. Goldstein. The intrinsic dimension of images
442 and its impact on learning. *arXiv preprint arXiv:2104.08894*, 2021.
- 443 K. Preechakul, N. Chatthee, S. Wizadwongsa, and S. Suwajanakorn. Diffusion autoencoders: Toward
444 a meaningful and decodable representation. In *Proceedings of the IEEE/CVF conference on*
445 *computer vision and pattern recognition*, pages 10619–10629, 2022.
- 446 A. Pretorius, S. Kroon, and H. Kamper. Learning dynamics of linear denoising autoencoders. In
447 *International Conference on Machine Learning*, pages 4141–4150. PMLR, 2018.

- 448 D. Roich, R. Mokady, A. H. Bermanto, and D. Cohen-Or. Pivotal tuning for latent-based editing of
449 real images. *ACM Transactions on Graphics (TOG)*, 42(1):1–13, 2022.
- 450 R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis
451 with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
452 and Pattern Recognition*, pages 10684–10695, 2022.
- 453 O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image
454 segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015:
455 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18,*
456 pages 234–241. Springer, 2015.
- 457 N. Ruiz, Y. Li, V. Jampani, Y. Pritch, M. Rubinstein, and K. Aberman. Dreambooth: Fine tuning
458 text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF
459 conference on computer vision and pattern recognition*, pages 22500–22510, 2023.
- 460 Y. Shi, C. Xue, J. H. Liew, J. Pan, H. Yan, W. Zhang, V. Y. Tan, and S. Bai. Dragdiffusion: Harnessing
461 diffusion models for interactive point-based image editing. In *Proceedings of the IEEE/CVF
462 Conference on Computer Vision and Pattern Recognition*, pages 8839–8849, 2024.
- 463 J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning
464 using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages
465 2256–2265. PMLR, 2015.
- 466 B. Song, S. M. Kwon, Z. Zhang, X. Hu, Q. Qu, and L. Shen. Solving inverse problems with latent
467 diffusion models via hard data consistency. In *The Twelfth International Conference on Learning
468 Representations*, 2024.
- 469 Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based gener-
470 ative modeling through stochastic differential equations. *International Conference on Learning
471 Representations*, 2021.
- 472 J. Stanczuk, G. Batzolis, T. Deveney, and C.-B. Schönlieb. Your diffusion model secretly knows the
473 dimension of the data manifold. *arXiv preprint arXiv:2212.12611*, 2022.
- 474 H. Steck. Autoencoders that don’t overfit towards the identity. In *Neural Information Processing
475 Systems*, 2020.
- 476 tanelp. tiny-diffusion. <https://github.com/tanelp/tiny-diffusion>, 2022.
- 477 L. Tang, M. Jia, Q. Wang, C. P. Phoo, and B. Hariharan. Emergent correspondence from image
478 diffusion. *Advances in Neural Information Processing Systems*, 36:1363–1389, 2023.
- 479 P. Vincent. A connection between score matching and denoising autoencoders. *Neural computation*,
480 23(7):1661–1674, 2011.
- 481 P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features
482 with denoising autoencoders. In *International Conference on Machine Learning*, 2008.
- 483 P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol. Stacked denoising autoencoders:
484 Learning useful representations in a deep network with a local denoising criterion. *J. Mach. Learn.
485 Res.*, 11:3371–3408, 2010. URL <https://api.semanticscholar.org/CorpusID:17804904>.
- 486 B. Wang and J. J. Vastola. Diffusion models generate images like painters: an analytical theory of
487 outline first, details later. *arXiv preprint arXiv:2303.02490*, 2023.
- 488 P. Wang, H. Zhang, Z. Zhang, S. Chen, Y. Ma, and Q. Qu. Diffusion models learn low-dimensional
489 distributions via subspace clustering. *arXiv preprint arXiv:2409.02426*, 2024.
- 490 Y. Wang, Y. Schiff, A. Gokaslan, W. Pan, F. Wang, C. De Sa, and V. Kuleshov. Infodiffusion: Repre-
491 sentation learning using information maximizing diffusion models. In *International Conference on
492 Machine Learning*, pages 36336–36354. PMLR, 2023.

- 493 W. Xiang, H. Yang, D. Huang, and Y. Wang. Denoising diffusion autoencoders are unified self-
494 supervised learners. In *Proceedings of the IEEE/CVF International Conference on Computer*
495 *Vision*, pages 15802–15812, 2023.
- 496 X. Yang and X. Wang. Diffusion model as representation learner. In *Proceedings of the IEEE/CVF*
497 *International Conference on Computer Vision*, pages 18938–18949, 2023.
- 498 H. Zhang, J. Zhou, Y. Lu, M. Guo, P. Wang, L. Shen, and Q. Qu. The emergence of reproducibility
499 and consistency in diffusion models. In *Forty-first International Conference on Machine Learning*,
500 2023.
- 501 H. Zhang, Y. Lu, I. Alkhouri, S. Ravishankar, D. Song, and Q. Qu. Improving training efficiency of
502 diffusion models via multi-stage framework and tailored multi-decoder architectures. In *Conference*
503 *on Computer Vision and Pattern Recognition 2024*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=YtptmpZQ0g)
504 [forum?id=YtptmpZQ0g](https://openreview.net/forum?id=YtptmpZQ0g).
- 505 Q. Zhang and L. Zhang. Convolutional adaptive denoising autoencoders for hierarchical feature
506 extraction. *Frontiers of Computer Science*, 12:1140 – 1148, 2018.

507 A Appendix / supplemental material

508 The Appendix is organized as follows: in Appendix A.1, we discuss related works; in Appendix A.2,
509 we present the detailed experimental setups for the empirical results in the paper; in Appendix A.3, we
510 provide complementary experiments. Lastly, in Appendix A.4, we provide proof details for Section 3.

511 A.1 Related Works

512 **Denoising auto-encoders.** Denoising autoencoders (DAEs) are trained to reconstruct corrupted
513 images to extract semantically meaningful information, which can be applied to various vision
514 [Vincent et al., 2008, 2010] and language downstream tasks [Lewis, 2019]. Related to our analysis
515 of the weight-sharing mechanism, several studies have shown that training with a noise scheduler
516 can enhance downstream performance [Chandra and Sharma, 2014, Geras and Sutton, 2014, Zhang
517 and Zhang, 2018]. On the theoretical side, prior works have studied the learning dynamics [Pretorius
518 et al., 2018, Steck, 2020] and optimization landscape [Kunin et al., 2019] through the simplified
519 linear DAE models.

520 **Diffusion-based representation learning.** Diffusion-based representation learning Fuest et al.
521 [2024] has demonstrated significant success in various downstream tasks, including image classi-
522 fication [Xiang et al., 2023, Mukhopadhyay et al., 2023], segmentation [Baranchuk et al., 2021],
523 correspondence [Tang et al., 2023], and image editing [Shi et al., 2024]. To further enhance the utility
524 of diffusion features, knowledge distillation methods [Yang and Wang, 2023, Li et al., 2023] have
525 been proposed, aiming to bypass the computationally expensive grid search for the optimal t in feature
526 extraction and improving downstream performance. Beyond directly using intermediate features from
527 pre-trained diffusion models, research efforts has also explored novel loss functions [Abstreiter et al.,
528 2021, Wang et al., 2023] and network modifications [Hudson et al., 2024, Preechakul et al., 2022] to
529 develop more unified generative and representation learning capabilities within diffusion models. Un-
530 like the aforementioned efforts, our work focuses more on understanding the representation learning
531 capabilities of diffusion models.

532 A.2 Experimental Details

533 In this section, we provide technical details for all the experiments in the main body of the paper.

534 **Experimental details for Figure 1 (a)-(b).** We utilize a minimal implementation of the original
535 DDPM model from an online public repository [BohaoZou, 2022], consisting of a 12-layer UNet
536 (including input/output embedding layers), and train it on the CIFAR10 dataset with $T = 1000$ time
537 steps for 200 epochs with an AdamW optimizer and learning rate 1×10^{-4} . Features are extracted as
538 512-dimensional vectors from the output of the 7th layer (i.e., the bottleneck layer) at time steps [1, 5,
539 10, 20, 30, 40, 60, 80, 100, 200, 400, 500, 600], each corresponding to a specific σ_t ranging from
540 0.01 to 6.17. Linear probing is applied to the extracted features, as in [Xiang et al., 2023], to plot
541 the feature probing accuracy curve in Figure 1(a). For the posterior estimation ($x_\theta(x_0, t)$) probing
542 accuracy curve, also shown in Figure 1(a), we use a two-layer MLP probe with ReLU activation. The
543 estimated posterior at these time steps is visualized in Figure 1(b).

544 **Experimental details for Figure 1 (c)-(d).** We train diffusion models based on the unified frame-
545 work proposed by Karras et al. [2022]. Specifically, we use the DDPM+ network, and use EDM
546 configuration for Figure 1 (c) while taking VP configuration Figure 1 (d). Karras et al. [2022]
547 has shown equivalence between VP configuration and the traditional DDPM setting, thus we call
548 the models in Figure 1 (d) as DDPM* models. For each of EDM and VP configuration, we train
549 two models on CIFAR10 and CIFAR100, respectively. After training, we conduct linear probe on
550 CIFAR10 and CIFAR100. At a specific noise level $\sigma(t)$, we either use clean image x_0 or noisy
551 image $x_t = x_0 + n$ as input to the EDM or the DDPM* models for extracting features after the
552 '8x8_block3' layer. Here, n represents random noise and $n \sim \mathcal{N}(\mathbf{0}, \sigma(t)^2 \mathbf{I})$. We train a logistic
553 regression on features in the train split and report the classification accuracy on the test split of the
554 dataset. We perform the linear probe for each of the following noise levels: [0.002, 0.008, 0.023,
555 0.060, 0.140, 0.296, 0.585, 1.088, 1.923, 3.257].

556 **Experimental details for Figure 3 and Figure 4.** For the MoLRG experiments, we train a 3-layer
 557 MLP with ReLU activation and a hidden dimension of 128, following the setup provided in an
 558 open-source repository [tanelp, 2022]. The MLP is trained for 200 epochs using DDPM scheduling
 559 with $T = 500$, employing the Adam optimizer with a learning rate of 1×10^{-3} . For feature extraction,
 560 we use the activations of the second layer of the MLP (dimension 128) as intermediate features for
 561 linear probing. For CSNR computation, we follow the definition in Equation (8) since we have access
 562 to the ground-truth basis for the MoLRG data. In Figure 3, we visualize the posterior estimations
 563 at time steps [1, 20, 80, 200, 260] by projecting them onto the union of U_1, U_2 , and U_3 (a 3D
 564 space), then further projecting onto the 2D plane along the (1, 1, 1) direction. The subtitles of each
 565 visualization show the corresponding probing accuracy and CSNR calculated as explained above. For
 566 Figure 4(a)(b), we plot the accuracy and CSNR at time steps [1, 5, 10, 20, 40, 60, 80, 100, 120, 140,
 567 160, 180, 220, 240, 260]. We perform linear probing using the features extracted from the training set
 568 and test on five different MoLRG datasets generated with five different random seeds, reporting the
 569 average accuracy.

570 **Experimental details for Figure 5.** We use the same experimental settings as in Figure 1(a)(b).
 571 Additionally, we train individual DAEs for each different time step. The accuracy curves in Figure 5(a)
 572 are plotted identically as in Figure 1(a). The CSNR metric in Figure 5(b) is calculated from the
 573 definition Equation (8), with the basis U_k for each CIFAR10 class estimated as the first five right
 574 singular vectors of the data from the k -th class.

575 **Experimental details for Figure 6.** We train individual DAEs using the DDPM++ net-
 576 work and VP configuration outlined in Karras et al. [2022] at the following noise scales:
 577 [0.002, 0.008, 0.023, 0.06, 0.14, 0.296, 0.585, 1.088, 1.923, 3.257]. Each model is trained for 500
 578 epochs using the Adam optimizer [Kingma, 2014] with a fixed learning rate of 1×10^{-4} . For the
 579 diffusion models, we reuse the model from Figure 1(d). The sliced Wasserstein distance is computed
 580 according to the implementation described in Doan et al. [2024].

581 **Experimental details for Figure 7.** We use the DDPM++ network and VP configuration to train
 582 diffusion models[Karras et al., 2022] on the CIFAR10 dataset, using two network configurations:
 583 UNet-64 and UNet-128, by varying the embedding dimension of the UNet. Training dataset sizes
 584 range exponentially from 2^6 to 2^{15} . For each dataset size, both UNet-64 and UNet-128 are trained on
 585 the same subset of the training data. All models are trained with a duration of 50K images following
 586 the EDM training setup. After training, we calculate the generalization score as described in Zhang
 587 et al. [2023], using 10K generated images and the full training subset to compute the score.

588 A.3 Additional Experiments

589 **Additional representation learning experiments on DDPM.** Apart from EDM and DDPM*
 590 models pre-trained using the framework proposed by Karras et al. [2022], we also experiment with
 591 the features extracted by classic DDPM models [Ho et al., 2020] to make sure the observations do not
 592 depend on the specific training framework. We use the same groups of noise levels and also test using
 593 clean or noisy images as input to extract features at the bottleneck layer, and then conduct the linear
 594 probe. The DDPM models we use are trained on the Flowers-102 [Nilsback and Zisserman, 2008]
 595 and the CIFAR10 dataset accordingly. Different from the framework proposed by Karras et al. [2022],
 596 the input to the classic DDPM model is the same as the input to the UNet inside. Therefore, we
 597 calculate the scaling factor $\sqrt{\bar{\alpha}_t} = 1/\sqrt{\sigma^2(t) + 1}$, and use $\sqrt{\bar{\alpha}_t}\mathbf{x}_0$ as the clean image input. Besides,
 598 for noisy input, we set $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}(\mathbf{x}_0 + \mathbf{n})$, with $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma(t)^2\mathbf{I})$. The linear probe results are
 599 presented in Figure 8, where we consistently see an unimodal curve, as well as compatible or even
 600 superior representation learning performance of clean input \mathbf{x}_0 .

601 **Validation of f^* approximation in Section 3.** In Section 3, we approximate the optimal posterior
 602 estimation function \mathbf{x}_θ using f^* by taking the expectation inside the softmax with respect to \mathbf{x}_0 . To
 603 validate this approximation, we compare the CSNR calculated from \mathbf{x}_θ and from f^* using (8) and
 604 (9), respectively. We use a fixed dataset size of 2400 and set the default parameters to $n = 50$, $d = 5$,
 605 $K = 3$, and $\delta = 0.1$ to generate MoLRG data. We then vary one parameter at a time while keeping
 606 the others constant, and present the computed CSNR in Figure 9. As shown, the approximated
 607 CSNR score consistently aligns with the actual score.

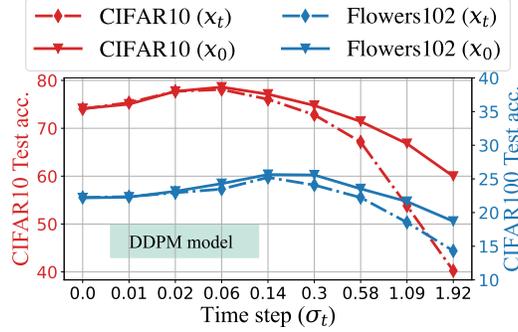


Figure 8: **Performance comparison: clean vs. noisy inputs.** We use pre-trained DDPM model on the Flowers-102 [Nilsback and Zisserman, 2008] and CIFAR10 dataset. The feature probing accuracy is plotted to compare the performance when using clean versus noisy inputs.

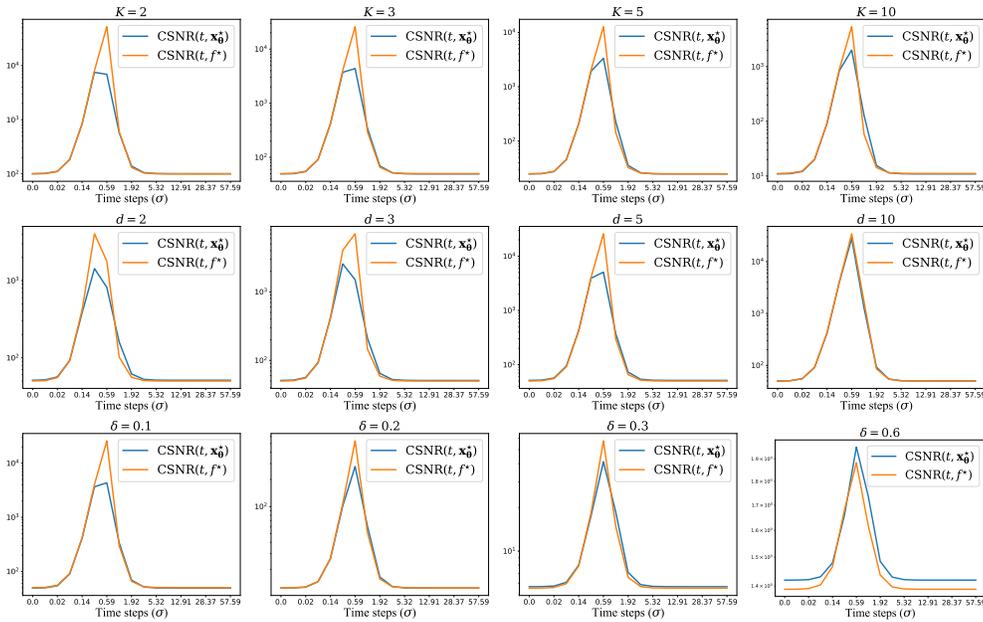


Figure 9: **Comparison between CSNR calculated using the optimal model x_θ^* and the CSNR calculated with our approximation in Theorem 1.** We generate MoLRG data and calculate CSNR using both the corresponding optimal posterior function x_θ^* and our approximation f^* from Theorem 1. Default parameters are set as $n = 50$, $d = 5$, $K = 3$, and $\delta = 0.1$. In each row, we vary one parameter while keeping the others fixed, comparing the actual and approximated CSNR.

608 **Mitigating the performance gap between DAE and diffusion models.** Throughout the empirical
609 results presented in this paper, we consistently observe a performance gap between individual DAEs
610 and diffusion models, especially in low-noise regions. Here, we use a DAE trained on the CIFAR-10
611 dataset with a single noise level $\sigma = 0.002$, using the NCSN++ architecture [Karras et al., 2022].
612 In the default setting, the DAE achieves a test accuracy of 32.3. We then explore three methods to
613 improve the test performance: (a) adding dropout, as noise regularization and dropout have been
614 effective in preventing autoencoders from learning identity functions [Steck, 2020]; (b) adopting
615 EDM-based preconditioning during training, including input/output scaling, loss weighting, etc.;
616 and (c) multi-level noise training, in which the DAE is trained simultaneously on three noise levels
617 $[0.002, 0.012, 0.102]$. Each modification is applied independently, and the results are reported in
618 Table 1. As shown, dropout helps improve performance, but even with a dropout rate of 0.95, the
619 improvement is minor. EDM-based preconditioning achieves moderate improvement, while multi-
620 level noise training yields the most promising results, demonstrating the benefit of incorporating the
621 diffusion process in DAE training.

Table 1: **Improve DAE representation performance at low noise region.** A vanilla DAE trained on the CIFAR-10 dataset with a single noise level of $\sigma = 0.002$ serves as the baseline. We evaluate the performance improvement of dropout regularization, EDM-based preconditioning, and multi-level noise training ($\sigma = \{0.002, 0.012, 0.102\}$). Each technique is applied independently to assess its contribution to performance enhancement.

Modifications	Test acc.
Vanilla DAE	32.3
+Dropout (0.5)	35.3
+Dropout (0.9)	36.4
+Dropout (0.95)	38.1
+EDM preconditioning	49.2
+Multi-level noise training	58.6

622 A.4 Proofs

623 A.4.1 Proof of Proposition 1

624 *Proof.* We follow the same proof steps as in [Wang et al., 2024] Lemma 1 with a change of variable.

625 Let $\mathbf{c}_k = \begin{bmatrix} \mathbf{a}_k \\ \mathbf{e}_k \end{bmatrix}$ and $\widetilde{\mathbf{U}}_k = [\mathbf{U}_k \quad \delta \mathbf{U}_k^\perp]$, we first compute

$$\begin{aligned}
& p_t(\mathbf{x}|Y = k) \\
&= \int p_t(\mathbf{x}|Y = k, \mathbf{c}_k) \mathcal{N}(\mathbf{c}_k; \mathbf{0}, \mathbf{I}_{d+D}) d\mathbf{c}_k \\
&= \int p_t(\mathbf{x}|\mathbf{x}_0 = \widetilde{\mathbf{U}}_k \mathbf{c}_k) \mathcal{N}(\mathbf{c}_k; \mathbf{0}, \mathbf{I}_{d+D}) d\mathbf{c}_k \\
&= \int \mathcal{N}(\mathbf{x}; s_t \widetilde{\mathbf{U}}_k \mathbf{c}_k, \gamma_t^2 \mathbf{I}_n) \mathcal{N}(\mathbf{c}_k; \mathbf{0}, \mathbf{I}_{d+D}) d\mathbf{c}_k \\
&= \frac{1}{(2\pi)^{n/2} (2\pi)^{(d+D)/2} \gamma_t^n} \int \exp\left(-\frac{1}{2\gamma_t^2} \|\mathbf{x} - s_t \widetilde{\mathbf{U}}_k \mathbf{c}_k\|^2\right) \exp\left(-\frac{1}{2} \|\mathbf{c}_k\|^2\right) d\mathbf{c}_k \\
&= \frac{1}{(2\pi)^{n/2} (2\pi)^{(d+D)/2} \gamma_t^n} \int \exp\left(-\frac{1}{2\gamma_t^2} \left(\mathbf{x}^T \mathbf{x} - 2s_t \mathbf{x}^T \widetilde{\mathbf{U}}_k \mathbf{c}_k + s_t^2 \mathbf{c}_k^T \widetilde{\mathbf{U}}_k^T \widetilde{\mathbf{U}}_k \mathbf{c}_k + \gamma_t^2 \mathbf{c}_k^T \mathbf{c}_k\right)\right) d\mathbf{c}_k \\
&= \frac{1}{(2\pi)^{n/2} \gamma_t^n} \left(\frac{s_t^2 + \gamma_t^2}{\gamma_t^2}\right)^{-d/2} \left(\frac{s_t^2 \delta^2 + \gamma_t^2}{\gamma_t^2}\right)^{-D/2} \exp\left(-\frac{1}{2\gamma_t^2} \mathbf{x}^T \left(\mathbf{I}_n - \frac{s_t^2}{s_t^2 + \gamma_t^2} \mathbf{U}_k \mathbf{U}_k^T - \frac{s_t^2 \delta^2}{s_t^2 \delta^2 + \gamma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T}\right) \mathbf{x}\right) \\
&\int \frac{1}{(2\pi)^{d/2}} \left(\frac{\gamma_t^2}{s_t^2 + \gamma_t^2}\right)^{-d/2} \exp\left(-\frac{s_t^2 + \gamma_t^2}{2\gamma_t^2} \left\|\mathbf{a}_k - \frac{s_t}{s_t^2 + \gamma_t^2} \mathbf{U}_k^T \mathbf{x}\right\|^2\right) d\mathbf{a}_k \\
&\int \frac{1}{(2\pi)^{D/2}} \left(\frac{\gamma_t^2}{s_t^2 \delta^2 + \gamma_t^2}\right)^{-D/2} \exp\left(-\frac{s_t^2 \delta^2 + \gamma_t^2}{2\gamma_t^2} \left\|\mathbf{e}_k - \frac{s_t \delta}{s_t^2 \delta^2 + \gamma_t^2} \mathbf{U}_k^{\perp T} \mathbf{x}\right\|^2\right) d\mathbf{e}_k \\
&= \frac{1}{(2\pi)^{n/2}} \frac{1}{(s_t^2 + \gamma_t^2)^{d/2} (s_t^2 \delta^2 + \gamma_t^2)^{D/2}} \exp\left(-\frac{1}{2\gamma_t^2} \mathbf{x}^T \left(\mathbf{I}_n - \frac{s_t^2}{s_t^2 + \gamma_t^2} \mathbf{U}_k \mathbf{U}_k^T - \frac{s_t^2 \delta^2}{s_t^2 \delta^2 + \gamma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T}\right) \mathbf{x}\right) \\
&= \frac{1}{(2\pi)^{n/2} \det^{1/2}(s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)} \\
&\exp\left(-\frac{1}{2} \mathbf{x}^T (s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)^{-1} \mathbf{x}\right) \\
&= \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n),
\end{aligned}$$

626 where we repeatedly apply the pdf of multi-variate Gaussian and the second last equality uses
627 $\det(s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n) = (s_t^2 + \gamma_t^2)^d (s_t^2 \delta^2 + \gamma_t^2)^D$ and $(s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} +$
628 $\gamma_t^2 \mathbf{I}_n)^{-1} = (\mathbf{I}_n - s_t^2 / (s_t^2 + \gamma_t^2) \mathbf{U}_k \mathbf{U}_k^T - s_t^2 \delta^2 / (s_t^2 \delta^2 + \gamma_t^2) \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T}) / \gamma_t^2$ because of the Wood-

629 bury matrix inversion lemma. Hence, with $\mathbb{P}(Y = k) = \pi_k$ for each $k \in [K]$, we have

$$p_t(\mathbf{x}) = \sum_{k=1}^K p_t(\mathbf{x}|Y = k)\mathbb{P}(Y = k) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n).$$

630 Now we can compute the score function

$$\begin{aligned} \nabla \log p_t(\mathbf{x}) &= \frac{\nabla p_t(\mathbf{x})}{p_t(\mathbf{x})} = \frac{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)} \\ &= -\frac{1}{\gamma_t^2} \left(\mathbf{x} - \frac{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)} \right). \end{aligned}$$

631 According to Tweedie's formula, we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] &= \frac{\mathbf{x}_t + \gamma_t^2 \nabla \log p_t(\mathbf{x}_t)}{s_t} \\ &= \frac{s_t}{s_t^2 + \gamma_t^2} \frac{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n) \mathbf{U}_k \mathbf{U}_k^T \mathbf{x}}{\mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)} \\ &\quad + \frac{s_t \delta^2}{s_t^2 \delta^2 + \gamma_t^2} \frac{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n) \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \mathbf{x}}{\mathcal{N}(\mathbf{x}; \mathbf{0}, s_t^2 \mathbf{U}_k \mathbf{U}_k^T + s_t^2 \delta^2 \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} + \gamma_t^2 \mathbf{I}_n)} \\ &= \frac{s_t}{s_t^2 + \gamma_t^2} \frac{\sum_{k=1}^K \pi_k \exp(\phi_t \|\mathbf{U}_k^T \mathbf{x}_t\|^2) \exp(\psi_t \|\mathbf{U}_k^{\perp T} \mathbf{x}_t\|^2) \mathbf{U}_k \mathbf{U}_k^T \mathbf{x}_t}{\sum_{k=1}^K \pi_k \exp(\phi_t \|\mathbf{U}_k^T \mathbf{x}_t\|^2) \exp(\psi_t \|\mathbf{U}_k^{\perp T} \mathbf{x}_t\|^2)} \\ &\quad + \frac{s_t \delta^2}{s_t^2 \delta^2 + \gamma_t^2} \frac{\sum_{k=1}^K \pi_k \exp(\phi_t \|\mathbf{U}_k^T \mathbf{x}_t\|^2) \exp(\psi_t \|\mathbf{U}_k^{\perp T} \mathbf{x}_t\|^2) \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \mathbf{x}_t}{\sum_{k=1}^K \pi_k \exp(\phi_t \|\mathbf{U}_k^T \mathbf{x}_t\|^2) \exp(\psi_t \|\mathbf{U}_k^{\perp T} \mathbf{x}_t\|^2)}, \end{aligned}$$

632 with $\phi_t = s_t^2 / (2\gamma_t^2 (s_t^2 + \gamma_t^2))$ and $\psi_t = s_t^2 \delta^2 / (2\gamma_t^2 (s_t^2 \delta^2 + \gamma_t^2))$. The final equality uses the pdf of
633 multi-variant Gaussian and the matrix inversion lemma discussed earlier.

634 Now since π_k is consistent for all k and $s_t = 1$, we have

$$\begin{aligned} \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] &= \sum_{k=1}^K w_k^*(\mathbf{x}_t) \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \right) \mathbf{x}_t \\ \text{where } w_k^*(\mathbf{x}_t) &:= \frac{\exp\left(\frac{1}{2\sigma_t^2(1+\sigma_t^2)} \|\mathbf{U}_k^T \mathbf{x}_t\|^2 + \frac{\delta^2}{2\sigma_t^2(\delta^2+\sigma_t^2)} \|\mathbf{U}_k^{\perp T} \mathbf{x}_t\|^2\right)}{\sum_{k=1}^K \exp\left(\frac{1}{2\sigma_t^2(1+\sigma_t^2)} \|\mathbf{U}_k^T \mathbf{x}_t\|^2 + \frac{\delta^2}{2\sigma_t^2(\delta^2+\sigma_t^2)} \|\mathbf{U}_k^{\perp T} \mathbf{x}_t\|^2\right)}. \end{aligned}$$

635

□

636 **A.4.2 Proof of Theorem 1**

637 *Proof.* Following Equation (8) and Lemma 1, we can write

$$\begin{aligned}
\text{CSNR}(t, f^*) &= \frac{\mathbb{E}_{\mathbf{x}_0}[\|\mathbf{U}_k \mathbf{U}_k^T f^*(\mathbf{x}_0, t)\|^2]}{\mathbb{E}_{\mathbf{x}_0}[\sum_{l \neq k} \|\mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}_0, t)\|^2]} = \frac{\mathbb{E}_{\mathbf{x}_0}[\|\mathbf{U}_k \mathbf{U}_k^T f^*(\mathbf{x}_0, t)\|^2]}{\sum_{l \neq k} \mathbb{E}_{\mathbf{x}_0}[\|\mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}_0, t)\|^2]} \\
&= \frac{\left(\frac{\hat{w}_k}{1+\sigma_t^2} + \frac{(K-1)\delta^2 \hat{w}_l}{\delta^2 + \sigma_t^2}\right)^2 d}{(K-1) \left(\frac{\hat{w}_l}{1+\sigma_t^2} + \frac{\delta^2(\hat{w}_k + (K-2)\hat{w}_l)}{\delta^2 + \sigma_t^2}\right)^2 \delta^2 d} \\
&= \frac{1}{(K-1)\delta^2} \cdot \left(\frac{\hat{w}_k \delta^2 + \hat{w}_k \sigma_t^2 + (K-1)\delta^2 \hat{w}_l + (K-1)\delta^2 \hat{w}_l \sigma_t^2}{\hat{w}_l \delta^2 + \hat{w}_l \sigma_t^2 + \delta^2 \hat{w}_k + (K-2)\delta^2 \hat{w}_l + \delta^2 \hat{w}_k \sigma_t^2 + (K-2)\delta^2 \hat{w}_l \sigma_t^2}\right)^2 \\
&= \frac{1}{(K-1)\delta^2} \cdot \left(\frac{\delta^2 + \sigma_t^2 (\hat{w}_k + (K-1)\delta^2 \hat{w}_l)}{\delta^2 + \sigma_t^2 (\hat{w}_l + \delta^2 \hat{w}_k + (K-2)\delta^2 \hat{w}_l)}\right)^2 \\
&= \frac{1}{(K-1)\delta^2} \cdot \left(\frac{1 + \frac{\sigma_t^2}{\delta^2} ((1-\delta^2)\hat{w}_k + \delta^2(\hat{w}_k + (K-1)\hat{w}_l))}{1 + \frac{\sigma_t^2}{\delta^2} ((1-\delta^2)\hat{w}_l + \delta^2(\hat{w}_l + \hat{w}_k + (K-2)\hat{w}_l))}\right)^2 \\
&= \frac{1}{(K-1)\delta^2} \cdot \left(\frac{1 + \frac{\sigma_t^2}{\delta^2} ((1-\delta^2)\hat{w}_k + \delta^2)}{1 + \frac{\sigma_t^2}{\delta^2} ((1-\delta^2)\hat{w}_l + \delta^2)}\right)^2 \\
&= \frac{1}{(K-1)\delta^2} \cdot \left(\frac{1 + \frac{\sigma_t^2}{\delta^2} h(\hat{w}_k, \delta)}{1 + \frac{\sigma_t^2}{\delta^2} h(\hat{w}_l, \delta)}\right)^2
\end{aligned}$$

638 where $h(w, \delta) := (1 - \delta^2)w + \delta^2$. □

639 **Lemma 1.** *With the set up of a K -class MoLRG data distribution as defined in (4), consider the*
640 *following the function:*

$$f^*(\mathbf{x}, t) = \sum_{k=1}^K \hat{w}_k(\mathbf{x}) \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \right) \mathbf{x}, \quad (11)$$

$$\text{where } \hat{w}_k(\mathbf{x}) := \frac{\exp(\mathbb{E}_{\mathbf{x}}[g_k(\mathbf{x}, t)])}{\sum_{k=1}^K \exp(\mathbb{E}_{\mathbf{x}}[g_k(\mathbf{x}, t)]), \quad (12)$$

$$\text{and } g_k(\mathbf{x}) = \frac{1}{2\sigma_t^2(1 + \sigma_t^2)} \|\mathbf{U}_k^T \mathbf{x}\|^2 + \frac{\delta^2}{2\sigma_t^2(\delta^2 + \sigma_t^2)} \|\mathbf{U}_k^{\perp T} \mathbf{x}\|^2. \quad (13)$$

641 *I.e., we consider a simplified version of the expected posterior mean as in (5) by taking expectation*
642 *of $g_k(\mathbf{x})$ prior to the softmax operation. Under this setting, for any clean \mathbf{x}_0 from class k (i.e.,*
643 *$\mathbf{x}_0 = \mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i$), we have:*

$$\mathbb{E}_{\mathbf{x}_0}[\|\mathbf{U}_k \mathbf{U}_k^T f^*(\mathbf{x}_0, t)\|^2] = \left(\frac{\hat{w}_k}{1 + \sigma_t^2} + \frac{(K-1)\delta^2 \hat{w}_l}{\delta^2 + \sigma_t^2}\right)^2 d \quad (14)$$

$$\mathbb{E}_{\mathbf{x}_0}[\|\mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}_0, t)\|^2] = \left(\frac{\hat{w}_l}{1 + \sigma_t^2} + \frac{\delta^2(\hat{w}_k + (K-2)\hat{w}_l)}{\delta^2 + \sigma_t^2}\right)^2 \delta^2 d \quad (15)$$

$$\mathbb{E}_{\mathbf{x}_0}[\|\mathbf{U}_\perp \mathbf{U}_\perp^T f^*(\mathbf{x}_0, t)\|^2] = \frac{\delta^6(n - kd)}{(\delta^2 + \sigma_t^2)^2} \quad (16)$$

$$\begin{aligned}
\mathbb{E}[\|f^*(\mathbf{x}_0, t)\|^2] &= \underbrace{\left(\frac{\hat{w}_k}{1 + \sigma_t^2} + \frac{(K-1)\delta^2\hat{w}_l}{\delta^2 + \sigma_t^2}\right)^2}_d \mathbb{E}[\|\mathbf{U}_k \mathbf{U}_k^T f^*(\mathbf{x}_0, t)\|^2] \\
&+ (K-1) \underbrace{\left(\frac{\hat{w}_l}{1 + \sigma_t^2} + \frac{\delta^2(\hat{w}_k + (K-2)\hat{w}_l)}{\delta^2 + \sigma_t^2}\right)^2}_{\mathbb{E}[\|\sum_{l \neq k}^K \mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}_0, t)\|^2]} \delta^2 d + \underbrace{\frac{\delta^6(n-Kd)}{(\delta^2 + \sigma_t^2)^2}}_{\mathbb{E}[\|\mathbf{U}_\perp \mathbf{U}_\perp^T f^*(\mathbf{x}_0, t)\|^2]}
\end{aligned} \tag{17}$$

644 and

$$\begin{aligned}
\hat{w}_k &:= \hat{w}_k(\mathbf{x}_0) = \frac{\exp\left(\frac{d}{2\sigma_t^2(1+\sigma_t^2)} + \frac{\delta^4 D}{2\sigma_t^2(\delta^2 + \sigma_t^2)}\right)}{\exp\left(\frac{d}{2\sigma_t^2(1+\sigma_t^2)} + \frac{\delta^4 D}{2\sigma_t^2(\delta^2 + \sigma_t^2)}\right) + (K-1) \exp\left(\frac{\delta^2 d}{2\sigma_t^2(1+\sigma_t^2)} + \frac{\delta^2 d + \delta^4(D-d)}{2\sigma_t^2(\delta^2 + \sigma_t^2)}\right)}, \\
\hat{w}_l &:= \hat{w}_l(\mathbf{x}_0) = \frac{\exp\left(\frac{\delta^2 d}{2\sigma_t^2(1+\sigma_t^2)} + \frac{\delta^2 d + \delta^4(D-d)}{2\sigma_t^2(\delta^2 + \sigma_t^2)}\right)}{\exp\left(\frac{d}{2\sigma_t^2(1+\sigma_t^2)} + \frac{\delta^4 D}{2\sigma_t^2(\delta^2 + \sigma_t^2)}\right) + (K-1) \exp\left(\frac{\delta^2 d}{2\sigma_t^2(1+\sigma_t^2)} + \frac{\delta^2 d + \delta^4(D-d)}{2\sigma_t^2(\delta^2 + \sigma_t^2)}\right)}
\end{aligned} \tag{18}$$

645 for all class index $l \neq k$.

646 *Proof.* Throughout the proof, we use the following notation for slices of vectors.

$\mathbf{e}_i[a : b]$ Slices of vector \mathbf{e}_i from a th entry to b th entry.

647 We begin with the softmax terms. Since each class has its unique disjoint subspace, it suffices to
648 consider $g_k(\mathbf{x}_0, t)$ and $g_l(\mathbf{x}_0, t)$ for any $l \neq k$. Let $a_t = \frac{1}{2\sigma_t^2(1+\sigma_t^2)}$ and $c_t = \frac{\delta^2}{2\sigma_t^2(\delta^2 + \sigma_t^2)}$, we have:

$$\begin{aligned}
\mathbb{E}[g_k(\mathbf{x}_0, t)] &= \mathbb{E}[a_t \|\mathbf{U}_k^T \mathbf{x}_0\|^2 + c_t \|\mathbf{U}_k^{\perp T} \mathbf{x}_0\|^2] \\
&= \mathbb{E}[a_t \|\mathbf{U}_k^T (\mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i)\|^2] + \mathbb{E}[c_t \|\mathbf{U}_k^{\perp T} (\mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i)\|^2] \\
&= \mathbb{E}[a_t \|\mathbf{a}_i\|^2] + \mathbb{E}[c_t \|b \mathbf{e}_i\|^2] \\
&= a_t d + c_t \delta^2 D
\end{aligned}$$

649 where the last equality follows from $\mathbf{a}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ and $\mathbf{e}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$.

650 Without loss of generality, assume the $j = k + 1$, we have:

$$\begin{aligned}
\mathbb{E}[g_l(\mathbf{x}_0, t)] &= \mathbb{E}[a_t \|\mathbf{U}_l^T \mathbf{x}_0\|^2 + c_t \|\mathbf{U}_l^{\perp T} \mathbf{x}_0\|^2] \\
&= \mathbb{E}[a_t \|\mathbf{U}_l^T (\mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i)\|^2] + \mathbb{E}[c_t \|\mathbf{U}_l^{\perp T} (\mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i)\|^2] \\
&= \mathbb{E}[a_t \|b \mathbf{e}_i[1 : d]\|^2] + \mathbb{E}\left[c_t \left\| \begin{bmatrix} \mathbf{a}_i \\ \mathbf{0} \in \mathbb{R}^{D-d} \end{bmatrix} + b \begin{bmatrix} \mathbf{0} \in \mathbb{R}^d \\ \mathbf{e}_i[d : D] \end{bmatrix} \right\|^2\right] \\
&= a_t \delta^2 d + c_t (d + \delta^2 (D - d))
\end{aligned}$$

651 Plug a_t and b_t back with the exponentials, we get \hat{w}_k and \hat{w}_l .

652

653 Now we prove (14):

$$\begin{aligned}
\mathbf{U}_k \mathbf{U}_k^T f^*(\mathbf{x}_0, t) &= \hat{w}_k \mathbf{U}_k \mathbf{U}_k^T \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \right) \mathbf{x}_0 \\
&\quad + \sum_{l \neq k} \hat{w}_l \mathbf{U}_k \mathbf{U}_k^T \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_l \mathbf{U}_l^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_l^\perp \mathbf{U}_l^{\perp T} \right) \mathbf{x}_0 \\
&= \hat{w}_k \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T \mathbf{x}_0 \right) + \sum_{l \neq k} \hat{w}_l \left(\frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T \mathbf{x}_0 \right) \\
&= \left(\frac{\hat{w}_k}{1 + \sigma_t^2} + \frac{(K-1)\delta^2 \hat{w}_l}{\delta^2 + \sigma_t^2} \right) \mathbf{U}_k \mathbf{U}_k^T (\mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i) \\
&= \left(\frac{\hat{w}_k}{1 + \sigma_t^2} + \frac{(K-1)\delta^2 \hat{w}_l}{\delta^2 + \sigma_t^2} \right) \mathbf{U}_k \mathbf{a}_i
\end{aligned}$$

654 Since $\mathbf{U}_k \in \mathcal{O}^{n \times d}$:

$$\mathbb{E}[\|\mathbf{U}_k \mathbf{U}_k^T f^*(\mathbf{x}_0, t)\|^2] = \left(\frac{\hat{w}_k}{1 + \sigma_t^2} + \frac{(K-1)\delta^2 \hat{w}_l}{\delta^2 + \sigma_t^2} \right)^2 d$$

655 and similarly for (15):

$$\begin{aligned}
\mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}_0, t) &= \hat{w}_k \mathbf{U}_l \mathbf{U}_l^T \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \right) \mathbf{x}_0 \\
&\quad + \hat{w}_l \mathbf{U}_l \mathbf{U}_l^T \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_l \mathbf{U}_l^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_l^\perp \mathbf{U}_l^{\perp T} \right) \mathbf{x}_0 \\
&\quad + \sum_{j \neq k, l} \hat{w}_j \mathbf{U}_l \mathbf{U}_l^T \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_j \mathbf{U}_j^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_j^\perp \mathbf{U}_j^{\perp T} \right) \mathbf{x}_0 \\
&= \hat{w}_k \left(\frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_l \mathbf{U}_l^T \mathbf{x}_0 \right) + \hat{w}_l \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_l \mathbf{U}_l^T \mathbf{x}_0 \right) + \sum_{j \neq k, l} \hat{w}_j \left(\frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_l \mathbf{U}_l^T \mathbf{x}_0 \right) \\
&= \left(\frac{\hat{w}_l}{1 + \sigma_t^2} + \frac{\delta^2(\hat{w}_k + (K-2)\hat{w}_j)}{\delta^2 + \sigma_t^2} \right) \mathbf{U}_l \mathbf{U}_l^T (\mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i) \\
&= \left(\frac{\hat{w}_l}{1 + \sigma_t^2} + \frac{\delta^2(\hat{w}_k + (K-2)\hat{w}_l)}{\delta^2 + \sigma_t^2} \right) b \mathbf{U}_l \mathbf{e}_i [1 : d]
\end{aligned}$$

656 where the third equality follows since $\hat{w}_j = \hat{w}_l$ for all $j \neq k, l$. Further, we have:

$$\mathbb{E}[\|\mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}_0, t)\|^2] = \left(\frac{\hat{w}_l}{1 + \sigma_t^2} + \frac{\delta^2(\hat{w}_k + (K-2)\hat{w}_l)}{\delta^2 + \sigma_t^2} \right)^2 \delta^2 d$$

657 Next, we consider (16):

$$\begin{aligned}
\mathbf{U}_\perp \mathbf{U}_\perp^T f^*(\mathbf{x}_0, t) &= \hat{w}_k \mathbf{U}_\perp \mathbf{U}_\perp^T \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_k \mathbf{U}_k^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_k^\perp \mathbf{U}_k^{\perp T} \right) \mathbf{x}_0 \\
&\quad + \sum_{l \neq k} \hat{w}_l \mathbf{U}_\perp \mathbf{U}_\perp^T \left(\frac{1}{1 + \sigma_t^2} \mathbf{U}_l \mathbf{U}_l^T + \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_l^\perp \mathbf{U}_l^{\perp T} \right) \mathbf{x}_0 \\
&= \hat{w}_k \left(\frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_\perp \mathbf{U}_\perp^T \mathbf{x}_0 \right) + \sum_{l \neq k} \hat{w}_l \left(\frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_\perp \mathbf{U}_\perp^T \mathbf{x}_0 \right) \\
&= \frac{\delta^2}{\delta^2 + \sigma_t^2} \mathbf{U}_\perp \mathbf{U}_\perp^T (\mathbf{U}_k \mathbf{a}_i + b \mathbf{U}_k^\perp \mathbf{e}_i) \\
&= \frac{\delta^3}{\delta^2 + \sigma_t^2} \mathbf{U}_\perp \mathbf{e}_i [(K-1)d : D]
\end{aligned}$$

658 Hence:

$$\mathbb{E}[\|\mathbf{U}_\perp \mathbf{U}_\perp^T f^*(\mathbf{x}_0, t)\|^2] = \frac{\delta^6(n - Kd)}{(\delta^2 + \sigma_t^2)^2}$$

659 Lastly, we prove (17). Given that the subspaces of all classes and the complement space are both
660 orthonormal and mutually orthogonal, we can write:

$$\mathbb{E}[\|f^*(\mathbf{x}_0, t)\|^2] = \mathbb{E}[\|\mathbf{U}_k \mathbf{U}_k^T f^*(\mathbf{x}_0, t)\|^2] + \mathbb{E}\left[\sum_{l \neq k} \|\mathbf{U}_l \mathbf{U}_l^T f^*(\mathbf{x}_0, t)\|^2\right] + \mathbb{E}[\|\mathbf{U}_\perp \mathbf{U}_\perp^T f^*(\mathbf{x}_0, t)\|^2]$$

661 Combine terms, we get:

$$\begin{aligned} \mathbb{E}[\|f^*(\mathbf{x}_0, t)\|^2] &= \left(\frac{\hat{w}_k}{1 + \sigma_t^2} + \frac{(K-1)\delta^2 \hat{w}_l}{\delta^2 + \sigma_t^2} \right)^2 d \\ &\quad + (K-1) \left(\frac{\hat{w}_l}{1 + \sigma_t^2} + \frac{\delta^2(\hat{w}_k + (K-2)\hat{w}_l)}{\delta^2 + \sigma_t^2} \right)^2 \delta^2 d + \frac{\delta^6(n - Kd)}{(\delta^2 + \sigma_t^2)^2}. \end{aligned}$$

662

□