
VarDiU: A Variational Diffusive Upper Bound for One-Step Diffusion Distillation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recently, diffusion distillation methods have compressed thousand-step teacher
2 diffusion models into one-step student generators while preserving sample quality.
3 Most existing approaches train the student model using a diffusive divergence
4 whose gradient is approximated via the student’s score function, learned through
5 denoising score matching (DSM). Since DSM training is imperfect, the resulting
6 gradient estimate is inevitably biased, leading to sub-optimal performance. In this
7 paper, we propose *VarDiU* (pronounced /va:rdju:/), a Variational Diffusive Upper
8 Bound that admits an unbiased gradient estimator and can be directly applied to
9 diffusion distillation. Using this objective, we compare our method with Diff-
10 Instruct and demonstrate that it achieves higher generation quality and enables a
11 more efficient and stable training procedure for one-step diffusion distillation.

12 1 Introduction

13 Diffusion models [38, 15, 42, 44] have achieved remarkable success in generating high-quality
14 samples. However, sampling typically requires a large number of denoising steps, making the process
15 computationally inefficient in practice. To address this limitation, many acceleration methods [15,
16 31, 39, 23, 25, 2, 1, 32, 5, 50, 54] have been proposed to reduce the number of function evaluations
17 (NFEs) in pre-trained diffusion models, lowering the cost from thousands of steps to only tens. More
18 recently, distillation-based approaches have further compressed the sampling process to as few as a
19 single step while largely preserving generation quality [61, 36, 4, 41, 14, 18, 21, 28, 37, 51, 57].

20 In this paper, we investigate a family of distillation methods based on divergence minimisation [28,
21 60, 57, 61]. These methods compress the full denoising process of a pre-trained teacher diffusion
22 model into a one-step latent variable model by minimising a diffusive divergence (e.g. Diffusive
23 Reverse KL Divergence [13, 57, 28]) between the student and teacher through their respective score
24 estimations. Although these approaches have demonstrated promising results, their gradients are
25 approximated using the student’s score function, which is learned through denoising score matching
26 (DSM). Because DSM training is inherently imperfect, the resulting gradient estimates are inevitably
27 biased, often leading to sub-optimal generation performance.

28 In this study, we introduce *Variational Diffusive Upper Bound Distillation (VarDiU)*, a method that
29 formulates a variational upper bound of the diffusive divergence whose gradient can be estimated
30 **unbiasedly**. We demonstrate that this unbiased gradient estimation leads to more stable and efficient
31 training, achieving better performance compared to the original diffusive divergence, which relies on
32 biased gradient approximations.

33 2 Background on Diffusion Models

34 Let $p_d(x_0)$ denote the data distribution. Diffusion models define a forward noising process as
 35 $q(x_{0:T}) = p_d(x_0) \prod_{t=1}^T q(x_t|x_{t-1})$ where the transition kernel is a fixed Gaussian distribution
 36 $q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{1-\beta_t}x_{t-1}, \beta_t I)$. This process admits the closed-form distribution,

$$q(x_t|x_0) = \mathcal{N}(x_t|\sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I), \quad (1)$$

37 with $\bar{\alpha}_t = \prod_{s=1}^t (1-\beta_s)$. As $T \rightarrow \infty$, the final state x_T converges to a standard Gaussian distribution,
 38 i.e., $q(x_T) \rightarrow \mathcal{N}(0, I)$. The generative process also uses a sequence of Gaussian distributions:

$$p(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}|\mu_{t-1}(x_t), \Sigma_{t-1}(x_t)), \quad (2)$$

39 where the covariance $\Sigma_{t-1}(x_t)$ may either be learned [31, 32, 1] or fixed [2, 15]. The mean is
 40 estimated via Tweedie’s formula [10, 34]:

$$\mu_{t-1}(x_t) = \frac{1}{\sqrt{1-\beta_t}} \left(x_t + \beta_t \nabla_{x_t} \log p_d(x_t) \right), \quad (3)$$

41 where the score function $\nabla_{x_t} \log p_d(x_t)$ is approximated with a network $\mathbf{s}_\psi(x_t; t)$ which is learned
 42 through denoising score matching (DSM) [48, 42]:

$$\mathcal{L}_{\text{DSM}}(\psi) = \int p_d(x_0) q(x_t|x_0) \left\| \mathbf{s}_\psi(x_t; t) - \nabla_{x_t} \log q(x_t|x_0) \right\|_2^2 dx_t dx_0. \quad (4)$$

43 The classical diffusion models require thousands of time steps to generate high-quality samples,
 44 making them computationally expensive. In the following section, we discuss a family of methods
 45 that can distill multi-step diffusion processes into one-step generative models.

46 2.1 Score-Based Distillation via Diffusive Reverse KL Minimisation

47 Score-based distillation methods aim to compress a teacher diffusion model, which we also denote as
 48 p_d , into a one-step implicit generative model [11, 16, 58], defined as

$$p_\theta(x_0) = \int \delta(x_0 - g_\theta(z)) p(z) dz, \quad (5)$$

49 where $\delta(\cdot)$ is the Dirac delta function, $p(z)$ is a standard Gaussian prior, and g_θ is a deterministic
 50 neural network that generates samples in a single step. The training objective is formulated as the
 51 *Diffusive KL divergence* (DiKL) [58, 49, 28, 13, 57], which measures the discrepancy between the
 52 student and teacher distributions across all intermediate noise levels:

$$\text{DiKL}(p_\theta||p_d) \equiv \int_0^1 \omega(t) \text{KL}(p_\theta^{(t)}(x_t)||p_d^{(t)}(x_t)) dt, \quad (6)$$

53 where $\omega(t)$ is a positive weighting function. Here, $p_\theta^{(t)}(x_t) = \int q(x_t|x_0) p_\theta(x_0) dx_0$, and $p_d^{(t)}(x_t)$ is
 54 the marginal distribution induced by the pre-trained teacher diffusion model under the same noise
 55 kernel $q(x_t|x_0)$. This divergence is well defined even when the underlying distributions have disjoint
 56 support or don’t even allow density functions [58]. However, direct estimation of this divergence is
 57 intractable. One can derive the gradient w.r.t. θ at time t as

$$\nabla_\theta \text{KL}(p_\theta^{(t)}||p_d^{(t)}) = \int \left[\left(\nabla_{x_t} \log p_\theta^{(t)}(x_t) - \nabla_{x_t} \log p_d^{(t)}(x_t) \right) \frac{\partial x_t}{\partial \theta} \right] dx_t, \quad (7)$$

58 see A.2 in [28] for a derivation. The teacher score $\nabla_{x_t} \log p_d^{(t)}(x_t)$ is available from the pre-trained
 59 diffusion model. However, the student score $\nabla_{x_t} \log p_\theta^{(t)}(x_t)$ is unknown. Fortunately, since we can
 60 sample from the student distribution, it can be estimated using denoising score matching (DSM):

$$\mathcal{L}_{\text{DSM}}(\psi') = \int p_\theta(x_0) q(x_t|x_0) \left\| \mathbf{s}_{\psi'}(x_t; t) - \nabla_{x_t} \log q(x_t|x_0) \right\|_2^2 dx_t dx_0. \quad (8)$$

61 This method is known as VSD [49] or Diff-Instruct [28] in the context of diffusion distillation. In
 62 practice, DSM approximations are imperfect due to the limited capacity of the score network $\mathbf{s}_{\psi'}$. As
 63 a result, *the direct gradient estimation of DiKL is biased*, which leads to suboptimal training of the
 64 distilled student model. In the next section, we introduce a variational upper bound on the reverse KL
 65 divergence that enables unbiased gradient estimation for diffusion distillation.

66 **3 Variational Diffusive Upper Bound Distillation**

67 Instead of relying on optimisation-based estimation of the DiKL objective, we propose a variational
 68 upper bound whose gradient can be estimated unbiasedly. Specifically, we propose to extend the
 69 RKL upper bound of [56] into the diffusion space,

$$\text{KL}(p_\theta^{(t)}(x)||p_d^{(t)}(x)) \leq \text{KL}(p_\theta^{(t)}(x_t|z)p(z)||p_d^{(t)}(x)q_\phi^{(t)}(z|x_t)) \equiv U^{(t)}(\theta, \phi).$$

70 The inequality follows by Jensen’s inequality; and the equality can be attained when the variational
 71 posterior is equal to the true posterior $q_\phi(z|x_t; t) = p_\theta(z|x_t) \propto p_\theta(x_t|z)p(z)$ (see Appendix A for a
 72 proof). We are then ready to define the diffusive upper bound as

$$\text{DiU}(\theta, \phi) \equiv \int_0^1 w(t)U^{(t)}(\theta, \phi)dt \geq \int_0^1 w(t)\text{KL}(p_\theta^{(t)}||p_d^{(t)})dt \equiv \text{DiKL}(p_\theta||p_d), \quad (9)$$

73 where the bound is tight when $q_\phi^{(t)}(z|x_t) = p_\theta^{(t)}(z|x_t) \propto p_\theta^{(t)}(x_t|z)p(z)$ for all t almost everywhere.

74 **3.1 Tractable Upper Bound Estimation**

75 We then discuss how to estimate this upper bound. Writing the model joint as $p_\theta^{(t)}(x_t, z) =$
 76 $p_\theta^{(t)}(x_t|z)p(z)$, the upper bound at time step t can be written as

$$U^{(t)}(\theta, \phi) = -\mathbb{H}(p_\theta^{(t)}(x_t, z)) - \iint p_\theta^{(t)}(x_t, z) \left(\log p_d^{(t)}(x_t)dx_t - \log q_\phi^{(t)}(z|x_t) \right) dzdx_t. \quad (10)$$

77 where $\mathbb{H}(p_\theta^{(t)}(x_t, z))$ denotes the joint entropy, i.e. $\mathbb{H}(p) = -\int p(x) \log p(x)dx$. In the original
 78 reverse KL upper bound in [56], the joint entropy is trained with the reparametrisation trick. In our
 79 case, when the student model is an implicit model, we present the following surprising fact that *the*
 80 *joint entropy is a constant*, which removes the requirement of joint entropy estimation.

81 **Proposition 3.1.** *Let $p_\theta(x, z) = p_\theta(x|z)p(z)$ with a fixed prior $p(z)$ and $p_\theta(x|z) =$
 82 $\mathcal{N}(x; g_\theta(z), \sigma^2 I)$. Then the joint entropy $\mathbb{H}(p_\theta(x, z))$ is a constant and independent of θ .*

83 See Appendix D for a proof. This should be distinguished from the marginal entropy $\mathbb{H}(p_\theta^{(t)}(x_t))$,
 84 which still depends on θ , also see Section 3.3 for a deeper discussion. This observation greatly
 85 simplifies training by eliminating the need for entropy estimation.

86 The second term of Equation 9 cannot be directly optimised via automatic differentiation because
 87 the teacher density $p_d^{(t)}(x_t)$ is unknown. However, in typical distillation settings, the score function
 88 $\nabla_{x_t} \log p_d^{(t)}(x_t)$ is assumed to be available. Leveraging this, we propose an alternative loss whose
 89 gradient is equivalent and can be computed via the score function:

$$\nabla_\theta \int p_\theta^{(t)}(x_t) \log p_d^{(t)}(x_t)dx_t = \nabla_\theta \int p_\theta^{(t)}(x_t) \left(x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} \right) dx_t, \quad (11)$$

90 where $[\cdot]_{\text{sg}}$ denotes the stop-gradient operator. The proof involves using the reparametrisation trick
 91 and chain rule; see Appendix B for details. The final loss objective at time t becomes

$$U^{(t)}(\theta, \phi) \doteq - \iint p_\theta^{(t)}(x_t, z) \left(x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} + \log q_\phi^{(t)}(z|x_t) \right) dzdx_t, \quad (12)$$

92 where \doteq denotes equivalence up to a constant. Therefore, when using a simple Gaussian variational
 93 family $q_\phi^{(t)}(z|x_t)$, this bound can be unbiasedly estimated. A loss used practically can be found in
 94 Appendix C. We also propose a novel noise schedule and variance reduction techniques for training
 95 with this objective, see Appendix F.2 for details.

96 **3.2 Improving Variational Inference with Normalising Flows**

97 The upper bound becomes tight, i.e., equal to the DiKL objective, when the variational posterior
 98 matches the true posterior: $q_\phi^{(t)}(z|x_t) = p_\theta^{(t)}(z|x_t)$. Simple Gaussian variational distributions often

99 struggle to approximate complex posteriors, particularly in cases where the true posterior is multi-
 100 modal or exhibits strong correlations. To address this limitation, we could use normalising flow
 101 [47, 46, 8, 33, 55, 9, 19] to increase the flexibility of the posterior [33]. Specifically, we define:

$$q_\phi^{(t)}(z|x_t) = \int \delta(z - f_\phi(a, x_t, t)) r_\phi(a|x_t; t) da, \quad (13)$$

102 where f is an invertible function and the base distribution is defined as a Gaussian indexed by t .

$$r_\phi(a|x_t; t) = \mathcal{N}(a; \mu_\phi(x_t; t), \sigma_t^2 \text{diag}(\sigma_\phi^2(x_t; t))). \quad (14)$$

103 The time-dependent scaling factor σ_t modulates the entropy of the base distribution to match the
 104 noise level at diffusion step t . The corresponding log-density of $q_\phi^{(t)}(z|x_t)$ can be computed using
 105 the change-of-variables formula:

$$\log q_\phi^{(t)}(z = f_\phi(a; t)|x_t; t) = \log r_\phi(a|x_t; t) - \log \left| \det \frac{\partial f_\phi}{\partial a} \right|, \quad (15)$$

106 We can then substitute this log-density into Equation (12) to obtain a tractable objective for distillation.

107 3.3 An Information Maximisation View of the Variational Diffusive Upper Bound

108 In the previous section, we showed that the joint entropy term in Equation (12) is constant with
 109 respect to θ . As a result, the variational upper bound at diffusion step t simplifies to:

$$\mathbb{U}^{(t)}(\theta, \phi) = - \int p_\theta^{(t)}(x_t) \log p_d^{(t)}(x_t) dx_t - \iint p_\theta^{(t)}(x_t, z) \log q_\phi^{(t)}(z|x_t) dz dx_t. \quad (16)$$

110 Now consider the original DiKL divergence at time step t , defined as:

$$\text{KL}(p_\theta^{(t)} \| p_d^{(t)}) = - \int p_\theta^{(t)}(x_t) \log p_d^{(t)}(x_t) dx_t + \int p_\theta^{(t)}(x_t) \log p_\theta^{(t)}(x_t) dx_t. \quad (17)$$

111 We observe that the first term in both objectives (Equations (16) and (17)) is identical. The second term
 112 in Equation (17) corresponds to the negative entropy of the marginal distribution, i.e., $-\mathbb{H}(p_\theta^{(t)}(x_t))$,
 113 which depends on θ . We can then obtain an alternative perspective for justifying the proposed upper
 114 bound, formalised in the following theorem:

115 **Theorem 3.2.** *By the chain rule of entropy, we have: $\mathbb{H}(p_\theta^{(t)}(x_t, z)) = \mathbb{H}(p_\theta^{(t)}(z|x_t)) + \mathbb{H}(p_\theta^{(t)}(x_t))$.
 116 Since $\mathbb{H}(p_\theta^{(t)}(x_t, z))$ is constant with respect to θ , it follows that:*

$$\max_\theta \mathbb{H}(p_\theta^{(t)}(x_t)) = \min_\theta \mathbb{H}(p_\theta^{(t)}(z|x_t)) \leq \min_{\theta, \phi} - \iint p_\theta^{(t)}(x_t, z) \log q_\phi^{(t)}(z|x_t) dz dx_t, \quad (18)$$

117 where $q_\phi^{(t)}(z|x_t)$ is a variational approximation to the true posterior $p_\theta^{(t)}(z|x_t)$.

118 See Appendix E for a proof. This result mirrors the variational approach in the Information Maximisa-
 119 tion (IM) algorithm [3], which maximises mutual information by minimising the conditional entropy
 120 using a variational decoder. Hence, this entropy bound can be viewed as a continuous latent-variable
 121 adaptation of the IM algorithm, where entropy maximisation serves to promote informative and
 122 diverse model outputs.

123 4 Experiments

124 To demonstrate the effectiveness of our approach, we consider a toy 2d mixture of 40 Gaussians
 125 (MoG-40) proposed in [30], where the means are uniformly distributed over $[-40, 40]^2$. We train
 126 VarDiU with two types of posteriors: **(a)** a Gaussian posterior with learnable mean and variance
 127 (VarDiU-Gaussian), and **(b)** a Neural Spline Flow (NSF; [9]) with a learnable-parameters Gaussian
 128 base distribution (VarDiU-NSF). We compare VarDiU to Diff-Instruct [28] using varying numbers of
 129 student score training steps (1, 5, 10). Experimental details can be found in Appendix F.

130 **Settings** We consider three practical settings: **(1)** access to the true analytical score; **(2)** access to the
 131 training dataset; **(3)** access to a teacher score pre-trained by a diffusion model. For the given training

Table 1: **Comparison across different settings.** (\uparrow) denotes the higher is better and (\downarrow) denotes the lower is better. Best performance is in bold. Second best is underlined. “True” indicates the samples from the target distribution. Diff-inst (k) means Diff-instruct trained with k score steps. We can find VarDiU achieve best results with true and empirical score and remain competitive with learned score.

(a) Comparison with true score

Metrics	True	Diff-Inst (1)	Diff-Inst (5)	Diff-Inst (10)	VarDiU Gaussian	VarDiU NSF
Log-Density (\uparrow)	-6.65	-10.17 \pm 1.77	-8.09 \pm 0.86	-7.86 \pm 0.56	<u>-7.00 \pm 0.02</u>	-6.99 \pm 0.01
Log-MMD (\downarrow)	/	-5.59 \pm 0.36	-6.09 \pm 0.40	-5.79 \pm 0.36	<u>-6.86 \pm 0.29</u>	-7.33 \pm 0.17

(b) Comparison with training data

Metrics	True	Diff-Inst (1)	Diff-Inst (5)	Diff-Inst (10)	VarDiU Gaussian	VarDiU NSF
Log-Density (\uparrow)	-6.65	-9.65 \pm 0.62	-8.38 \pm 0.42	-7.96 \pm 0.62	<u>-7.09 \pm 0.03</u>	-7.01 \pm 0.02
Log-MMD (\downarrow)	/	-5.51 \pm 0.26	-5.60 \pm 0.30	-5.63 \pm 0.68	<u>-6.36 \pm 0.22</u>	-6.81 \pm 0.30

(c) Comparison with learned score

Metrics	True	EDM	Diff-Inst (1)	Diff-Inst (5)	Diff-Inst (10)	VarDiU Gaussian	VarDiU NSF
Log-Density (\uparrow)	-6.65	-6.69	-10.88 \pm 1.55	-9.21 \pm 0.70	-8.47 \pm 0.56	<u>-8.13 \pm 0.19</u>	-7.89 \pm 0.15
Log-MMD (\downarrow)	/	-6.67	-5.24 \pm 0.30	-5.57 \pm 0.27	-5.82 \pm 0.27	-5.64 \pm 0.25	<u>-5.68 \pm 0.30</u>

132 data case, we use the empirical score, defined by $\nabla_{x_t} \log \hat{p}_d(x_t) = \frac{1}{N} \sum_{n=1}^N \nabla_{x_t} \log q(x_t | x_0^{(n)})$ with
 133 the given dataset $\mathcal{D} = \{x_0^{(n)}\}_{n=1}^N$, where $q(x_t | x_0)$ matches the diffusion schedule of the VarDiU.
 134 This empirical score is a consistent estimator of the true data score [42, 43, 56], see Appendix F.3 for
 135 details. The learned score estimator is obtained by training EDM [17], which is a popular class of
 136 diffusion models, on the observed data. For each setting, we report results over 10 independent runs.
 137 The noise schedule and weighting strategy are detailed in Appendix F.2.

138 **Metrics** We report the *log-density* of generated samples under the true data distribution, which reflects
 139 sample quality but not sample diversity. We additionally compute the Maximum Mean Discrepancy
 140 (MMD; [12]) with five kernels which can measure diversity, see for details in Appendix F.1.

141 **Sample Quality** We compare samples generated by one-step models trained with VarDiU and
 142 Diff-Instruct across three practical settings: using the true score, the training dataset, and a learned
 143 score (see Figure 1 for visualizations). For Diff-Instruct, we train the model with 10 score-matching
 144 steps per generator gradient step to obtain the best sample quality we can, whereas the original
 145 Diff-Instruct paper [28] uses only a single step. Both VarDiU-Gaussian and VarDiU-NSF produce
 146 cleaner and sharper samples than Diff-Instruct when using the true score or the training dataset. With
 147 a learned score, however, all methods yield different sub-optimal results due to biases in the provided
 148 scores. We also compare log-density and log-MMD across the samples (see Tables 1a to 1c). We
 149 can find VarDiU consistently performs best under reliable evaluations (true score and training data),
 150 demonstrating high accuracy and stability with tight confidence intervals, while Diff-Instruct lags
 151 despite larger steps. Under learned scores, evaluations are noisier: VarDiU achieves higher log-density
 152 but slightly worse MMD, making comparisons unreliable due to bias in the given pre-trained scores.
 153 Overall, VarDiU is more accurate and stable across settings.

154 **Training Stability and Efficiency** We present plots of log-MMD against both generator gradient
 155 iterations and total training time in Figure 2; additional results are provided in Appendix Figures 4
 156 to 6. Each curve shows the mean of 10 independent runs, with shaded regions indicating one standard
 157 deviation. For VarDiU, we apply a noise scheme with an annealing schedule. For Diff-Instruct, we
 158 use both the same annealing schedule (denoted as Diff-Instruct-a) and a non-annealing variant (see
 159 Appendix F.2 for details). In the log-MMD vs. generator gradient step plot, we find that under the
 160 same generator gradient budget (1M steps), VarDiU methods significantly outperform Diff-Instruct
 161 in both MMD values and stability (smaller variance). Moreover, incorporating a normalizing flow
 162 further improves performance compared to the Gaussian posterior, highlighting the importance of
 163 accurate generator gradient estimation. In the log-MMD vs. training time plot, we observe that under
 164 the same training hours, VarDiU-Gaussian clearly outperforms all other methods.

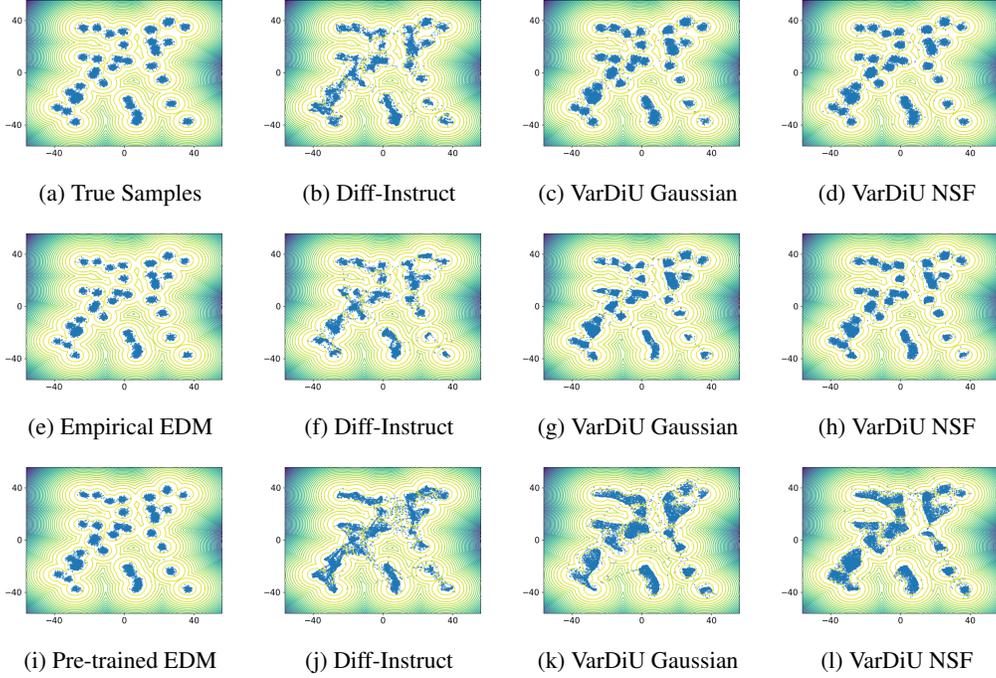


Figure 1: The three rows correspond to the true score, training data, and learned score settings. Diff-Instruct models use 10 score-matching steps per generator step. Empirical EDM indicates the samples generated by EDM sampler with empirical scores. See Tables 1a to 1c for evaluation results.

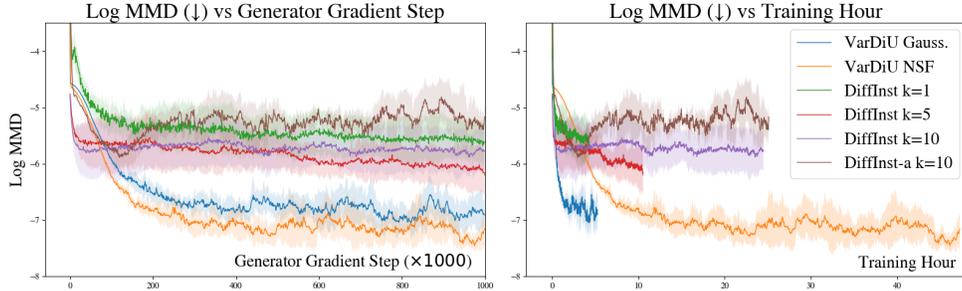


Figure 2: The left figure shows log-MMD over 1M generator steps, reflecting the efficiency of generator gradient estimation across different methods. The right figure shows total training time for 1M generator steps. In both, VarDiU clearly surpasses Diff-Instruct in quality, stability, and efficiency.

165 5 Conclusion and Future Work

166 In this paper, we introduced a novel method for diffusion distillation by proposing a variational
 167 diffusive upper bound on the DiKL divergence whose gradient can be estimated without bias. Our
 168 bound extends the Reverse KL upper bound from [56] to the diffusion setting and further eliminates
 169 the need to estimate the joint entropy. We also adapt this bound naturally to scenarios where only the
 170 teacher model’s score function is accessible. Additionally, we establish a connection between our
 171 method and the information maximisation (IM [3]) framework, see Appendix G for a discussion on
 172 related works. Empirically, we demonstrate on a toy problem that our method achieves improved
 173 one-step generation quality while significantly reducing training time.

174 As future work, we are extending this approach to higher-dimensional data, such as image or video
 175 generation, which we believe to be a promising direction.

References

- 176
- 177 [1] F. Bao, C. Li, J. Sun, J. Zhu, and B. Zhang. Estimating the optimal covariance with imperfect mean in
178 diffusion probabilistic models. *arXiv preprint arXiv:2206.07309*, 2022.
- 179 [2] F. Bao, C. Li, J. Zhu, and B. Zhang. Analytic-dpm: an analytic estimate of the optimal reverse variance in
180 diffusion probabilistic models. *arXiv preprint arXiv:2201.06503*, 2022.
- 181 [3] D. Barber and F. Agakov. The im algorithm: a variational approach to information maximization. *Advances*
182 *in neural information processing systems*, 2004.
- 183 [4] D. Berthelot, A. Autef, J. Lin, D. A. Yap, S. Zhai, S. Hu, D. Zheng, W. Talbott, and E. Gu. Tract: Denoising
184 diffusion models with transitive closure time-distillation. *arXiv preprint arXiv:2303.04248*, 2023.
- 185 [5] V. D. Bortoli, A. Galashov, J. S. Guntupalli, G. Zhou, K. Murphy, A. Gretton, and A. Doucet. Distributional
186 diffusion models with scoring rules, 2025.
- 187 [6] Z. Botev and A. Ridder. Variance reduction. *Wiley statsRef: Statistics reference online*, 136:476, 2017.
- 188 [7] V. De Bortoli, J. Thornton, J. Heng, and A. Doucet. Diffusion schrödinger bridge with applications to
189 score-based generative modeling. *Advances in Neural Information Processing Systems*, 2021.
- 190 [8] L. Dinh, D. Krueger, and Y. Bengio. Nice: Non-linear independent components estimation. *arXiv preprint*
191 *arXiv:1410.8516*, 2014.
- 192 [9] C. Durkan, A. Bekasov, I. Murray, and G. Papamakarios. Neural spline flows. *Advances in Neural*
193 *Information Processing Systems*, 2019.
- 194 [10] B. Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*,
195 106(496):1602–1614, 2011.
- 196 [11] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio.
197 Generative adversarial nets. *Advances in neural information processing systems*, 2014.
- 198 [12] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *The*
199 *journal of machine learning research*, (1):723–773, 2012.
- 200 [13] J. He, W. Chen, M. Zhang, D. Barber, and J. M. Hernández-Lobato. Training neural samplers with reverse
201 diffusive kl divergence. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2025.
- 202 [14] J. Heek, E. Hoogeboom, and T. Salimans. Multistep consistency models. *arXiv preprint arXiv:2403.06807*,
203 2024.
- 204 [15] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information*
205 *processing systems*, 2020.
- 206 [16] F. Huszár. Variational inference using implicit distributions. *arXiv preprint arXiv:1702.08235*, 2017.
- 207 [17] T. Karras, M. Aittala, T. Aila, and S. Laine. Elucidating the design space of diffusion-based generative
208 models. *Advances in Neural Information Processing Systems*, 2022.
- 209 [18] D. Kim, C.-H. Lai, W.-H. Liao, N. Murata, Y. Takida, T. Uesaka, Y. He, Y. Mitsufuji, and S. Ermon.
210 Consistency trajectory models: Learning probability flow ODE trajectory of diffusion. In *International*
211 *Conference on Learning Representations*, 2024.
- 212 [19] I. Kobyzev, S. J. Prince, and M. A. Brubaker. Normalizing flows: An introduction and review of current
213 methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020.
- 214 [20] D. P. Kroese, T. Taimre, and Z. I. Botev. *Handbook of monte carlo methods*. John Wiley & Sons, 2013.
- 215 [21] L. Li and J. He. Bidirectional consistency models. *arXiv preprint arXiv:2403.18035*, 2024.
- 216 [22] Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le. Flow matching for generative modeling. In
217 *International Conference on Learning Representations*, 2023.
- 218 [23] L. Liu, Y. Ren, Z. Lin, and Z. Zhao. Pseudo numerical methods for diffusion models on manifolds. *arXiv*
219 *preprint arXiv:2202.09778*, 2022.
- 220 [24] X. Liu, C. Gong, et al. Flow straight and fast: Learning to generate and transfer data with rectified flow. In
221 *The Eleventh International Conference on Learning Representations*, 2023.

- 222 [25] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic
223 model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 2022.
- 224 [26] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu. Dpm-solver++: Fast solver for guided sampling of
225 diffusion probabilistic models. *Machine Intelligence Research*, pages 1–22, 2025.
- 226 [27] S. Luo, Y. Tan, L. Huang, J. Li, and H. Zhao. Latent consistency models: Synthesizing high-resolution
227 images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023.
- 228 [28] W. Luo, T. Hu, S. Zhang, J. Sun, Z. Li, and Z. Zhang. Diff-instruct: A universal approach for transferring
229 knowledge from pre-trained diffusion models. *Advances in Neural Information Processing Systems*, 2023.
- 230 [29] C. Meng, R. Rombach, R. Gao, D. Kingma, S. Ermon, J. Ho, and T. Salimans. On distillation of guided
231 diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
232 2023.
- 233 [30] L. I. Midgley, V. Stimper, G. N. Simm, B. Schölkopf, and J. M. Hernández-Lobato. Flow annealed
234 importance sampling bootstrap. In *International Conference on Learning Representations*, 2023.
- 235 [31] A. Q. Nichol and P. Dhariwal. Improved denoising diffusion probabilistic models. In *International
236 conference on machine learning*. PMLR, 2021.
- 237 [32] Z. Ou, M. Zhang, A. Zhang, T. Z. Xiao, Y. Li, and D. Barber. Improving probabilistic diffusion models
238 with optimal covariance matching. *International Conference on Learning Representations*, 2025.
- 239 [33] D. Rezende and S. Mohamed. Variational inference with normalizing flows. In *International Conference
240 on Machine Learning*. PMLR, 2015.
- 241 [34] H. E. Robbins. An empirical bayes approach to statistics. In *Breakthroughs in Statistics: Foundations and
242 basic theory*, pages 388–394. Springer, 1992.
- 243 [35] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with
244 latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
245 recognition*, 2022.
- 246 [36] T. Salimans and J. Ho. Progressive distillation for fast sampling of diffusion models. In *International
247 Conference on Learning Representations*, 2022.
- 248 [37] T. Salimans, T. Mensink, J. Heek, and E. Hoogeboom. Multistep distillation of diffusion models via
249 moment matching. *arXiv preprint arXiv:2406.04103*, 2024.
- 250 [38] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using
251 nonequilibrium thermodynamics. In *International conference on machine learning*. PMLR, 2015.
- 252 [39] J. Song, C. Meng, and S. Ermon. Denoising diffusion implicit models. In *International Conference on
253 Learning Representations*, 2021.
- 254 [40] Y. Song and P. Dhariwal. Improved techniques for training consistency models. In *International Conference
255 on Learning Representations*, 2024.
- 256 [41] Y. Song, P. Dhariwal, M. Chen, and I. Sutskever. Consistency models. In *International Conference on
257 Machine Learning*. PMLR, 2023.
- 258 [42] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in
259 Neural Information Processing Systems*, 2019.
- 260 [43] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative
261 modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- 262 [44] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative
263 modeling through stochastic differential equations. 2021.
- 264 [45] V. Stimper, D. Liu, A. Campbell, V. Berenz, L. Ryll, B. Schölkopf, and J. M. Hernández-Lobato. normflows:
265 A pytorch package for normalizing flows. *Journal of Open Source Software*, 8(86):5361, 2023.
- 266 [46] E. G. Tabak and C. V. Turner. A family of nonparametric density estimation algorithms. *Communications
267 on Pure and Applied Mathematics*, 66(2):145–164, 2013.
- 268 [47] E. G. Tabak and E. Vanden-Eijnden. Density estimation by dual ascent of the log-likelihood. *Commun.
269 Math. Sci.*, 8(1):217–233, 2010.

- 270 [48] P. Vincent. A connection between score matching and denoising autoencoders. *Neural computation*,
271 23(7):1661–1674, 2011.
- 272 [49] Z. Wang, C. Lu, Y. Wang, F. Bao, C. Li, H. Su, and J. Zhu. Prolificdreamer: High-fidelity and diverse
273 text-to-3d generation with variational score distillation. *Advances in Neural Information Processing*
274 *Systems*, 2024.
- 275 [50] Z. Xiao, K. Kreis, and A. Vahdat. Tackling the generative learning trilemma with denoising diffusion gans.
276 *arXiv preprint arXiv:2112.07804*, 2021.
- 277 [51] S. Xie, Z. Xiao, D. Kingma, T. Hou, Y. N. Wu, K. P. Murphy, T. Salimans, B. Poole, and R. Gao. Em
278 distillation for one-step diffusion models. *Advances in Neural Information Processing Systems*, 2024.
- 279 [52] Y. Xu, W. Nie, and A. Vahdat. One-step diffusion models with f -divergence distribution matching. *arXiv*
280 *preprint arXiv:2502.15681*, 2025.
- 281 [53] T. Yin, M. Gharbi, T. Park, R. Zhang, E. Shechtman, F. Durand, and B. Freeman. Improved distribution
282 matching distillation for fast image synthesis. *Advances in neural information processing systems*, 2024.
- 283 [54] L. Yu, T. Xie, Y. Zhu, T. Yang, X. Zhang, and C. Zhang. Hierarchical semi-implicit variational inference
284 with application to diffusion model acceleration. *Advances in Neural Information Processing Systems*,
285 2024.
- 286 [55] S. Zhai, R. Zhang, P. Nakkiran, D. Berthelot, J. Gu, H. Zheng, T. Chen, M. A. Bautista, N. Jaitly, and
287 J. Susskind. Normalizing flows are capable generative models. *arXiv preprint arXiv:2412.06329*, 2024.
- 288 [56] M. Zhang, T. Bird, R. Habib, T. Xu, and D. Barber. Variational f -divergence minimization. *arXiv preprint*
289 *arXiv:1907.11891*, 2019.
- 290 [57] M. Zhang, W. Chen, J. He, Z. Ou, J. M. Hernández-Lobato, B. Schölkopf, and D. Barber. Towards training
291 one-step diffusion models without distillation. *arXiv preprint arXiv:2502.08005*, 2025.
- 292 [58] M. Zhang, P. Hayes, T. Bird, R. Habib, and D. Barber. Spread divergence. In *International Conference on*
293 *Machine Learning*. PMLR, 2020.
- 294 [59] L. Zhou, S. Ermon, and J. Song. Inductive moment matching. In *International Conference on Machine*
295 *Learning*, 2025.
- 296 [60] M. Zhou, H. Zheng, Y. Gu, Z. Wang, and H. Huang. Adversarial score identity distillation: Rapidly
297 surpassing the teacher in one step. In *International Conference on Learning Representations*, 2025.
- 298 [61] M. Zhou, H. Zheng, Z. Wang, M. Yin, and H. Huang. Score identity distillation: Exponentially fast
299 distillation of pretrained diffusion models for one-step generation. In *International Conference on Machine*
300 *Learning*, 2024.

301 **Appendix for “VarDiU: A Variational Diffusive Upper**
 302 **Bound for One-Step Diffusion Distillation”**

303 **Contents**

304 **A Proof of Variational Upper Bound** **10**

305 **B Proof of VarDiU Gradient Estimator** **11**

306 **C A Practically Equivalent Variational Diffusive Upper Bound Loss** **11**

307 **D Joint Entropy is Constant** **12**

308 **E Proof of Theorem 3.2** **13**

309 **F Experimental Details** **14**

310 F.1 Implementation Details 14

311 F.2 Noise Weighting and Scheduling 14

312 F.3 Empirical Score Estimation 15

313 F.4 Symmetric Sampling 16

314 **G Related Works** **17**

315 **H More Convergence and Efficiency Results** **18**

316 **A Proof of Variational Upper Bound**

317 **Proposition A.1** (Reverse KL Variational Upper Bound [56]). *For any choice of auxiliary variational*
 318 *distribution $q_\phi^{(t)}(z | x_t)$,*

$$\text{KL}\left(p_\theta^{(t)}(x_t) \parallel p_d^{(t)}(x_t)\right) \leq \text{KL}\left(p_\theta^{(t)}(x_t | z)p(z) \parallel p_d^{(t)}(x_t)q_\phi^{(t)}(z | x_t)\right) =: \mathcal{U}^{(t)}(\theta, \phi). \quad (19)$$

319 *Moreover, equality holds iff $q_\phi^{(t)}(z | x_t) = p_\theta^{(t)}(z | x_t)$ for $p_\theta^{(t)}(x_t)$ -almost every x_t .*

320 *Proof.* We start from the joint KL:

$$\text{KL}\left(p_\theta^{(t)}(x_t | z)p(z) \parallel p_d^{(t)}(x_t)q_\phi^{(t)}(z | x_t)\right) \quad (20)$$

$$= \iint p_\theta^{(t)}(x_t | z)p(z) \log \frac{p_\theta^{(t)}(x_t | z)p(z)}{p_d^{(t)}(x_t)q_\phi^{(t)}(z | x_t)} dx_t dz \quad (21)$$

$$= \iint p_\theta^{(t)}(x_t | z)p(z) [\log p_\theta^{(t)}(x_t | z) - \log p_d^{(t)}(x_t)] dx_t dz$$

$$+ \iint p_\theta^{(t)}(x_t | z)p(z) [\log p(z) - \log q_\phi^{(t)}(z | x_t)] dx_t dz \quad (22)$$

$$= \int p_\theta^{(t)}(x_t) \log \frac{p_\theta^{(t)}(x_t)}{p_d^{(t)}(x_t)} dx_t + \iint p_\theta^{(t)}(x_t | z)p(z) \log \frac{p_\theta^{(t)}(z | x_t)}{q_\phi^{(t)}(z | x_t)} dx_t dz \quad (23)$$

$$= \text{KL}\left(p_\theta^{(t)}(x_t) \parallel p_d^{(t)}(x_t)\right) + \mathbb{E}_{x_t \sim p_\theta^{(t)}(x_t)} \left[\text{KL}\left(p_\theta^{(t)}(z | x_t) \parallel q_\phi^{(t)}(z | x_t)\right) \right]. \quad (24)$$

321 Since the conditional KL is non-negative, it follows that

$$\text{KL}\left(p_\theta^{(t)}(x_t) \parallel p_d^{(t)}(x_t)\right) \leq \text{KL}\left(p_\theta^{(t)}(x_t | z)p(z) \parallel p_d^{(t)}(x_t)q_\phi^{(t)}(z | x_t)\right) =: U^{(t)}(\theta, \phi). \quad (25)$$

322 Equality holds iff $q_\phi^{(t)}(z | x_t) = p_\theta^{(t)}(z | x_t)$ for $p_\theta^{(t)}(x_t)$ almost every x_t . \square

323 B Proof of VarDiU Gradient Estimator

324 **Proposition B.1.** *The gradient can be estimated by*

$$\nabla_\theta \int p_\theta^{(t)}(x_t) \log p_d^{(t)}(x_t) dx_t = \nabla_\theta \int p_\theta^{(t)}(x_t) \left(x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} \right) dx_t, \quad (26)$$

325 where $[\cdot]_{\text{sg}}$ denotes the stop-gradient operator.

326 *Proof.* Since x_t is generated via a differentiable transformation

$$x_t = \mathcal{F}_\theta(g_\theta(z), \epsilon, t) = g_\theta(z) + \epsilon\sigma_t, \quad (27)$$

327 where ϵ is drawn from a noise distribution $p_\epsilon(\epsilon)$ that is independent of θ . Then we have

$$\begin{aligned} & \iint p_\theta^{(t)}(x_t|z)p(z) \log p_d^{(t)}(x_t) dx_t dz \\ &= \iint p(z)p_\epsilon(\epsilon) \log p_d^{(t)}\left(\mathcal{F}_\theta(g_\theta(z), \epsilon, t)\right) d\epsilon dz = \mathbb{E}_{z \sim p(z), \epsilon \sim p_\epsilon(\epsilon)} \left[\log p_d^{(t)}\left(\mathcal{F}_\theta(g_\theta(z), \epsilon, t)\right) \right] \end{aligned} \quad (28)$$

328 Since $\mathcal{F}_\theta(g_\theta(z), \epsilon, t)$ is differentiable in θ and under appropriate regularity conditions, we have

$$\nabla_\theta \iint p_\theta^{(t)}(x_t|z)p(z) \log p_d^{(t)}(x_t) dz dx_t = \mathbb{E}_{z \sim p(z), \epsilon \sim p_\epsilon(\epsilon)} \left[\nabla_\theta \log p_d^{(t)}\left(\mathcal{F}_\theta(g_\theta(z), \epsilon, t)\right) \right] \quad (29)$$

329 Using the chain rule, the gradient inside the expectation becomes

$$\nabla_\theta \log p_d^{(t)}\left(\mathcal{F}_\theta(g_\theta(z), \epsilon, t)\right) = \nabla_{x_t} \log p_d^{(t)}(x_t) \Big|_{x_t = \mathcal{F}_\theta(g_\theta(z), \epsilon, t)} \cdot \frac{\partial \mathcal{F}_\theta(g_\theta(z), \epsilon, t)}{\partial \theta} \quad (30)$$

330 Substituting back into our expression for the gradient, we obtain

$$\begin{aligned} & \nabla_\theta \iint p_\theta^{(t)}(x_t|z)p(z) \log p_d^{(t)}(x_t) dx_t dz \\ &= \mathbb{E}_{z \sim p(z), \epsilon \sim p_\epsilon(\epsilon)} \left[\nabla_{x_t} \log p_d^{(t)}(x_t) \Big|_{x_t = \mathcal{F}_\theta(g_\theta(z), \epsilon, t)} \cdot \frac{\partial \mathcal{F}_\theta(g_\theta(z), \epsilon, t)}{\partial \theta} \right] \end{aligned} \quad (31)$$

$$= \mathbb{E}_{z \sim p(z), x_t \sim p_\theta^{(t)}(x_t|z)} \left[\nabla_{x_t} \log p_d^{(t)}(x_t) \frac{\partial x_t}{\partial \theta} \right] = \iint p_\theta^{(t)}(x_t|z)p(z) \nabla_{x_t} \log p_d^{(t)}(x_t) \frac{\partial x_t}{\partial \theta} dx_t dz, \quad (32)$$

331 as required. \square

332 C A Practically Equivalent Variational Diffusive Upper Bound Loss

333 **Tweedie's formula [10, 34].** For the Gaussian corruption kernel $p_t(x_t | x) = \mathcal{N}(x_t; x, \sigma_t^2 \mathbf{I})$,
334 Tweedie's identity gives

$$\nabla_{x_t} \log p_d^{(t)}(x_t) = \frac{\mu(x_t; t) - x_t}{\sigma_t^2}, \quad \mu(x_t; t) := \mathbb{E}_{p_d(x|x_t)}[x | x_t]. \quad (33)$$

335 We have from Equation (12)

$$\begin{aligned} \mathcal{L}(\phi, \theta) &= - \int_0^1 \omega(t) p_\theta^{(t)}(x_t, z) \left[x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} + \log q_\phi^{(t)}(z|x_t; t) \right] dt + \text{const} \\ &= - \int_0^1 \omega(t) p_\theta^{(t)}(x_t, z) \left[g_\theta(z)^\top \left[\frac{\mu(x_t; t) - g_\theta(z)}{\sigma_t^2} \right]_{\text{sg}} + \log q_\phi^{(t)}(z|x_t; t) \right] dt + \text{const}. \end{aligned} \quad (34)$$

336 *Proof.* We are given the VarDiU objective:

$$\mathcal{L}(\phi, \theta) = - \int_0^1 \omega(t) p_\theta^{(t)}(x_t, z) \left[x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} + \log q_\phi^{(t)}(z | x_t; t) \right] dt + \text{const.} \quad (35)$$

337 With the stop-gradient applied to the score term,

$$\begin{aligned} x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} &= \frac{1}{\sigma_t^2} x_t^\top [\mu(x_t; t) - x_t]_{\text{sg}} \\ &= \frac{1}{\sigma_t^2} \left\{ x^\top [\mu(x_t; t) - x]_{\text{sg}} + (x_t - x)^\top [\mu(x_t; t) - x_t]_{\text{sg}} \right\}. \end{aligned} \quad (36)$$

338 Write $x_t = x + \sigma_t \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \mathbf{I})$, and note that $[\mu(x_t; t) - x_t]_{\text{sg}}$ is treated as a constant w.r.t.
339 (ϕ, θ) . Taking the expectation over $x_t | x$ of (36) yields

$$\begin{aligned} \mathbb{E}_{x_t|x} \left[x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} \right] \\ = \frac{1}{\sigma_t^2} \mathbb{E}_{x_t|x} \left[x^\top [\mu(x_t; t) - x]_{\text{sg}} \right] + \frac{1}{\sigma_t^2} \mathbb{E}_{x_t|x} \left[(x_t - x)^\top [\mu(x_t; t) - x_t]_{\text{sg}} \right]. \end{aligned} \quad (37)$$

340 The first term on the right is exactly $x^\top [(\mu(x_t; t) - x)/\sigma_t^2]_{\text{sg}}$. The second term decomposes as

$$\mathbb{E}_{x_t|x} \left[(x_t - x)^\top [\mu(x_t; t)]_{\text{sg}} \right] - \mathbb{E}_{x_t|x} [\|x_t - x\|^2].$$

341 Because $(x_t - x) = \sigma_t \varepsilon$ with ε independent of (ϕ, θ) and $[\mu(\cdot; t)]_{\text{sg}}$ blocks all gradients through its
342 argument, this entire bracket has *zero gradient* w.r.t. (ϕ, θ) ; in particular, $\mathbb{E}_{x_t|x} [\|x_t - x\|^2] = d \sigma_t^2$
343 (with d the data dimension) and $\mathbb{E}_{x_t|x} [(x_t - x)^\top [\mu(x_t; t)]_{\text{sg}}]$ does not contribute to parameter
344 gradients under reparametrisation. Hence it can be absorbed into “const”.

345 Combining, we obtain (up to a parameter-independent constant)

$$\mathbb{E}_{x_t|x} \left[x_t^\top [\nabla_{x_t} \log p_d^{(t)}(x_t)]_{\text{sg}} \right] = x^\top \left[\frac{\mu(x_t; t) - x}{\sigma_t^2} \right]_{\text{sg}} + \text{const.} \quad (38)$$

346 Substituting $x = g_\theta(z)$ and keeping the $\log q_\phi^{(t)}(z | x_t; t)$ term unchanged gives

$$\mathcal{L}(\phi, \theta) = - \int_0^1 \omega(t) p_\theta^{(t)}(x_t, z) \left[g_\theta(z)^\top \left[\frac{\mu(x_t; t) - g_\theta(z)}{\sigma_t^2} \right]_{\text{sg}} + \log q_\phi^{(t)}(z | x_t; t) \right] dt + \text{const.}, \quad (39)$$

347 as required. \square

348 **D Joint Entropy is Constant**

349 **Proposition D.1.** *Let*

$$p_\theta^{(t)}(x, z) = p(z) p_\theta^{(t)}(x|z) \quad (40)$$

350 *with a fixed prior $p(z)$ and conditional $\mathcal{N}(x; g_\theta(z), \sigma_t^2 I)$. Then the joint entropy $\mathbb{H}(p_\theta^{(t)}(x_t, z))$ is*
351 *independent of θ .*

352 *Proof.* First, write the joint entropy as the sum of marginal and conditional entropies:

$$\mathbb{H}(p_\theta^{(t)}(x_t, z)) = - \iint p_\theta^{(t)}(x_t, z) \log p_\theta^{(t)}(x_t, z) dx_t dz = \mathbb{H}(p(z)) + \mathbb{H}(p_\theta^{(t)}(x_t|z)), \quad (41)$$

353 where

$$\mathbb{H}(p(z)) = - \int p(z) \log p(z) dz, \quad \mathbb{H}(p_\theta^{(t)}(x_t|z)) = - \iint p(z) p_\theta^{(t)}(x_t|z) \log p_\theta^{(t)}(x_t|z) dx_t dz. \quad (42)$$

354 By assumption, the prior $p(z)$ does not depend on θ , so $\mathbb{H}(p(z))$ is a fixed number.

355 For each z , $x_t|z \sim \mathcal{N}(g_\theta(z), \sigma_t^2 I)$. The entropy of a Gaussian depends only on its covariance:

$$\mathbb{H}(\mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \log[(2\pi e)^d \det \Sigma]. \quad (43)$$

356 where d is independent of θ , Here $\Sigma = \sigma_t^2 I$, so

$$\mathbb{H}(\mathcal{N}(g_\theta(z), \sigma_t^2 I)) = \frac{d}{2} \log[(2\pi e)\sigma_t^2], \quad (44)$$

357 which is independent of $g_\theta(z)$. Hence

$$\mathbb{H}(p_\theta^{(t)}(x_t|z)) = \int p(z) \mathbb{H}(\mathcal{N}(g_\theta(z), \sigma_t^2 I)) dz = \frac{d}{2} \log[(2\pi e)\sigma_t^2], \quad (45)$$

358 also constant. Since both $\mathbb{H}(p(z))$ and $\mathbb{H}(p_\theta^{(t)}(x_t|z))$ are independent of θ , their sum $\mathbb{H}(p_\theta^{(t)}(x_t, z))$
359 is a constant. \square

360 E Proof of Theorem 3.2

361 **Theorem E.1** (Restatement of Theorem 3.2). *By the chain rule of entropy, we have: $\mathbb{H}(p_\theta^{(t)}(x_t, z)) =$
362 $\mathbb{H}(p_\theta^{(t)}(z|x_t)) + \mathbb{H}(p_\theta^{(t)}(x_t))$. Since $\mathbb{H}(p_\theta^{(t)}(x_t, z))$ is constant with respect to θ , it follows that:*

$$\max_\theta \mathbb{H}(p_\theta^{(t)}(x_t)) = \min_\theta \mathbb{H}(p_\theta^{(t)}(z|x_t)) \leq \min_{\theta, \phi} - \iint p_\theta^{(t)}(x_t, z) \log q_\phi^{(t)}(z|x_t) dz dx_t, \quad (46)$$

363 where $q_\phi^{(t)}(z|x_t)$ is a variational approximation to the true posterior $p_\theta^{(t)}(z|x_t)$.

364 *Proof.* We have

$$\mathbb{H}(p_\theta^{(t)}(z|x_t)) = - \iint \log p_\theta(z|x_t) p_\theta(z|x_t) p_\theta^{(t)}(x_t) dz dx_t \quad (47)$$

$$= - \iint \log \left(\frac{p_\theta(z|x_t)}{q_\phi^{(t)}(z|x_t)} q_\phi^{(t)}(z|x_t) \right) p_\theta^{(t)}(x_t|z) p(z) dz dx_t \quad (48)$$

$$= - \iint \log \frac{p_\theta(z|x_t)}{q_\phi^{(t)}(z|x_t)} p_\theta^{(t)}(x_t|z) p(z) dz dx_t - \iint \log q_\phi^{(t)}(z|x_t) p_\theta^{(t)}(z|x_t) p_\theta^{(t)}(x_t) dz dx_t \quad (49)$$

$$= - \underbrace{\text{KL}(p_\theta(z|x_t) || q_\phi^{(t)}(z|x_t))}_{\geq 0} - \int \log q_\phi^{(t)}(z|x_t) p_\theta(z|x_t) p_\theta^{(t)}(x_t) dz dx_t. \quad (50)$$

Therefore,

$$\mathbb{H}(p_\theta^{(t)}(z|x_t)) \leq - \iint \log q_\phi^{(t)}(z|x_t) p_\theta(z|x_t) p_\theta^{(t)}(x_t) dz dx_t$$

365 and the bound is tight when $\text{KL}(p_\theta(z|x_t) || q_\phi^{(t)}(z|x_t)) = 0$. To minimise $\text{KL}(p_\theta(z|x_t) || q_\phi^{(t)}(z|x_t))$
366 wrt ϕ , we have

$$\begin{aligned} \min_\phi \text{KL}(p_\theta^{(t)}(z|x_t) || q_\phi^{(t)}(z|x_t)) &= \min_\phi \iint \log q_\phi(z|x_t) p_\theta^{(t)}(x_t|z) p(z) dz dx_t. \\ \Rightarrow \max_\theta \mathbb{H}(p_\theta^{(t)}(x_t)) &= \min_\theta \mathbb{H}(p_\theta^{(t)}(z|x_t)) \leq \min_\theta \min_\phi - \iint \log q_\phi^{(t)}(z|x_t) p(z) p_\theta^{(t)}(x_t|z) dz dx_t \end{aligned} \quad (51)$$

367 where the last line since Proposition D.1

$$\underbrace{\mathbb{H}(p_\theta^{(t)}(x_t, z))}_{\text{const.}} = \mathbb{H}(p_\theta^{(t)}(x_t)) + \mathbb{H}(p_\theta^{(t)}(z|x_t)) \Rightarrow \mathbb{H}(p_\theta^{(t)}(x_t)) = -\mathbb{H}(p_\theta^{(t)}(z|x_t)) + \text{const.}$$

368 \square

369 F Experimental Details

370 In this section, we present the details of the experimental setting.

371 F.1 Implementation Details

372 For the generator g_θ , we use a 5-layer MLP with latent dimension 2, hidden dimension 400 and SiLU
373 activation. We use Adam to train the generator with learning rate 10^{-4} , batch size 1024, and gradient
374 norm clip 10.0.

375 We employ the same 5-layer MLP with hidden dimension 400, and SiLU activation as well to predict
376 the mean and log-variance of the Gaussian posterior. We use Adam to train the generator and posterior
377 network with learning rate 10^{-4} , batch size 1024, and gradient norm clip 10.0.

378 For Diff-instruct [28], we use the same 5-layer MLP with hidden dimension 400, and SiLU activation
379 to parametrise the student score network. We use Adam to train the student score network with
380 learning rate $5 * 10^{-5}$, batch size 1024, and gradient norm clip 10.0.

381 Both methods are trained for 1,000,000 epochs until convergence. For the learned score, we train an
382 EDM [17] on 10,000 observed samples, while the empirical score is estimated from 10,000 observed
383 samples. All experiments are conducted on NVIDIA RTX 3090 GPUs.

384 **Maximum Mean Discrepancy Implementation Details** For evaluation, we employ the Maximum
385 Mean Discrepancy (MMD) [12] with five Gaussian kernels of bandwidths $\{2^{-2}, 2^{-1}, 2^0, 2^1, 2^2\}$ for
386 evaluation across multiple scales. Given two sets of samples $\{x_i\}_{i=1}^n \sim p_d$ and $\{y_j\}_{j=1}^m \sim p_\theta$, the
387 squared MMD is defined as

$$\text{MMD}^2(p_d, p_\theta; k) = \mathbb{E}_{x, x' \sim p_d} [k(x, x')] + \mathbb{E}_{y, y' \sim p_\theta} [k(y, y')] - 2\mathbb{E}_{x \sim p_d, y \sim p_\theta} [k(x, y)], \quad (52)$$

388 where $k(\cdot, \cdot)$ is the positive definite kernel. A smaller MMD indicates closer alignment between
389 the generated and real distributions. For the results in Table 1, we perform 10 independent training
390 runs. Evaluation metrics are recorded every 1,000 epochs, and at each evaluation step, we compute
391 the mean and standard deviation across runs. We use 10,000 samples for MMD evaluation. To
392 summarise performance, we report the average of these statistics over the last 50 evaluations.

393 **VarDiU Posterior Parametrisation** For VarDiU, we parametrise the Gaussian posterior as

$$q_\phi^{(t)}(z|x_t; t) = \mathcal{N}\left(z; \mu_\phi(x_t; t), \sigma_t^2 \text{diag}(\sigma_\phi^2(x_t; t))\right). \quad (53)$$

394 The time-dependent scaling factor σ_t modulates the entropy of the base distribution to match the
395 noise level at diffusion step t . This allows the posterior to capture fine-grained detail at small t , while
396 maintaining flexibility to denoise highly corrupted samples at larger t .

397 For flow-based posterior, we implement Neural Splines Flows [9] using an open-source package
398 `normflows`¹ [45], we choose the flow length to be 4. We use Adam to train the generator and
399 flow-based posterior network with learning rate 10^{-4} , batch size 1024, and gradient norm clip 10.0
400 as well.

401 F.2 Noise Weighting and Scheduling

402 The one-step generation procedure is given by

$$x = \int \delta(x - g_\theta(z)) dx, \quad (54)$$

403 where the latent prior is sampled as $z \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{init}}^2 I)$. We set $\sigma_{\text{init}}^2 = 1$ in all toy data experiments.

404 For the forward diffusion distribution, we set

$$q(x_t|x_0) = \mathcal{N}(x_t; \alpha_t x_0, \sigma_t^2 I) \quad (55)$$

¹<https://github.com/VincentStimper/normalizing-flows>

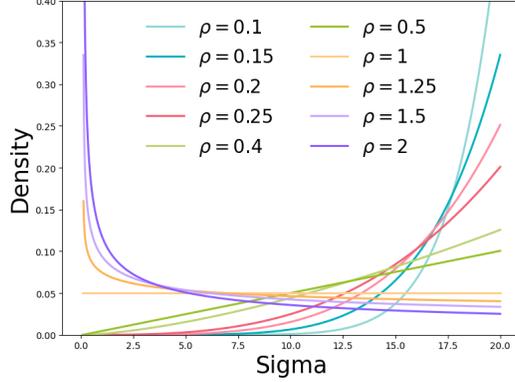


Figure 3: Density of sigma v.s different ρ . Here $\sigma_{\min} = 0.1, \sigma_{\max} = 20$.

405 Unlike Variance Preserving (VP; [42, 15]) schedule which must satisfy $\alpha_t^2 + \sigma_t^2 = 1$, we employ the
 406 schedule used in EDM [17] which is a popular class of diffusion models, where we set $\alpha_t = 1$ for all
 407 t and we sample the time variable $t \sim \text{Unif}[0, 1]$ and define the corresponding noise levels as

$$\sigma_t = \sigma_{\min} + t^\rho (\sigma_{\max} - \sigma_{\min}). \quad (56)$$

408 A visualisation of the resulting density of σ_t under different values of ρ is provided in Figure 3. For
 409 Diff-Instruct, we fix $\rho = 1.5$ across all three score estimation settings. We found an annealing ρ does
 410 not help with the training dynamics of Diff-Instruct, so we use a fixed schedule.

411 **Annealing ρ .** In VarDiU, we adopt an annealing strategy in which ρ gradually increases with
 412 training iterations. Specifically, ρ is initialised at $\rho_{\text{init}} = 0.1$ and incremented by 0.01 every 1000
 413 epochs until it reaches ρ_{end} .

414 The annealed noise schedule plays a crucial role in stabilising training. The central idea is to initially
 415 place more emphasis on larger noise levels (σ), where the learning problem is smoother and less
 416 sensitive to estimation errors. As training progresses, the schedule gradually shifts towards smaller
 417 σ values, allowing the model to refine details and capture higher-frequency structures. This coarse-
 418 to-fine strategy not only mitigates instability at the early stage but also improves convergence and
 419 sample quality in later stages.

420 **(1) Given The True Score.** We set $\sigma_{\min} = 0.1$ and $\sigma_{\max} = 20$ for both Diff-Instruct and VarDiU.
 421 For VarDiU, we set $\rho_{\text{end}} = 5.0$.

422 **(2) Given A Pre-trained DM.** For Diff-Instruct, we set $\sigma_{\min} = 1.1$ and $\sigma_{\max} = 40$. For VarDiU,
 423 we set $\sigma_{\min} = 1.5, \sigma_{\max} = 40$, and $\rho_{\text{end}} = 2.0$.

424 **(3) Given A Dataset.** We employ 10,000 samples for empirical score estimation, with details
 425 provided in Appendix F.3. For both Diff-Instruct and VarDiU, we set $\sigma_{\min} = 0.65$ and $\sigma_{\max} = 40$,
 426 and for VarDiU, we additionally set $\rho_{\text{end}} = 2.0$.

427 **Noise Weighting.** For both Diff-Instruct and VarDiU, we adopt the weighting function

$$\omega(t) = \frac{\sigma_t^2}{\sigma_{\max}^2}.$$

428 The noise schedule aligns with [57, 28]. We further divide the weighting by the σ_{\max} for not changing
 429 the learning rate of training.

430 F.3 Empirical Score Estimation

431 **Score of a Gaussian-smoothed empirical distribution.** Let $\{x_0^{(n)}\}_{n=1}^N \subset \mathbb{R}^d$ be data points, and
 432 define the Gaussian kernel

$$\phi_\sigma(x) := (2\pi\sigma^2)^{-d/2} \exp\left(-\frac{\|x\|^2}{2\sigma^2}\right).$$

433 The σ -smoothed empirical distribution is

$$\hat{p}_\theta(x) := \frac{1}{N} \sum_{n=1}^N \mathcal{N}(x|x_0^{(n)}, \sigma^2 I) = \frac{1}{N} \sum_{n=1}^N \phi_\sigma(x - x_0^{(n)}).$$

434 **Lemma F.1** (Log-sum gradient identity). *For positive functions $\{a_n(x)\}_{n=1}^N$,*

$$\nabla_x \log\left(\sum_{n=1}^N a_n(x)\right) = \sum_{n=1}^N w_n(x) \nabla_x \log a_n(x), \quad w_n(x) := \frac{a_n(x)}{\sum_{m=1}^N a_m(x)}.$$

435 *Proof.* By the chain and quotient rules,

$$\nabla_x \log\left(\sum_n a_n\right) = \frac{\sum_n \nabla_x a_n}{\sum_m a_m} = \sum_n \frac{a_n}{\sum_m a_m} \frac{\nabla_x a_n}{a_n} = \sum_n w_n \nabla_x \log a_n. \quad \square$$

436 Applying the lemma with $a_n(x) = \phi_\sigma(x - x_0^{(n)})$ gives

$$\nabla_x \log \hat{p}_\theta(x) = \sum_{n=1}^N w_n(x; \sigma) \nabla_x \log \phi_\sigma(x - x_0^{(n)}),$$

437 with responsibilities

$$w_n(x; \sigma) = \frac{\exp(-\|x - x_0^{(n)}\|^2 / (2\sigma^2))}{\sum_{m=1}^N \exp(-\|x - x_0^{(m)}\|^2 / (2\sigma^2))}.$$

438 The Gaussian score is

$$\nabla_x \log \phi_\sigma(x - x_0^{(n)}) = -\frac{1}{\sigma^2} (x - x_0^{(n)}) = \frac{x_0^{(n)} - x}{\sigma^2}.$$

439 Therefore,

$$\nabla_x \log \hat{p}_\theta(x) = \sum_{n=1}^N w_n(x; \sigma) \frac{x_0^{(n)} - x}{\sigma^2} = \frac{\sum_{n=1}^N w_n(x; \sigma) x_0^{(n)} - x}{\sigma^2}.$$

440 **Final form.** For any query point x_t , the score of the Gaussian-KDE is

$$\nabla_{x_t} \log \hat{p}_\theta(x_t) = \frac{\bar{x}(x_t; \sigma) - x_t}{\sigma^2}; \quad \bar{x}(x; \sigma) := \sum_{n=1}^N w_n(x; \sigma) x_0^{(n)}. \quad (57)$$

441 We employ Equation (57) to compute the empirical score using the given dataset practically. The
 442 accuracy of empirical score estimation tends to be small when a large number of samples are given,
 443 providing useful teacher guidance.

444 **F.4 Symmetric Sampling**

445 To reduce the variance of the stochastic estimator, we employ a symmetric sampling scheme [6, 20]
 446 for both VarDiU and Diff-Instruct.

447 We first draw a half-batch of latent codes,

$$z^{(i)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\text{init}} \mathbf{I}), \quad i = 1, \dots, \frac{B}{2},$$

448 where B is the batch size. We then duplicate the samples to construct a full batch,

$$z = (z^{(1)}, \dots, z^{(B/2)}, z^{(1)}, \dots, z^{(B/2)}).$$

449 Passing these latents through the generator yields samples

$$x_0 = g_\theta(z).$$

450 Next, we simulate forward diffusion process

$$\epsilon^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad i = 1, \dots, \frac{B}{2},$$

451 and extend them symmetrically as

$$\epsilon = (\epsilon^{(1)}, \dots, \epsilon^{(B/2)}, -\epsilon^{(1)}, \dots, -\epsilon^{(B/2)}).$$

452 Thus, each noise vector is paired with its negative, ensuring that the perturbations are exactly centred.

453 We also sample noise levels from the power distribution

$$\sigma^{(i)} \sim p_\rho(\sigma), \quad i = 1, \dots, \frac{B}{2},$$

454 and replicate them symmetrically:

$$\sigma = (\sigma^{(1)}, \dots, \sigma^{(B/2)}, \sigma^{(1)}, \dots, \sigma^{(B/2)}).$$

455 **Noisy inputs.** The perturbed inputs are then constructed as

$$x_t = x_0 + \sigma \odot \epsilon,$$

456 where \odot denotes element-wise multiplication.

457 We found symmetric sampling can reduce the variance during training, further enhance the training

458 stability of VarDiU, and can help with stabilising the training with Diff-Instruct.

459 **G Related Works**

460 Diffusion acceleration has received significant attention as a means of reducing the number of
461 steps required in the reverse process. At the core of this problem lies the posterior distribution
462 $q(x_0 | x_t)$, which connects a noisy intermediate x_t to the original data x_0 . Accurate modelling of
463 this posterior enables faithful reconstruction of x_0 in fewer steps, thereby lowering the sampling
464 cost. Existing approaches can be broadly grouped into strategies such as designing alternative
465 samplers [25, 26, 39, 7], relaxing the Gaussian posterior assumption [5, 22, 24], or projecting the
466 dynamics onto lower-dimensional subspaces [35]. Whilst these techniques substantially reduce
467 sampling time, they still typically require a non-trivial number of steps to achieve high-quality
468 generation.

469 Another important direction focuses on compressing the reverse diffusion process into shorter chains.
470 This includes progressive distillation [36, 29], consistency models and their improvements [41, 40, 18],
471 as well as extensions to latent diffusion [27, 35]. More recently, Inductive Moment Matching has
472 been proposed as a single-stage framework that trains few-step models from scratch by matching
473 distributional moments [59].

474 In parallel, several works have investigated single-step diffusion models. Some fine-tune pre-trained
475 diffusion models by minimising divergences such as the Diffusive KL or Fisher divergence [28, 51,
476 61, 52], while others incorporate adversarial losses to enhance perceptual fidelity [53, 60]. Our work
477 builds on this line by eliminating the need for a separate student score network, thereby avoiding bias
478 from inaccurate score estimation and enabling more stable training with improved generation quality.

479 Finally, recent work has shown that teacher guidance is not strictly necessary for effective diffusion
480 distillation, highlighting that the multi-scale representations learned by diffusion models are the
481 crucial factor behind the success of one-step generators [57].

482 H More Convergence and Efficiency Results

483 We present the log-MMD plots with respect to both generator gradient iterations and total training
484 time, using the pre-trained DM and the given dataset in Figures 5 and 6, respectively.

485 When relying on the pre-trained DM, both methods exhibit high variance and unstable curves due to
486 inaccurate score predictions from the pre-trained model. The noticeable jump of the VarDiU variants
487 in the middle of training arises from the annealing schedule, where the density concentrates on a
488 range of σ values for which the pre-trained DM suffers from large estimation errors. Nevertheless,
489 the VarDiU variants maintain relatively low variance in the log-density curves, highlighting their
490 robustness.

491 In contrast, when using the given dataset, we estimate the score empirically. The empirical score
492 becomes increasingly accurate with larger sample sizes. As shown in Figure 6, the VarDiU variants
493 achieve performance comparable to that obtained with the true analytical score. Furthermore, VarDiU
494 demonstrates substantially improved training stability and efficiency when employing a Gaussian-
495 parameterised posterior.

496 Overall, these results highlight the superiority of VarDiU, which consistently achieves more stable
and efficient convergence across different settings compared to the baseline.

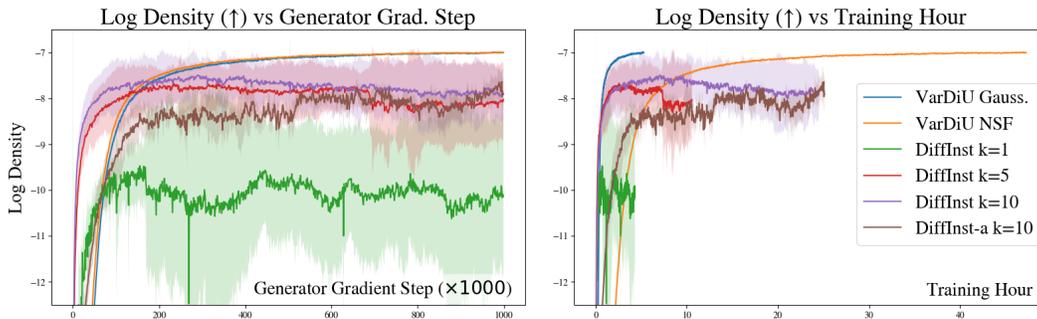
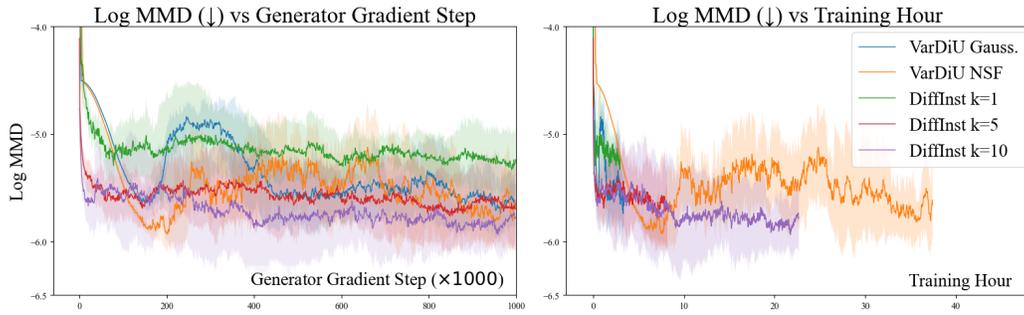
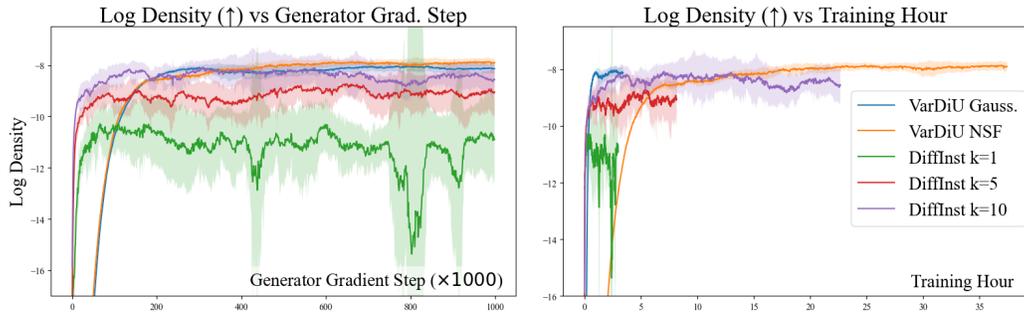


Figure 4: Log-density trajectory with true score setting.

497

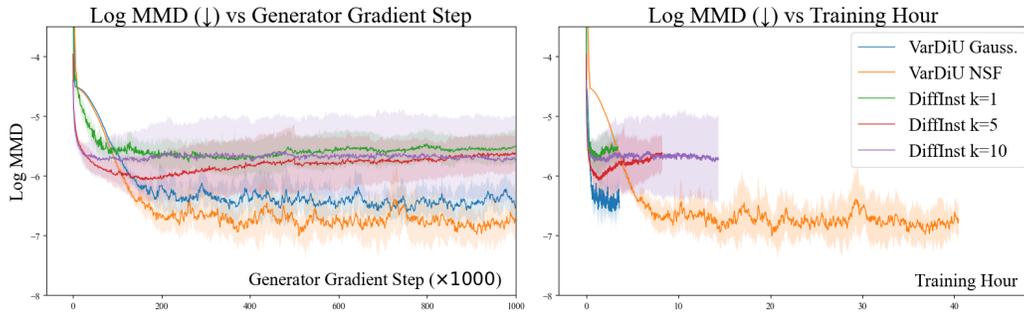


(a) Log-MMD trajectory with pre-trained DM.

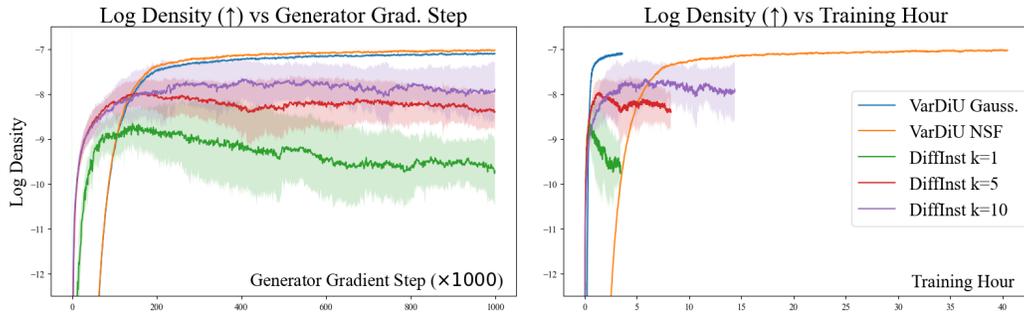


(b) Log-density trajectory with pre-trained DM..

Figure 5: Comparison of samples log-density and log-MMD trajectories with learned score.



(a) Log-MMD trajectory with a given dataset.



(b) Log-density trajectory with a given dataset.

Figure 6: Comparison of samples log-density and log-MMD trajectories under setting of given dataset.