Your Absorbing Discrete Diffusion Secretly Models the Conditional Distributions of Clean Data

Anonymous Author(s) Affiliation Address email

Abstract

Discrete diffusion models with absorbing processes have shown promise in lan-1 guage modeling. The key quantities to be estimated are the ratios between the 2 marginal probabilities of two transitive states at all timesteps, called the concrete З score. In this paper, we reveal that the concrete score in absorbing diffusion can be 4 expressed as conditional probabilities of clean data, multiplied by a time-dependent 5 scalar in an analytic form. Motivated by the finding, we propose reparameterized 6 absorbing discrete diffusion (RADD), a dedicated diffusion model that character-7 izes the time-independent conditional probabilities. Besides its simplicity, RADD 8 can reduce the number of function evaluations (NFEs) by caching the output of 9 the time-independent network when the noisy sample remains unchanged in a 10 sampling interval. Empirically, RADD is up to 3.5 times faster while consistently 11 achieving a better performance than the strongest baseline. Built upon the new 12 13 factorization of the concrete score, we further prove a surprising result that the exact likelihood of absorbing diffusion can be rewritten to a simple form (named 14 15 denoise cross-entropy) and then estimated efficiently by the Monte Carlo method. The resulting approach also applies to the original parameterization of the concrete 16 score. It significantly advances the state-of-the-art discrete diffusion on 5 zero-shot 17 language modeling benchmarks (measured by perplexity) at the GPT-2 scale. 18

19 1 Introduction

Auto-regressive models [1, 2, 3] have dominated the area of language modeling for many years. In particular, such models significantly benefit from large-scale transformers [4] and training data and have achieved remarkable progress [5, 6, 7, 8]. From a probabilistic perspective, the sequential sampling process of auto-regressive models is inefficient and limits the reasoning ability in nonsequential orders [9, 10]. Intrinsically, this is because such models characterize the joint distribution by the chain rule of probability, motivating research on developing other types of generative models for text.

Diffusion models [11, 12, 13] generate data in a coarse-to-fine manner efficiently [14, 15, 16, 17, 18]
and all dimensions simultaneously, providing an appealing alternative to auto-regressive models.
Among other efforts [19, 20, 21, 22, 23, 24, 20, 25, 26, 27, 28, 29] (see Section 5 for a comprehensive discussion), score entropy discrete diffusion (SEDD) [29] has shown promise in text generation. In
particular, SEDD has achieved comparable results to auto-regressive models on 5 zero-shot language modeling benchmarks at the GPT-2 scale. Meanwhile, SEDD can reduce the number of function
evaluations (NFEs) in sampling and fulfill text conditioned on prompts at different positions.

Technically, SEDD employs a discrete-state (absorbing) Markov process that adds noises to data by randomly replacing a token with a mask token [M] and then learns a reverse process to denoise from an entirely masked sentence. The key quantities to be estimated in SEDD are the ratios between the marginal probabilities of two transitive states at all timesteps, called the **concrete score**. SEDD also

Submitted to 38th Conference on Neural Information Processing Systems (NeurIPS 2024). Do not distribute.

proposes a "scaling trick" (see details in Section 3) that scales the output of the score estimation by a factor. The trick has been proven very effective in practice yet not fully understood in theory [29].

One of our main contributions is to reveal that the concrete score in absorbing diffusion can be ex-39 pressed as conditional probabilities of clean data, multiplied by a time-dependent scalar in an analytic 40 form. Our finding theoretically explains the benefits of the scaling trick as a reparameterization 41 for better optimization. Motivated by the finding, we propose reparameterized absorbing discrete 42 diffusion (RADD), a dedicated diffusion model that characterizes the time-independent conditional 43 probabilities by removing the time embedding from the score estimation in SEDD. Besides its 44 simplicity, RADD can significantly reduce the NFEs by caching the output of the time-independent 45 network when the noisy sample remains unchanged in a sampling interval (see Fig. 1). 46

Built upon the new factorization of the concrete score, we further prove a surprising result that the 47 exact likelihood of absorbing diffusion can be rewritten to a simple form (named denoise cross-48 entropy, DCE) and then estimated efficiently by the Monte Carlo method. To establish the theory, 49 we apply a change of variable from the time t to the probability that a single-dimensional token is 50 masked at time t in the forward process. By integrating the probability variable analytically, we show 51 that DCE enumerates all orders to decompose the joint distribution auto-regressively and accumulates 52 log densities of all conditional distributions in every order, finishing the proof. Such theoretical 53 findings enable exact likelihood evaluation and optimization for both the original parameterization of 54 absorbing diffusion [29] and the proposed RADD. 55

Empirically, RADD is up to 3.5 times faster while consistently achieving a better performance than
 the strongest baseline, i.e. SEDD with the scaling trick [29]. Further, the DCE loss applies to both
 RADD and SEDD for precise likelihood evaluation. It significantly advances the state-of-the-art
 discrete diffusion (i.e. SEDD [29]) on 5 zero-shot language modeling benchmarks (measured by
 perplexity) at the GPT-2 scale. The empirical evidence validates our theoretical findings.

61 In summary, this paper has several contributions:

- Deeper understanding of discrete diffusion: Both the factorization form of the concrete
 score and DCE loss for the exact likelihood computation reveal important yet overlooked
 theoretical properties of absorbing discrete diffusion, which explain the mysterious scaling
 trick, provide practice guidance, and may inspire future work.
- **Simplification**: By removing the time conditions, we reparameterize the model to focus on a time-independent conditional probability, simplifying the existing model.
- Efficient sampling: Leveraging the reparameterized form, RADD with a caching strategy is consistently faster while achieving a better performance than the strongest competitor.
- Improved likelihood evaluation: The exact likelihood evaluation approach significantly advances the state-of-the-art discrete diffusion on 5 zero-shot language modeling benchmarks (measured by perplexity) at the GPT-2 scale.

73 2 Background

In this section, we present preliminaries on continuous-time discrete diffusion models. We start with the one-dimensional case in Section 2.1, followed by the multi-dimensional case in Section 2.2.

76 2.1 Single dimension

⁷⁷ Let x denote a single dimensional sample with possible values in $\{1, ..., N\}$. A continuous-time ⁷⁸ discrete Markov chain at time t is characterized by a transition rate matrix Q_t as follows

$$p_{t+\Delta t|t}(\hat{x}|x) = \begin{cases} \boldsymbol{Q}_t(x,\hat{x})\Delta t + o(\Delta t), & \hat{x} \neq x, \\ 1 + \boldsymbol{Q}_t(x,x)\Delta t + o(\Delta t), & \hat{x} = x, \end{cases}$$
(2.1)

where $Q_t(x, \hat{x})$ is the (x, \hat{x}) element of transition rate matrix Q_t , denoting the transition rate from state x to state \hat{x} at time t. Equivalently, we can directly define $Q_t(x, \hat{x})$ as

$$\boldsymbol{Q}_{t}(x,\hat{x}) = \begin{cases} \lim_{\Delta t \to 0} \frac{p_{t+\Delta t|t}(\hat{x}|x)}{\Delta t}, & \hat{x} \neq x, \\ \lim_{\Delta t \to 0} \frac{p_{t+\Delta t|t}(x|x)-1}{\Delta t}, & \hat{x} = x. \end{cases}$$
(2.2)

Given the above definition, denote $P_{s \to t}(x, \hat{x}) := p_{t|s}(\hat{x}|x)$. The following Kolmogorov's forward 81

equation holds [26, 30]: 82

$$\frac{d}{dt}\boldsymbol{P}_{s\to t} = \boldsymbol{P}_{s\to t}\boldsymbol{Q}_t.$$
(2.3)

In practice [26, 29], Q_t is parameterized as $\sigma(t)Q$, where $\sigma(t)$ is a scalar function and Q is a constant matrix. In this case, the solution to Eq. (2.3) can be solved analytically as $P_{s\to t} = \exp((\bar{\sigma}(t) - \bar{\sigma}(s))Q)$, where $\bar{\sigma}(t) = \int_0^t \sigma(s)ds$ and exp is the matrix exponential. Therefore, we can directly sample x_t from x_s in one step for any t > s. 83 84 85 86

Further, Q is often designed to diffuse towards a uniform distribution or an absorbing state [M]. 87 Recent work [20, 26] suggests that the absorbing matrix achieves better empirical performance. 88 Besides, as detailed in Section 3, the specific structure of the absorbing matrix can be leveraged 89 to improve performance and accelerate sampling. Therefore, we focus on the absorbing matrix as 90 follows: 91

$$\boldsymbol{Q}^{\text{absorb}} = \begin{bmatrix} -1 & 0 & \cdots & 0 & 1\\ 0 & -1 & \cdots & 0 & 1\\ \vdots & \vdots & \ddots & \vdots & \vdots\\ 0 & 0 & \cdots & -1 & 1\\ 0 & 0 & \cdots & 0 & 0 \end{bmatrix}.$$
 (2.4)

- The time reversal of the forward process is characterized by a reverse transition rate matrix \tilde{Q}_t [31, 32], 92
- whose element from state x to state \hat{x} is given by 93

$$\tilde{Q}_{t}(x,\hat{x}) = \begin{cases} \frac{p_{t}(\hat{x})}{p_{t}(x)} Q_{t}(\hat{x}, x), & \hat{x} \neq x, \\ -\sum_{k \neq x} \tilde{Q}_{t}(x, k), & \hat{x} = x. \end{cases}$$
(2.5)

- 94
- Simulating the reverse process requires to learn the reverse transition rate $\tilde{Q}_t(x, \hat{x})$. As $Q_t(x_t, \hat{x}_t)$ is known, it is sufficient to estimate the concrete score $\frac{p_t(\hat{x}_t)}{p_t(x_t)}$ by a score network $s_\theta(x_t, t) \approx$ 95
- $[\frac{p_t(\hat{x}_t)}{p_t(x_t)}]_{\hat{x}_t \in \mathcal{X}}$ [28]. Denoising score entropy (DSE) [29] is an effective objective to train the score 96 97 network

$$\int_{0}^{T} \mathbb{E}_{\tilde{x} \sim p_{t|0}(\cdot|x_{0})} \sum_{y \neq \tilde{x}} \boldsymbol{Q}_{t}\left(\tilde{x}, y\right) \left(s_{\theta}\left(\tilde{x}, t\right)_{y} - \frac{p_{t|0}\left(y \mid x_{0}\right)}{p_{t|0}\left(\tilde{x} \mid x_{0}\right)} \log s_{\theta}\left(\tilde{x}, t\right)_{y} + K\left(\frac{p_{t|0}\left(y \mid x_{0}\right)}{p_{t|0}\left(\tilde{x} \mid x_{0}\right)}\right)\right) dt,$$
(2.6)

where $K(a) := a \log a - a$. In particular, the DSE loss in Eq. (2.6) is an evidence lower bound 98 (ELBO) of the negative log-likelihood with an unknown gap. Nevertheless, existing work [29] still 99 employs it for training and likelihood evaluation. 100

After training, sampling from the model can be understood as discretizing the following process 101

$$\frac{d}{dt}\boldsymbol{P}_{s\to t} = \boldsymbol{P}_{s\to t}\tilde{\boldsymbol{Q}}_t, \qquad (2.7)$$

where dt is an infinitesimal negative timestep and the concrete score is replaced by the score network. 102 Existing samplers include the Euler method, Gillespie method, and Tweedie τ -leaping, as detailed in 103 Appendix D. 104

2.2 Multi-dimension 105

The multi-dimensional cases consider a state space of size d like $\mathcal{X}^d = \{1, \ldots, n\}^d$. We denote the sample as a sequence of one-dimensional data, i.e. $\boldsymbol{x} = x^1 \dots x^d$. The transition matrix 106 107 $Q_t \in \mathbb{R}^{n^d \times n^d}$ has an exponential number of possible states, making it expensive to reverse. To 108 alleviate this issue, existing work [26, 29] assumes independence between dimensions and each 109 dimension is a one-dimensional diffusion process with the same transition rate matrix $Q_t^{\text{tok}} \in \mathbb{R}^{n \times n}$. 110

Under the independent assumption, Q_t assigns zero values [26, 29] for all sequences with a Hamming distance larger than 1. According to Eq. (2.4), it is sufficient to model the concrete score between sequences that differ by a Hamming distance of 1, such as $\hat{x}_t = x_t^1 \dots \hat{x}_t^i \dots x_t^d$ given $x_t = x_t^1 \dots x_t^d$. Therefore, the score network $s_{\theta}(\cdot, t) : \{1, \dots, n\}^d \to \mathbb{R}^{d \times n}$ is defined as

$$\boldsymbol{s}_{\theta}\left(\boldsymbol{x}_{t},t\right)_{\boldsymbol{\hat{x}}_{t}}=\boldsymbol{s}_{\theta}\left(\boldsymbol{x}_{t}^{1}\ldots\boldsymbol{x}_{t}^{i}\ldots\boldsymbol{x}_{t}^{d},t\right)\left[\boldsymbol{i},\widehat{\boldsymbol{x}}_{t}^{i}\right]\approx\frac{p_{t}\left(\boldsymbol{x}_{t}^{1}\ldots\widehat{\boldsymbol{x}}_{t}^{i}\ldots\boldsymbol{x}_{t}^{d}\right)}{p_{t}\left(\boldsymbol{x}_{t}^{1}\ldots\boldsymbol{x}_{t}^{i}\ldots\boldsymbol{x}_{t}^{d}\right)},$$

which leads to the following expression to estimate the reverse transition rate matrix \hat{Q}_t : 111

$$\tilde{\boldsymbol{Q}}_t\left(\boldsymbol{x}_t^1 \dots \boldsymbol{x}_t^i \dots \boldsymbol{x}_t^d, \boldsymbol{x}_t^1 \dots \widehat{\boldsymbol{x}}_t^i \dots \boldsymbol{x}_t^d\right) = \boldsymbol{Q}_t^{\text{tok}}\left(\widehat{\boldsymbol{x}}_t^i, \boldsymbol{x}_t^i\right) \frac{p_t\left(\boldsymbol{x}_t^1 \dots \widehat{\boldsymbol{x}}_t^i \dots \boldsymbol{x}_t^d\right)}{p_t\left(\boldsymbol{x}_t^1 \dots \boldsymbol{x}_t^i \dots \boldsymbol{x}_t^d\right)}$$
(2.8)

$$\approx \boldsymbol{Q}_t^{\text{tok}}\left(\widehat{x}_t^i, x_t^i\right) \boldsymbol{s}_{\theta}\left(x_t^1 \dots x_t^i \dots x_t^d, t\right) [i, \widehat{x}_t^i].$$
(2.9)

Existing samplers assume that each dimension is independent within a small interval Δt and update 112 each dimension in parallel for efficiency [29, 26]. 113

Reparameterized absorbing discrete diffusion 3 114

In Section 3.1, we reveal that the concrete score of absorbing discrete diffusion can be reparameterized 115 as conditional distributions of clean data, which enables efficient sampling by caching the output of 116 time-independent network (see Section 3.2) and exact likelihood computation (see Section 3.3) by 117 applying the change of variable from time to the probability of being masked in a single dimension. 118

3.1 Parameterizing the concrete score as conditional distributions of clean data 119

A key observation is that only the transition from the masked token to an unmasked token is valid in 120 the reverse process of an absorbing discrete diffusion. In particular, according to the definition of 121 the transition matrix of the absorbing process (see Eq. (2.4)), we have $Q^{absorb}(\hat{x}_t^i, x_t^i) = 0$ for any 122 unmasked $x_t^i \neq [\mathbf{M}]$ and $\hat{x}_t^i \neq x_t^i$. Therefore, the corresponding element in the transition matrix of 123 the reverse process \tilde{Q}_t (see Eq. (2.5)) equals zero. Namely, 124

$$\tilde{\boldsymbol{Q}}_t\left(x_t^1 \dots x_t^i \dots x_t^d, x_t^1 \dots \widehat{x}_t^i \dots x_t^d\right) = \sigma(t) \boldsymbol{Q}^{\text{absorb}}\left(\widehat{x}_t^i, x_t^i\right) \frac{p_t\left(x_t^1 \dots \widehat{x}_t^i \dots x_t^d\right)}{p_t\left(x_t^1 \dots x_t^i \dots x_t^d\right)} = 0, \quad (3.1)$$

for any unmasked state $x_t^i \neq [\mathbf{M}]$ and $\hat{x}_t^i \neq x_t^i$ and it is unnecessary to model the corresponding concrete score $\frac{p_t(x_t^1...\hat{x}_t^i...x_t^d)}{p_t(x_t^1...x_t^i...x_t^d)}$. Also, note that the concrete score always takes the value of one if $\hat{x}_t^i = x_t^i$. Therefore, we only need to characterize the concrete score for $x_t^i = [\mathbf{M}]$ and $\hat{x}_t^i \neq [\mathbf{M}]$. 125 126

127

Interestingly, in this case, we discover that the concrete score has a simple analytic form w.r.t. to the 128 conditional distributions of clean data, as summarized in the following Theorem 1. 129

Theorem 1. (Analytic concrete score in absorbing case, proof in Appendix B) For $\mathbf{x}_t = x_t^1 \dots x_t^i \dots x_t^d$ and $\hat{\mathbf{x}}_t = x_t^1 \dots \hat{\mathbf{x}}_t^i \dots x_t^d$, if $x_t^i = [\mathbf{M}]$ and $\hat{x}_t^i \neq [\mathbf{M}]$, the concrete score at time t can be expressed as a time-independent conditional distribution at time zero multiplied by an 130 131 132 analytic time-dependent term: 133

$$\frac{p_t\left(x_t^1\dots\widehat{x}_t^i\dots x_t^d\right)}{p_t\left(x_t^1\dots x_t^i\dots x_t^d\right)} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} p_0(\widehat{x}_t^i | \boldsymbol{x}_t^{U\!M}),$$

where x_{t}^{UM} is the vector consists of all unmasked tokens of x_{t} . 134

One immediate implication of Theorem 1 is to theoretically explain the benefit of the "scaling 135 trick" in existing work [29] (see Appendix C.2 therein), which significantly improves the practical 136

performance of discrete diffusion (see Table 2) but has not been fully understood. 137

In particular, the scaling trick divides the output of the score network s_{θ} by a factor of $e^{\bar{\sigma}(t)} - 1$. 138 Equivalently, it reparameterizes $s_{\theta}(x_t, t)$ as: 139

$$\boldsymbol{s}_{\theta}(\boldsymbol{x}_{t},t) = \frac{1}{e^{\bar{\sigma}(t)} - 1} \tilde{\boldsymbol{s}}_{\theta}(\boldsymbol{x}_{t},t) = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \tilde{\boldsymbol{s}}_{\theta}(\boldsymbol{x}_{t},t),$$





Figure 1: Expected number of function evaluations (E-NFE) over a different number of sampling steps. E-NFE is measured by Tweedie τ -leaping method with log-linear noise schedule.

Figure 2: Sample quality measured by perplexity (\downarrow). We compare SEDD with Euler and Tweedie τ -leaping (abbr. T- τ) samplers, and RADD with Euler sampler. We show E-NFE for RADD with caching and NEF otherwise.

where the scaling factor coincides with the time-dependent term in Theorem 1. In the original parameterization, the score network s_{θ} must model the whole time-dependent concrete score. In contrast, with the scaling trick, the reparameterized score $\tilde{s}_{\theta}(\boldsymbol{x}_t, t)$ can focus on capturing the clean data distribution $p_0(\hat{x}^i | \boldsymbol{x}_t^{\text{UM}})$ and simplifies learning, according to Theorem 1.

Further, Theorem 1 suggests that it is unnecessary to incorporate the time t in the reparameterized score, and the reparameterized score $\tilde{s}_{\theta}(x_t, t)$ should output a valid probability distribution. Motivated by the insights, we propose reparameterized absorbing discrete diffusion (RADD), which

147 employs a network $c_{\theta}(x_t)$ that removes the time condition from the input and takes the softmax as

final nonlinearity. Formally, we can write our reparameterization as:

$$\boldsymbol{s}_{\theta}(\boldsymbol{x}_{t},t) = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \boldsymbol{c}_{\theta}(\boldsymbol{x}_{t}).$$
(3.2)

In practice, we make a minimal modification of the score network in SEDD [29] for simplicity andfairness, detailed in Appendix F.1.

Moreover, RADD also enjoys a more efficient sampling process than SEDD [29] (with or without the scaling trick) based on its simplified parameterization, as presented below.

153 **3.2** Efficient samplers to reduce NFE by caching $c_{\theta}(x_t)$

For the reverse process of an absorbing discrete diffusion, once a token is generated from [M] to 154 an unmasked token, it never transits to another token. Therefore, for a sequence x_t of length d, x_t 155 changes at most d times, irrespective of the number of sampling steps D. In the other steps, x_t 156 157 remains in all d dimensions. We highlight that we can cache $c_{\theta}(x_t)$ naturally without evaluating the time-independent c_{θ} to reduce the NFE compared to SEDD (see Appendix E for the pseudo-code,). 158 As shown in Fig. 2, RADD with the caching strategy is more efficient than SEDD given any number 159 of sampling steps, especially given large sampling steps. This is as expected because the NFE is 160 limited within the generating sequence length. 161

Note that the NFEs with the caching strategy is a random variable. To quantify it, we calculate the expected NFEs (abbr. E-NFEs) required in an analytic form, conditioned on the sampling method, time steps, and noise schedule. Specifically, denote l as the generating sequence length, which does not equal d generally. Given the sampling time steps $\{t_0 = 0, \dots, t_n = T\}$, let $N_k \in \{0, \dots, l\}$ denote the number of changed dimensions of x in $[t_{k-1}, t_k)$. Since we perform function evaluation 167 in $[t_{k-1}, t_k)$ only when x changes (i.e. $N_k \neq 0$), the NFEs and E-NFEs can expressed as:

$$NFEs(n) = \sum_{k=1}^{n} \mathbb{I}(N_k \neq 0), \qquad (3.3)$$

$$\text{E-NFEs}(n) = \sum_{k=1}^{n} \mathbb{E}[\mathbb{I}(N_k \neq 0)] = \sum_{k=1}^{n} P(N_k \neq 0).$$
(3.4)

For each dimension *i*, let r_k represent the probability that x^i changes within the interval $[t_{k-1}, t_k)$. As the probability is independent in different dimensions (proof in Appendix D.3), N_k follows a

binomial distribution with parameters l and r_k . Therefore, Eq. (3.4) can be further simplified as:

E-NFEs
$$(n) = \sum_{k=1}^{n} P(N_k \neq 0) = \sum_{k=1}^{n} (1 - (1 - r_k)^l),$$
 (3.5)

which applies to all samplers. Further, r_k can be analytically expressed w.r.t. the time steps and

¹⁷² noise schedule for both Euler and Tweedie τ -leaping samplers, as detailed in Appendix D.3. Taking ¹⁷³ Tweedie τ -leaping method with log-linear noise schedule [29] for example, its E-NFEs is given by:

$$\text{E-NFEs}(n) = \sum_{k=1}^{n} (1 - (1 - \frac{1}{n})^{l}) = n(1 - (1 - \frac{1}{n})^{l}).$$
(3.6)

Appendix D.3 provides the proof. As shown in Fig. 1, we plot the curve of Eq. (3.6) in blue, which agrees with our experiments (the red stars).

176 3.3 Denoise cross-entropy for exact likelihood evaluation and training

As illustrated in Theorem 1, the concrete score can be understood as a rescaled conditional distribution on clean data. From this perspective, it is natural to wonder: **is it possible to evaluate and optimize the exact likelihood of the model instead of the ELBO?** Surprisingly, the answer is yes for both the original parameterization [29] and our new parameterization.

Let $q_{\theta}(\boldsymbol{x}_0)$ denote the model distribution at time zero defined by \boldsymbol{s}_{θ} , or our \boldsymbol{c}_{θ} , which approximates the true distribution $p_0(\boldsymbol{x}_0)$. Inspired by the cross-entropy loss in auto-regressive models, we define the denoising cross-entropy loss $\mathcal{L}_{\text{DCE}}^T(\boldsymbol{x}_0)$ as:

$$\mathcal{L}_{\text{DCE}}^{T}(\boldsymbol{x}_{0}) := \int_{0}^{T} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{t\mid0}(\cdot\mid\boldsymbol{x}_{0})} \sum_{\boldsymbol{y} \neq \tilde{\boldsymbol{x}}} \boldsymbol{Q}_{t}\left(\tilde{\boldsymbol{x}}, \boldsymbol{y}\right) \left(-\frac{p_{t\mid0}\left(\boldsymbol{y}\mid\boldsymbol{x}_{0}\right)}{p_{t\mid0}\left(\tilde{\boldsymbol{x}}\mid\boldsymbol{x}_{0}\right)} \log \boldsymbol{s}_{\theta}\left(\tilde{\boldsymbol{x}}, t\right)_{\boldsymbol{y}}\right) dt,$$
(3.7)
$$= \int_{0}^{T} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{t\mid0}(\cdot\mid\boldsymbol{x}_{0})} \sum_{\boldsymbol{y} \neq \tilde{\boldsymbol{x}}} \boldsymbol{Q}_{t}\left(\tilde{\boldsymbol{x}}, \boldsymbol{y}\right) \left(-\frac{p_{t\mid0}\left(\boldsymbol{y}\mid\boldsymbol{x}_{0}\right)}{p_{t\mid0}\left(\tilde{\boldsymbol{x}}\mid\boldsymbol{x}_{0}\right)} \log \left(\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \boldsymbol{c}_{\theta}(\tilde{\boldsymbol{x}})_{\boldsymbol{y}}\right)\right) dt.$$
(3.8)

Compared with the DSE loss in Eq. (2.6), our DCE loss simply removed the terms $s_{\theta}(\tilde{x}, t)_{y}$ and $K\left(\frac{p_{t|0}(y|x_{0})}{p_{t|0}(\tilde{x}|x_{0})}\right)$, however, it shows that DCE loss exactly equals the negative log-likelihood of $q_{\theta}(x_{0})$ with a sufficiently long process in absorbing discrete diffusion.

Theorem 2. Suppose $\{X_t\}$ is a continuous time Markov chain with transition rate matrix $Q_t =$

¹⁸⁶ **Theorem 2.** Suppose $\{X_t\}$ is a commons time matrix of that with transition the matrix $\mathbf{G}_t = \sigma(t)\mathbf{Q}^{absorb}$. For a given data \mathbf{x}_0 , if $\sigma(t)$ satisfies $\int_0^\infty \sigma(\tau)d\tau = \infty$, then the denoising cross-entropy loss defined in Eq. (3.8) with $T \to \infty$ exactly equals the negative log-likelihood of \mathbf{x}_0 .

$$\mathcal{L}_{DCE}^{\infty}(\boldsymbol{x}_0) = -\log q_{\theta}(\boldsymbol{x}_0). \tag{3.9}$$

The proof of Theorem 2 consists of three key steps, detailed in Appendix C.1, Appendix C.2 and Appendix C.3 respectively. In the first step, we apply a change of variable from t to $\lambda(t) = 1 - e^{-\bar{\sigma}(t)}$, which is the probability of a token is masked from 0 to t in the forward process. Further, inspired by the factorization form discovered in Theorem 1, the denoising cross-entropy loss for both parameterizations can then be rewritten as an integral of λ

$$\mathcal{L}_{\text{DCE}}^{\infty}(\boldsymbol{x}_{0}) = \int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}}|\boldsymbol{x}_{0})} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} -\log q_{\theta}(\boldsymbol{x}_{0}^{i}|\tilde{\boldsymbol{x}}^{\text{UM}}) \right] d\lambda,$$
(3.10)

where $p_{\lambda}(\tilde{x}|x_0)$ is the joint distribution induced by masking each dimension in x_0 independently with a probability λ .

In the second step, we demonstrate that the integral w.r.t. λ in Eq. (3.10) can be integrated analytically,

and the DSE loss can be rewritten as expectations over the number and positions of masks as follows:

$$\mathcal{L}_{\text{DCE}}^{\infty}(\boldsymbol{x}_{0}) = d\mathbb{E}_{k \sim U(\{1, \cdots, d\})} \frac{1}{k} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim U(\tilde{\mathcal{X}}_{k})} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} -\log q_{\theta}(\boldsymbol{x}_{0}^{i} | \tilde{\boldsymbol{x}}^{\text{UM}}) \right], \quad (3.11)$$

where we denote $\tilde{\mathcal{X}}_k := \{ \tilde{\boldsymbol{x}} : \tilde{\boldsymbol{x}} \in \tilde{\mathcal{X}} \text{ and } \tilde{\boldsymbol{x}} \text{ has exact } k \text{ dimensions masked by } [\mathbf{M}] \}$ and $U(\cdot)$ as uniform distribution.

Finally, in the third step, we prove that Eq. (3.11) enumerates all orders to decompose the joint distribution auto-regressively and accumulates log densities of all conditional distributions in every order. Therefore, it is equivalent to the negative log-likelihood of q_{θ} :

$$\mathcal{L}_{\text{DCE}}^{\infty}(\boldsymbol{x}_0) = -\log q_{\theta}(\boldsymbol{x}_0). \tag{3.12}$$

Theorem 2 enables exact likelihood computation for both the original model s_{θ} and our c_{θ} , providing

a more accurate measure of model performance. Take c_{θ} for example, Eq. (3.8) can be rewritten as a form of expectation on t:

$$\mathcal{L}_{\text{DCE}}^{T}(\boldsymbol{x}_{0}) = \frac{1}{T} \mathbb{E}_{t \sim U([0,T])} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{t|0}(\cdot \mid \boldsymbol{x}_{0})} \sum_{\boldsymbol{y} \neq \tilde{\boldsymbol{x}}} \boldsymbol{Q}_{t}\left(\tilde{\boldsymbol{x}}, \boldsymbol{y}\right) \left(-\frac{p_{t|0}\left(\boldsymbol{y} \mid \boldsymbol{x}_{0}\right)}{p_{t|0}\left(\tilde{\boldsymbol{x}} \mid \boldsymbol{x}_{0}\right)} \log\left(\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \boldsymbol{c}_{\theta}(\tilde{\boldsymbol{x}})_{\boldsymbol{y}}\right) \right)$$
(3.13)

Naturally, we can take the Monte Carlo estimation of $\mathcal{L}_{DCE}^{T}(\boldsymbol{x}_{0})$ by sampling *t* to approximate $-\log q_{\theta}(\boldsymbol{x}_{0})$ according to Eq. (3.13). In addition, it can be used as an efficient and valid training target for discrete diffusion models, as an alternative to the ELBO (i.e. DSE loss). For pseudo-code of training, see Appendix E.

211 4 Experiments

We present the experimental setups in Section 4.1. We then evaluate the performance of accelerated generation in Section 4.2 and zero-shot perplexity on various language datasets in Section 4.3.

214 4.1 Settings

Model. We use RADD model c_{θ} reparameterzied as described in Section 3.1. Compared with SEDD small model, RADD model has 7M fewer parameters due to the removal of time-condition, which equates to an 8% decrease from the original 90M non-embedding parameters. We trained our RADD model c_{θ} using denoising score entropy and denoising cross entropy, abbreviated as RADD-DSE and RADD-DCE. For SEDD small model, we employed their pre-trained model.

Data. In line with the methodology outlined by SEDD, we trained on the OpenWebText [33] dataset and tested on the LAMBADA, WikiText2, PTB, WikiText103, and One Billion Words datasets [34, 35, 36]. For data splits and data processing, we adopted the same settings and techniques as SEDD, which involves packing sentences to generate uniform-length blocks as model input.

Training setup. We used the same training setup for RADD and SEDD. Specifically, we used a log-linear noise schedule where the expectation of the number of changed tokens at time t is linear with t. For simplicity, we also used the same optimization configuration as SEDD, which can be suboptimal for our RADD model and DCE loss. For more details see Appendix F.

Metric. Following previous work [29], we conduct experiments on unconditional generation and language modeling tasks. For generation, we use perplexity (PPL) on unconditional samples measured by an additional larger language model (i.e. GPT-2 large) to evaluate sample quality. To access inference efficiency, we computed the inference time on a single NVIDIA 4090 GPU with a batch size of 8 and averaged over 1024 samples. For language modeling tasks, we report the perplexity calculated on the dataset with different models.

Table 1: Avarage inference time of a single sample with varying sampling steps. The table compares the average inference time (in seconds) for the SEDD small model using both Euler and Tweedie τ -leaping (abbreviated as T- τ) sampling methods, and the RADD small model using the Euler method with a caching strategy.

Methods	Metrics	32	64	128	256	512	1024	2048	4096
SEDD (euler)	Time(s)	0.48	0.87	1.67	3.25	6.41	12.74	25.42	50.86
	PPL	155	105	81	66	53	43	35	28
SEDD (T- τ)	Time(s)	0.38	0.68	1.28	2.47	4.85	9.61	19.14	38.20
	PPL	151	104	81	65	52	42	34	28
RADD	Time(s)	0.33	0.54	0.94	1.68	2.97	5.15	8.73	14.88
	PPL	135	94	72	58	46	37	30	26

234 **4.2 Efficient sampling**

We compare the sample quality measured by perplexity between SEDD and our RADD-DCE model, as shown in Fig. 2. For a fixed NFE, RADD-DCE with the Euler sampler outperforms SEDD with multiple samplers. It suggests that RADD with caching accelerates the sampling process and benefits sample quality at the same time. Besides, the acceleration by cache strategy is particularly significant with large sampling steps, as analyzed in Section 3.2.

We further compare the running time for the methods in Table 1. Across all sampling steps, RADD consistently requires the shortest sampling time and outperforms SEDD with different samplers. Quantitatively, RADD achieves a speed-up of $2.5 \sim 3.5$ times as shown in Table 1. These results agree with the analysis of the E-NFEs in Fig. 1, validate the effectiveness of RADD and caching strategy, and demonstrate the practical implications of our Theorem 1.

According to Eq. (3.11), we can also use RADD as an auto-regressive model to generate samples in
different orders, leading to worse performance as a discrete diffusion, as detailed in Appendix F.4.
We present more sampling details in Appendix F.3. and the generated samples in Appendix G.1.

248 4.3 Improved zero-shot perplexity on language modeling

Following SEDD, we present zero-shot perplexities on the LAMBADA, WikiText2, PTB, Wiki-Text103, and 1 Billion Words datasets [37] in Table 2 and compare the zero-shot perplexity of our model with other baseline models [20, 38, 29].

Firstly, we conduct an ablation study of the scaling trick in the middle of the Table 2. With an absorbing process, the perplexity of the scaled version of SEDD outperforms its unscaled version, which matches our theoretical discovery in Theorem 1.

Secondly, without any modification of the model, we estimate the exact likelihood of the baseline model SEDD [29] based on Theorem 2 in Table 2. We observe that perplexity is consistently better than the ELBO of the strongest discrete diffusion models, which validates our Theorem 2.

Lastly, we report the maximum likelihood training results of RADD in the last row in Table 2. We observed that RADD-DCE outperforms RADD-DSE, but their performances are slightly worse than SEDD. This discrepancy could be because we did not search the hyperparameters and directly applied identical optimization configures as SEDD, which may be suboptimal.

262 5 Related work

Continuous-state diffusion models for text generation. Several works have been proposed to apply continuous diffusion to text [19, 21, 22, 23]. Li et al. [19] use an embedding layer to map discrete tokens to a latent space and learn a continuous-state diffusion on it. Bit Diffusion [22] learns a continuous diffusion model to generate binary bits of discrete tokens. However, transforming between these continuous representations and discrete tokens by thresholding may lose information. Bayesian Flow Network [23] achieves competitive log-likelihood on character-level language modeling tasks

Table 2: **Zero-shot language modeling perplexity** (\downarrow) **on five datasets.** [†] labels the results based on ELBO which is taken from [20, 38, 29] and * labels the results based on the exact likelihood implemented by us. In this table, SEDD-U / SEDD-S refer to the unscaled and scaled absorbing models respectively.

Method	LAMBADA	WikiText2	PTB	WikiText103	1BW
GPT-2	45.04	42.43	138.43	41.60	75.20
D3PM [†] PLAID [†] SEDD-Uniform [†]	$\leq 93.47 \\ \leq 57.28 \\ \leq 65.40$	$\leq 77.28 \\ \leq 51.80 \\ \leq 50.27$	≤ 200.82 ≤ 142.60 ≤ 140.12	$\leq 75.16 \\ \leq 50.86 \\ \leq 49.60$	≤ 138.92 ≤ 91.12 ≤ 101.37
SEDD-U [†] SEDD-S [†]	≤52.21 ≤50.92	$\leq 44.75 \\ \leq 41.84$	≤130.49 ≤114.24	$\leq 43.14 \\ \leq 40.62$	$\leq 80.70 \\ \leq 79.29$
SEDD-S* (Ours)	50.44	39.91	110.01	39.91	78.01
RADD-DSE* (Ours) RADD-DCE* (Ours)	96.62 56.67	43.35 42.83	125.03 116.74	40.34 41.02	80.11 79.00

and is proven equivalent to continuous stochastic differential equations trained by denoising score matching [24]. Such models underperform auto-regressive models on standard text generation tasks.

Discrete-state diffusion models for text generation. Several discrete-state diffusion models have 271 been proposed [11, 39, 20]. D3PM [20] proposed a diffusion framework based on any probability 272 transition matrix and trained with a lower bound of log-likelihood. DiffusionBERT [25] utilizes a 273 pre-trained BERT [40] as an initialization of diffusion. Furthermore, [26] generalizes the framework 274 to continuous time by introducing a rate matrix. It is difficult to apply the score matching in such 275 models because the gradient of the data distribution is undefined. Several works try to generalize the 276 score matching on discrete data [29, 28, 26, 27]. Meng et al. [28] introduce the concrete score and the 277 denoising concrete score matching loss. Furthermore, SEDD bridges the discrete state diffusion and 278 the concrete score by introducing a denoising score entropy loss [29]. By incorporating an absorbing 279 280 process, SEDD achieves competitive performance with the auto-regressive models, especially, GPT-2.

281 6 Conclusion

We introduce RADD, a dedicated discrete diffusion model that characterizes the time-independent conditional probabilities, built upon a new factorization form of the concrete score. RADD is much more efficient by reducing the NFEs with a cache strategy while retaining a better performance than strong baselines. Furthermore, we propose DCE loss and prove it is equivalent to the negative log-likelihood of absorbing diffusion. When applied to SEDD, DCE significantly advances the state-of-the-art discrete diffusion on 5 zero-shot language modeling benchmarks at the GPT-2 scale.

Limitaition. Our model has been trained and evaluated primarily on the GPT-2 scale. For broader applicability, it is essential to explore the effects of scaling on the performance [41], which is left as future work. The success of diffusion transformers on images [42, 43, 44] and videos [45] suggests that diffusion models can be scaled up by incorporating transformers.

Another limitation is that our model can only generate full-length outputs, unlike auto-regressive
 models that can produce variable-length outputs. This restricts the flexibility of our model in certain
 applications. We leave the investigation on this issue as future work.

Social impact. For the current theoretical and experimental scope of this paper, we have not found any direct social impacts. However, considering future developments, the paper potentially contributes to the next-generation large language models. In this context, this work could significantly reduce the inference cost of language models but may also lead to hallucinations, amplify biases and discrimination in the data, and pose risks of misuse. As with other generative models, addressing these issues requires further advancements in the field.

301 **References**

- [1] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language
 understanding by generative pre-training. 2018.
- [2] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [5] OpenAI. ChatGPT: Optimizing Language Models for Dialogue. November 2022. URL
 https://openai.com/blog/chatgpt/.
- [6] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni
 Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4
 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [7] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,
 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
 language models, 2023.
- [8] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos,
 Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report.
 arXiv preprint arXiv:2305.10403, 2023.
- [9] Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz
 Korbak, and Owain Evans. The reversal curse: Llms trained on" a is b" fail to learn" b is a".
 arXiv preprint arXiv:2309.12288, 2023.
- [10] Ang Lv, Kaiyi Zhang, Shufang Xie, Quan Tu, Yuhan Chen, Ji-Rong Wen, and Rui Yan. Are we
 falling in a middle-intelligence trap? an analysis and mitigation of the reversal curse. *arXiv preprint arXiv:2311.07468*, 2023.
- [11] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsuper vised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances
 in neural information processing systems, 33:6840–6851, 2020.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and
 Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
- [14] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [15] Fan Bao, Chongxuan Li, Jun Zhu, and Bo Zhang. Analytic-dpm: an analytic estimate of the
 optimal reverse variance in diffusion probabilistic models. *arXiv preprint arXiv:2201.06503*,
 2022.
- [16] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential
 integrator. *arXiv preprint arXiv:2204.13902*, 2022.
- [17] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver:
 A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022.

- [18] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022.
- [19] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B. Hashimoto.
 Diffusion-Im improves controllable text generation, 2022.
- [20] Jacob Austin, Daniel D. Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. Structured denoising diffusion models in discrete state-spaces. In *Advances in Neural Information Processing Systems*, 2021.
- [21] Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin,
 Pierre H. Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, Curtis
 Hawthorne, Rémi Leblond, Will Grathwohl, and Jonas Adler. Continuous diffusion for categor ical data, 2022.
- [22] Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using
 diffusion models with self-conditioning, 2023.
- [23] Alex Graves, Rupesh Kumar Srivastava, Timothy Atkinson, and Faustino Gomez. Bayesian
 flow networks, 2024.
- [24] Kaiwen Xue, Yuhao Zhou, Shen Nie, Xu Min, Xiaolu Zhang, Jun Zhou, and Chongxuan Li.
 Unifying bayesian flow networks and diffusion models through stochastic differential equations,
 2024.
- [25] Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Diffusion bert: Improving generative masked language models with diffusion models. *arXiv preprint arXiv:2211.15029*, 2022.
- [26] Andrew Campbell, Joe Benton, Valentin De Bortoli, Tom Rainforth, George Deligiannidis, and
 A. Doucet. A continuous time framework for discrete denoising models. In *Advances in Neural Information Processing Systems*, 2022.
- [27] Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time
 discrete diffusion models, 2023.
- [28] Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching:
 Generalized score matching for discrete data, 2023.
- [29] Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the
 ratios of the data distribution, 2024.
- [30] William J Anderson. *Continuous-time Markov chains: An applications-oriented approach*.
 Springer Science & Business Media, 2012.
- [31] Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time
 discrete diffusion models. In *The Eleventh International Conference on Learning Representa- tions*, 2023.
- [32] Frank Kelly. Reversibility and stochastic networks. 1980. URL https://api.
 semanticscholar.org/CorpusID:125211322.
- [33] Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007.github.io/
 OpenWebTextCorpus, 2019.
- [34] Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi,
 Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernandez. The LAMBADA
 dataset: Word prediction requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages
 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics. URL
- 393 http://www.aclweb.org/anthology/P16-1144.
- [35] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
 models, 2016.

- [36] Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and
 Tony Robinson. One billion word benchmark for measuring progress in statistical language
 modeling, 2014.
- [37] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language
 models are unsupervised multitask learners. 2019. URL https://api.semanticscholar.
 org/CorpusID:160025533.
- [38] Ishaan Gulrajani and Tatsunori Hashimoto. Likelihood-based diffusion language models. In
 Advances in Neural Information Processing Systems, 2023.
- [39] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax
 flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems*, 34:12454–12465, 2021.
- [40] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of
 deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*,
 2018.
- [41] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom
 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia
 Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and L. Sifre.
 Training compute-optimal large language models. *ArXiv*, abs/2203.15556, 2022. URL https:
 //api.semanticscholar.org/CorpusID:247778764.
- [42] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth
 words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22669–22679, 2023.
- [43] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.
- [44] Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao,
 Hang Su, and Jun Zhu. One transformer fits all distributions in multi-modal diffusion at scale.
 arXiv preprint arXiv:2303.06555, 2023.
- [45] Fan Bao, Chendong Xiang, Gang Yue, Guande He, Hongzhou Zhu, Kaiwen Zheng, Min
 Zhao, Shilong Liu, Yaole Wang, and Jun Zhu. Vidu: a highly consistent, dynamic and skilled
 text-to-video generator with diffusion models. *arXiv preprint arXiv:2405.04233*, 2024.
- [46] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,
 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Neural Information Processing Systems*, 2017. URL https://api.semanticscholar.org/CorpusID:13756489.
- [47] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
 deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Confer- ence of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational
 Linguistics, 2019.
- [48] William S. Peebles and Saining Xie. Scalable diffusion models with transformers. In *Interna- tional Conference on Computer Vision*, 2023.
- [49] Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer
 with rotary position embedding. *Neurocomputing*, 2021.

439 A Proof of Proposition 1

Since the different dimensions of the forward diffusion process are independent of each other, we can first analyze the conditional distribution in one dimension. This can be derived directly from Eq. (2.3), but for a better understanding, here we provide a more intuitive proof in the case when $Q_t = Q^{absorb}$.

Lemma 1. (Analytic conditional distribution for absorbing case) Suppose $\{X_t\}$ is a continuous time Markov chain with transition rate matrix $Q_t = \sigma(t)Q^{absorb}$, given the value x at time zero, the conditional distribution $p_{t|0}(x_t|x)$ has the following analytic form:

$$p_{t|0}(x_t|x) = \begin{cases} e^{-\bar{\sigma}(t)}, & x_t = x, \\ 1 - e^{-\bar{\sigma}(t)}, & x_t = [\mathbf{M}], \\ 0. & x_t \neq [\mathbf{M}] \text{ and } x_t \neq [\mathbf{M}]. \end{cases}$$
(A.1)

447 *Proof.* Given the initial value $x \in \mathcal{X} = \{1, \dots, N\}$, we have

$$x_t = \begin{cases} x, & t < T_h, \\ [\mathbf{M}], & t \ge T_h, \end{cases}$$

- where T_h is the holding time before the system transitions to the next state.
- 449 Based on the properties of the Q^{absorb} :

$$p_{t+\Delta t|t}(x|x) = 1 - (-\sigma(t)\boldsymbol{Q}^{\text{absorb}}(x,x))\Delta t + o(\Delta t).$$
(A.2)

- Partitioning the interval [0, t] into $\{s_k\}_{k=0}^n$, make use of Memoryless Property of Continuous-Time
- 451 Markov Chains:

$$p_{t|0}(x|x) = \prod_{k=1}^{n} p_{s_k|s_{k-1}}(x|x)$$
(A.3)

$$=\prod_{k=1}^{n} (1 - (-\sigma(t_k)\boldsymbol{Q}^{\text{absorb}}(x,x))(s_k - s_{k-1}) + o((s_k - s_{k-1})))$$
(A.4)

$$= \exp(\sum_{k=1}^{n} \ln(1 - (-\sigma(t_k)\boldsymbol{Q}^{\text{absorb}}(x, x))(s_k - s_{k-1}) + o((s_k - s_{k-1}))))$$
(A.5)

$$= \exp(\sum_{k=1}^{n} -(-\sigma(t_k)\boldsymbol{Q}^{\text{absorb}}(x,x))(s_k - s_{k-1}) + o((s_k - s_{k-1}))).$$
(A.6)

452 Let $\max(s_k - s_{k-1}) \to 0$, we have:

$$p_{t|0}(x|x) = \exp\left(-\int_0^t -\sigma(s)\boldsymbol{Q}^{\text{absorb}}(x,x)ds\right) = \exp\left(-\left(-\boldsymbol{Q}^{\text{absorb}}(x,x)\bar{\sigma}(t)\right)\right).$$
(A.7)

453 As $\boldsymbol{Q}^{\mathrm{absorb}}(x,x) = -1,$ we have

$$p_{t|0}(x|x) = P(T_h > t) = e^{-\bar{\sigma}(t)},$$
 (A.8)

454 455

$$p_{t|0}([\mathbf{M}]|x) = P(T_h > t) = 1 - e^{-\bar{\sigma}(t)},$$
 (A.9)

$$p_{t|0}(k|x) = 0 \quad \text{if } k \neq [\mathbf{M}] \text{ and } k \neq x.$$
(A.10)

456

457 **Proposition 1.** (Analytic joint distribution for absorbing case)

Suppose $\{X_t\}$ is a continuous time Markov chain with transition rate matrix $Q_t = \sigma(t)Q^{absorb}$. For $\mathbf{x}_t = x_t^1 \cdots x_t^d$ with N_1 components as $[\mathbf{M}]$ and $N_2 = d - N_1$ components as specific value, $p_t(\mathbf{x}_t)$ can be expressed as Eq. (A.11):

$$p_t(\boldsymbol{x}_t) = [1 - e^{-\bar{\sigma}(t)}]^{N_1} [e^{-\bar{\sigma}(t)}]^{N_2} p_0(\boldsymbol{x}_t^{UM}),$$
(A.11)

461 where $x_t^{UM} := \{ x^k | x^k \neq [M] \}$ represents unmasked part of x_t^{UM} .

Proposition 1 shows that the joint distribution $p_t(x_t)$ can be expressed as the multiplication of two 462 terms. One is an analytic term only depending on time, the other is a joint distribution of clean data 463 $p_0(x_t^{\text{UM}})$ with N_2 dimensions independent of time.

464

Proof. Without loss of generality, let's assume that the preceding N_1 terms of x are all [**M**], and the remaining N_2 terms are fixed at specific values. That is, $x_t = [\mathbf{M}] \cdots [\mathbf{M}] x_t^{N_1+1} \cdots x_t^d$, and here x^k 465 466 is a fixed value in \mathcal{X} . 467

Use the law of total probability and Lemma 1, along with independent property: 468

$$\begin{split} p_t([\mathbf{M}]\cdots[\mathbf{M}]x_t^{N_1+1}\cdots x_t^d) \\ &= \sum_{x_0\in\mathcal{X}^d} p_{t|0}([\mathbf{M}]\cdots[\mathbf{M}]x_t^{N_1+1}\cdots x_t^d|x_0)p_0(x_0) \\ &= \sum_{x_0^1\in\mathcal{X},\cdots,x_0^d\in\mathcal{X}} p_{t|0}([\mathbf{M}]\cdots[\mathbf{M}]x_t^{N_1+1}\cdots x_t^d|x_0^1\cdots x_0^d)p_0(x_0^1\cdots x_0^d) \\ &= \sum_{x_0^1\in\mathcal{X},\cdots,x_0^d\in\mathcal{X}} \prod_{k=1}^{N_1} p_{t|0}^k([\mathbf{M}]|x_0^k) \prod_{k=N_1+1}^d p_{t|0}^k(x_t^k|x_0^k)p_0(x_0^1\cdots x_0^d) \\ &= \sum_{x_0^1\in\mathcal{X},\cdots,x_0^{N_1}\in\mathcal{X}} \prod_{k=1}^{N_1} p_{t|0}^k([\mathbf{M}]|x_0^k)[e^{-\bar{\sigma}(t)}]^{N_2}p_0(x_0^1\cdots x_0^{N_1}x_t^{N_1+1}\cdots x_t^d) \\ &= \sum_{x_0^1\in\mathcal{X},\cdots,x_0^{N_1}\in\mathcal{X}} [1-e^{-\bar{\sigma}(t)}]^{N_1}[e^{-\bar{\sigma}(t)}]^{N_2}p_0(x_0^1\cdots x_0^{N_1}x_t^{N_1+1}\cdots x_t^d) \\ &= [1-e^{-\bar{\sigma}(t)}]^{N_1}[e^{-\bar{\sigma}(t)}]^{N_2}\sum_{x_0^1\in\mathcal{X},\cdots,x_0^{N_1}\in\mathcal{X}} p_0(x_0^1\cdots x_0^{N_1}x_t^{N_1+1}\cdots x_t^d) \\ &= [1-e^{-\bar{\sigma}(t)}]^{N_1}[e^{-\bar{\sigma}(t)}]^{N_2}p_0(x_t^{N_1+1}\cdots x_t^d). \end{split}$$

In the general case, we have: 469

$$p_t(\boldsymbol{x}_t) = [1 - e^{-\bar{\sigma}(t)}]^{N_1} [e^{-\bar{\sigma}(t)}]^{N_2} p_0(\boldsymbol{x}_t^{\text{UM}}),$$

which shows that the likelihood of noisy data x at time t equals the likelihood of unmasked part of x470 at time 0 multiplied by a analytic time-dependent term. 471

Proof of Theorem 1 B 472

Theorem 1. (Analytic concrete score in absorbing case, proof in Appendix B) For $\mathbf{x}_t = x_t^1 \dots x_t^i$ and $\hat{\mathbf{x}}_t = x_t^1 \dots \hat{\mathbf{x}}_t^i \dots x_t^d$, if $x_t^i = [\mathbf{M}]$ and $\hat{x}_t^i \neq [\mathbf{M}]$, the concrete score at time t can be expressed as a time-independent conditional distribution at time zero multiplied by an 473 474 475 analytic time-dependent term: 476

$$\frac{p_t\left(x_t^1\dots\widehat{x}_t^i\dots x_t^d\right)}{p_t\left(x_t^1\dots x_t^i\dots x_t^d\right)} = \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} p_0(\hat{x}_t^i | \boldsymbol{x}_t^{UM}),$$

where x_{t}^{UM} is the vector consists of all unmasked tokens of x_{t} . 477

Proof. According to Proposition 1, if $x_t^i = [\mathbf{M}]$ and $\hat{x}_t^i \neq [\mathbf{M}], \hat{x}_t^{\mathrm{UM}} = (\boldsymbol{x}_t^{\mathrm{UM}}, \hat{x}_t^i),$ 478

$$\begin{aligned} \frac{p_t(\hat{\boldsymbol{x}}_t)}{p_t(\boldsymbol{x}_t)} = & \frac{[1 - e^{-\sigma(t)}]^{N_1 - 1} [e^{-\sigma(t)}]^{N_2 + 1} p_0(\hat{\boldsymbol{x}}_t^{\text{UM}})}{[1 - e^{-\bar{\sigma}(t)}]^{N_1} [e^{-\bar{\sigma}(t)}]^{N_2} p_0(\boldsymbol{x}_t^{\text{UM}})} \\ = & \frac{[1 - e^{-\bar{\sigma}(t)}]^{N_1 - 1} [e^{-\bar{\sigma}(t)}]^{N_2 + 1} p_0(\boldsymbol{x}_t^{\text{UM}}, \hat{\boldsymbol{x}}_t^i)}{[1 - e^{-\bar{\sigma}(t)}]^{N_1} [e^{-\bar{\sigma}(t)}]^{N_2} p_0(\boldsymbol{x}_t^{\text{UM}})} \\ = & \frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} p_0(\hat{\boldsymbol{x}}_t^i | \boldsymbol{x}_t^{\text{UM}}). \end{aligned}$$

479

480 C Proof of Theorem 2

481 C.1 Denoising cross-entropy loss by λ

482 According to definition of Q_t we can simplify Eq. (3.8) as:

$$\begin{split} \mathcal{L}_{\text{DCE}}^{\infty}(\boldsymbol{x}_{0}) \\ &= \int_{0}^{\infty} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{t|0}(\tilde{\boldsymbol{x}}|\boldsymbol{x}_{0})} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}], j \neq [\mathbf{M}]} \sigma(t) \left(-\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} I(x_{0}^{i} = j) \log \left(\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} c_{\theta}(\tilde{\boldsymbol{x}})[i, j] \right) \right) \right] dt \\ &= \int_{0}^{\infty} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{t|0}(\tilde{\boldsymbol{x}}|\boldsymbol{x}_{0})} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} \sigma(t) \left(-\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \log \left(\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} c_{\theta}(\tilde{\boldsymbol{x}})[i, x_{0}^{i}] \right) \right) \right] dt \\ &= \int_{0}^{\infty} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{t|0}(\tilde{\boldsymbol{x}}|\boldsymbol{x}_{0})} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} \sigma(t) \left(-\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} \log \left(\frac{e^{-\bar{\sigma}(t)}}{1 - e^{-\bar{\sigma}(t)}} q_{\theta}(x_{0}^{i}|\tilde{\boldsymbol{x}}^{\text{UM}}) \right) \right) \right] dt. \end{split}$$

483 Define $\lambda(t) = 1 - e^{-\bar{\sigma}(t)}$, $d\lambda = \sigma(t)e^{-\bar{\sigma}(t)}dt$. As $\bar{\sigma}(t) = \int_0^t \sigma(\tau)d\tau$, we have $\lambda(0) = 0$, 484 $\lim_{t\to\infty} \lambda(t) = 1$. By a change of variables for the integration variable from t to λ , we can 485 rewrite the above equation as:

$$\int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}} | \boldsymbol{x}_{0})} \left[\sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(\frac{1-\lambda}{\lambda} q_{\theta}(x_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right] d\lambda$$
$$= \int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}} | \boldsymbol{x}_{0})} \sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(x_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) d\lambda + \int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}} | \boldsymbol{x}_{0})} \sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(\frac{1-\lambda}{\lambda}) \right) d\lambda$$

486 By independence of forward process and Lemma 1, $p_{t|0}(\tilde{x}|x_0) = \prod_{i=1}^d p_{t|0}(\tilde{x}^i|x_0^i)$ where

$$p_{t|0}(\tilde{x}^{i}|x_{0}^{i}) = \begin{cases} 1 - e^{-\bar{\sigma}(t)} & \tilde{x}^{i} = [\mathbf{M}], \\ e^{-\bar{\sigma}(t)} & \tilde{x}^{i} = x_{0}^{i}, \\ 0 & \text{else.} \end{cases}$$
(C.1)

487 Therefore, $p_{\lambda}(\tilde{\pmb{x}}|\pmb{x}_0) = \prod_{i=1}^d p_{\lambda}(\tilde{x}^i|x_0^i)$ where

$$p_{\lambda}(\tilde{x}^{i}|x_{0}^{i}) = \begin{cases} \lambda & \tilde{x}^{i} = [\mathbf{M}], \\ 1 - \lambda & \tilde{x}^{i} = x_{0}^{i}, \\ 0 & \text{else.} \end{cases}$$
(C.2)

488 Consider the second term, we have:

$$\begin{split} &\int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}} \mid \boldsymbol{x}_{0})} \left[\sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(\frac{1-\lambda}{\lambda}) \right) \right] d\lambda \\ &= \int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}} \mid \boldsymbol{x}_{0})} \left[\sum_{i=1}^{d} \mathbb{I}(\tilde{x}^{i} = [\mathbf{M}]) \left(-\log(\frac{1-\lambda}{\lambda}) \right) \right] d\lambda \\ &= \int_{0}^{1} \frac{1}{\lambda} \left[\sum_{i=1}^{d} p_{\lambda}(\tilde{x}^{i} = [\mathbf{M}] \mid \boldsymbol{x}_{0}) \left(-\log(\frac{1-\lambda}{\lambda}) \right) \right] d\lambda \\ &= d \int_{0}^{1} -\log(\frac{1-\lambda}{\lambda}) d\lambda \\ &= d \left(\lambda \log \lambda + (1-\lambda) \log(1-\lambda) \right) |_{0}^{1}. \end{split}$$

Note that:

$$\lim_{\lambda \to 0} \lambda \log \lambda = \lim_{\lambda \to 1} (1 - \lambda) \log(1 - \lambda) = 0,$$

 $\text{ therefore, } (\lambda \log \lambda + (1-\lambda) \log(1-\lambda))|_0^1 = 0.$

$$\mathcal{L}_{\text{DCE}}^{\infty}(\boldsymbol{x}_{0}) = \int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}}|\boldsymbol{x}_{0})} \left[\sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(x_{0}^{i}|\tilde{\boldsymbol{x}}^{\text{UM}})) \right) \right] d\lambda.$$
(C.3)

490 C.2 Denoising cross-entropy loss by k

491 By Eq. (C.3), we can express the loss in terms of λ . Given x_0 , we denote $\tilde{\mathcal{X}} := \{x_0^1, [\mathbf{M}]\} \times \cdots \{x_0^d, [\mathbf{M}]\}$ as the sample space of $\tilde{\boldsymbol{x}}$, and define $\tilde{\mathcal{X}}_k := \{\tilde{\boldsymbol{x}} : \tilde{\boldsymbol{x}} \in \tilde{\mathcal{X}} \land 493 \quad \tilde{\boldsymbol{x}}$ has exact k dimensions with values $[\mathbf{M}]\}$. Obviously, $|\tilde{\mathcal{X}}| = 2^d$ and $|\tilde{\mathcal{X}}_k| = C_d^k$. We have:

$$\int_{0}^{1} \frac{1}{\lambda} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim p_{\lambda}(\tilde{\boldsymbol{x}} | \boldsymbol{x}_{0})} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(\boldsymbol{x}_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right] d\lambda$$
(C.4)

$$= \int_{0}^{1} \frac{1}{\lambda} \sum_{\tilde{x} \in \tilde{\mathcal{X}}} p_{\lambda}(\tilde{x} | \boldsymbol{x}_{0}) \left[\sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(x_{0}^{i} | \tilde{x}^{\mathrm{UM}})) \right) \right] d\lambda$$
(C.5)

$$= \int_{0}^{1} \frac{1}{\lambda} \sum_{k=0}^{d} \sum_{\tilde{\boldsymbol{x}} \in \tilde{\mathcal{X}}_{k}} \lambda^{k} (1-\lambda)^{d-k} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(\boldsymbol{x}_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right] d\lambda$$
(C.6)

$$= \int_{0}^{1} \frac{1}{\lambda} \sum_{k=1}^{d} \sum_{\tilde{\boldsymbol{x}} \in \tilde{\mathcal{X}}_{k}} \lambda^{k} (1-\lambda)^{d-k} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(\boldsymbol{x}_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right] d\lambda$$
(C.7)

$$=\sum_{k=1}^{d}\int_{0}^{1}\lambda^{k-1}(1-\lambda)^{d-k}d\lambda\sum_{\tilde{\boldsymbol{x}}\in\tilde{\mathcal{X}}_{k}}\left[\sum_{\tilde{x}^{i}=[\mathbf{M}]}\left(-\log(q_{\theta}(x_{0}^{i}|\tilde{\boldsymbol{x}}^{\mathrm{UM}}))\right)\right]$$
(C.8)

$$=\sum_{k=1}^{d} \frac{(k-1)!(d-k)!}{d!} \sum_{\tilde{\boldsymbol{x}} \in \tilde{\mathcal{X}}_{k}} \left[\sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(x_{0}^{i}|\tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right]$$
(C.9)

$$= \sum_{k=1}^{d} \frac{1}{kC_{d}^{k}} \sum_{\tilde{\boldsymbol{x}} \in \tilde{\mathcal{X}}_{k}} \left[\sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(x_{0}^{i}|\tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right].$$
(C.10)

⁴⁹⁴ This can be reformulated in the form of expectation:

$$\sum_{k=1}^{d} \frac{1}{kC_{d}^{k}} \sum_{\tilde{\boldsymbol{x}} \in \tilde{\mathcal{X}}_{k}} \left[\sum_{\tilde{x}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(x_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right]$$
(C.11)

$$= \sum_{k=1}^{d} \frac{1}{k} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim U(\tilde{\mathcal{X}}_{k})} \left[\sum_{\tilde{\boldsymbol{x}}^{i} = [\mathbf{M}]} \left(-\log(q_{\theta}(\boldsymbol{x}_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right]$$
(C.12)

$$= d\mathbb{E}_{k \sim U(\{1, \cdots, d\})} \frac{1}{k} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim U(\tilde{\boldsymbol{X}}_k)} \left[\sum_{\tilde{\boldsymbol{x}}^i = [\mathbf{M}]} \left(-\log(q_{\theta}(\boldsymbol{x}_0^i | \tilde{\boldsymbol{x}}^{\mathrm{UM}})) \right) \right].$$
(C.13)

495 C.3 Exact negative likelihood

Let S_d represent the set of all permutations of the integers $1, \dots, d$, and let $\pi \in S_d$ be one of these permutations. Then, we can express $\log q_\theta(x_0)$ as follows: $\log q_{\theta}(\boldsymbol{x}_0) = \mathbb{E}_{\pi \sim U(S_d)} \log q_{\theta}(\boldsymbol{x}_0)$ (C.14)

$$= \mathbb{E}_{\pi \sim U(S_d)} \sum_{l=1}^d \log q_\theta(x_0^{\pi(l)} | x_0^{\pi((C.15)$$

$$= \sum_{l=1}^{d} \mathbb{E}_{\pi \sim U(S_d)} \log q_{\theta}(x_0^{\pi(l)} | x_0^{\pi((C.16)$$

$$= \sum_{l=1}^{d} \frac{1}{d-l+1} \mathbb{E}_{\pi \sim U(S_d)} \sum_{r=l}^{d} \log q_{\theta}(x_0^{\pi(r)} | x_0^{\pi((C.17)$$

$$= \sum_{k=1}^{d} \frac{1}{k} \mathbb{E}_{\pi \sim U(S_d)} \sum_{r=d-k+1}^{d} \log q_{\theta}(x_0^{\pi(r)} | x_0^{\pi((C.18)$$

$$= d\mathbb{E}_{k \sim U(\{1, \cdots, d\})} \frac{1}{k} \mathbb{E}_{\pi \sim U(S_d)} \sum_{r=d-k+1}^d \log q_\theta(x_0^{\pi(r)} | x_0^{\pi((C.19)$$

In this context, for a fixed k, the condition $x_0^{\pi(< d-k+1)}$ can be understood as the unmasked part of noisy data \tilde{x}^{UM} . For $r = d - k + 1, \dots, d$, $x_0^{\pi(r)}$ corresponds to the k items of the masked part. Therefore, we have:

$$\mathbb{E}_{\pi \sim U(S_d)} \sum_{r=d-k+1}^{d} \log q_{\theta}(x_0^{\pi(r)} | x_0^{\pi((C.20)$$

⁵⁰¹ Thus, substituting back, we have:

$$-\log q_{\theta}(\boldsymbol{x}_{0}) = d\mathbb{E}_{k \sim U(\{1, \cdots, d\})} \frac{1}{k} \mathbb{E}_{\tilde{\boldsymbol{x}} \sim U(\tilde{\mathcal{X}}_{k})} \sum_{\tilde{x}^{i} = [\mathbf{M}]} -\log q_{\theta}(x_{0}^{i} | \tilde{\boldsymbol{x}}^{\mathrm{UM}}), \quad (C.21)$$

⁵⁰² which is exactly Eq. (C.13).

This concludes the proof of the exact negative likelihood, showing the equivalence between the expected negative log-likelihood and the denoising cross-entropy formulation.

505 **D** Sampling methods of discrete diffusion

506 D.1 Original form in discrete diffusion

Euler discrete method According to the Eq. (2.7), take $t = s - \Delta s$ and use the Euler method, we can simulate the reverse process by iteratively taking small Δt Euler steps at time s, calculate the reverse transition rate based on $s_{\theta}(x_s, s)$, and randomly sampling $x_{s-\Delta s}$.

510 Left term:

$$\frac{d}{dt} \mathbf{P}_{s \to t}{}_{|t=s-\Delta s} \approx \frac{\mathbf{P}_{s \to s-\Delta s} - \mathbf{P}_{s \to s}}{\Delta s}.$$
(D.1)

511 Right term:

$$\boldsymbol{P}_{s \to t} \boldsymbol{\tilde{Q}}_{t|t=s-\Delta s} \approx \boldsymbol{P}_{s \to s} \boldsymbol{\tilde{Q}}_s. \tag{D.2}$$

512 As $P_{s \rightarrow s} = I$, we have

$$P_{s \to s - \Delta s} \approx I + \dot{Q}_s \Delta s.$$
 (D.3)

Rewrite in t and consider a specific input x_t , $x_{t-\Delta t}$ is sampled from the following transition probabilities:

$$p_{t-\Delta t|t}(x_{t-\Delta t}|x_t) \approx \delta_{x_t x_{t-\Delta t}} + \tilde{Q}_t(x_t, x_{t-\Delta t})\Delta t + O(\Delta t)$$
(D.4)

$$\approx \delta_{x_t x_{t-\Delta t}} + \tilde{\boldsymbol{Q}}_t(x_t, x_{t-\Delta t}) \Delta t, \tag{D.5}$$

515 where

$$\tilde{\boldsymbol{Q}}_t(x_t, x_{t-\Delta t}) \approx \begin{cases} \boldsymbol{Q}_t(x_{t-\Delta t}, x_t) \boldsymbol{s}_{\boldsymbol{\theta}}(x_t, t)_{x_{t-\Delta t}} & x_t \neq x_{t-\Delta t}, \\ -\sum_{k \neq x_t} \tilde{\boldsymbol{Q}}_t(x_t, k) & x_t = x_{t-\Delta t}. \end{cases}$$
(D.6)

Tweedie τ -leaping If we know the analytic form of $P_{s \to t}$, it is possible to get the closed form of reverse probability $P_{t \to s}$ for any s < t. According to the conditional decomposition of total probability, we have:

diag
$$(P_t^T) \mathbf{P}_{t \to s} = \left(\text{diag} \left(P_s^T \right) \mathbf{P}_{s \to t} \right)^T$$
. (D.7)

519 As $P_s^T P_{s \to t} = P_t^T$, the following equation holds:

$$\boldsymbol{P}_{t \to s} = \operatorname{diag}\left(\boldsymbol{P}_{t}^{T}\right)^{-1} \boldsymbol{P}_{s \to t}^{T} \operatorname{diag}\left(\boldsymbol{P}_{s}^{T}\right) = \operatorname{diag}\left(\boldsymbol{P}_{t}^{T}\right)^{-1} \boldsymbol{P}_{s \to t}^{T} \operatorname{diag}\left(\boldsymbol{P}_{t}^{T} \boldsymbol{P}_{s \to t}^{-1}\right).$$
(D.8)

Given x_t , to get $p_{s|t}(x_s|x_t)$, we only need to calculate row x_t of $P_{t\to s}$:

$$\boldsymbol{P}_{t \to s}(x_t, \cdot) = \frac{1}{p_t(x_t)} \boldsymbol{P}_{s \to t}^T(x_t, \cdot) \odot \left(\boldsymbol{P}_t^T \boldsymbol{P}_{s \to t}^{-1} \right)$$
(D.9)

$$= \boldsymbol{P}_{s \to t}^{T}(x_{t}, \cdot) \odot \left(\frac{P_{t}^{T}}{p_{t}(x_{t})} \boldsymbol{P}_{s \to t}^{-1}\right) \approx \boldsymbol{P}_{s \to t}^{T}(x_{t}, \cdot) \odot \left(\boldsymbol{s}_{\theta}(x_{t}, t)^{T} \boldsymbol{P}_{s \to t}^{-1}\right).$$
(D.10)

521 D.2 Simplified form in reparameterized absorbing discrete diffusion

Euler discrete method For $x_t = [\mathbf{M}]$, given the value of \hat{x}_t , use the $Q_t(\hat{x}_t, x_t) = \sigma(t) Q^{\text{absorb}}(\hat{x}_t, x_t)$ and $s_{\theta}(x_t, t)_{\hat{x}_t} = \frac{e^{-\sigma(t)}}{1 - e^{-\sigma(t)}} c_{\theta}(x_t)_{\hat{x}_t}$. Eq. (D.5) can be simplified as:

$$p_{t-\Delta t|t}(\hat{x}_t|[\mathbf{M}]) = \begin{cases} \sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1-e^{-\bar{\sigma}(t)}} \Delta t \mathbf{c}_{\theta}(x_t)_{\hat{x}_t} & \text{if } \hat{x}_t \neq [\mathbf{M}], \\ 1 - \sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1-e^{-\bar{\sigma}(t)}} \Delta t & \text{if } \hat{x}_t = [\mathbf{M}]. \end{cases}$$
(D.11)

524 For multi-dimension cases, similar results can be obtained:

$$p_{t-\Delta t|t}(x_{t-\Delta t}^{i}|\boldsymbol{x}_{t}) = \begin{cases} \sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1-e^{-\bar{\sigma}(t)}} \Delta t \boldsymbol{c}_{\theta}(\boldsymbol{x}_{t})[i, x_{t-\Delta t}^{i}] & \text{if } x_{t-\Delta t}^{i} \neq [\mathbf{M}], \\ 1-\sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1-e^{-\bar{\sigma}(t)}} \Delta t & \text{if } x_{t-\Delta t}^{i} = [\mathbf{M}]. \end{cases}$$
(D.12)

525 for all $x_t^i = [\mathbf{M}]$.

Tweedie τ -**leaping** Suppose $x_t = x_t^1 \cdots x_t^d$ has N_1 components as $[\mathbf{M}]$ and $N_2 = d - N_1$ components as specific values. Without loss of generality, let's assume that the preceding N_1 terms of x_t are all $[\mathbf{M}]$, and the remaining N_2 terms are fixed at specific values. For $1 \le i \le d$, given the value of $x_{t-\Delta t}^i \ne [\mathbf{M}]$, :

$$p_{t-\Delta t|t}(x_{t-\Delta t}^{i}|\boldsymbol{x}_{t}) = \frac{p_{t,t-\Delta t}(\boldsymbol{x}_{t}, x_{t-\Delta t}^{i})}{p_{t}(\boldsymbol{x}_{t})}.$$
(D.13)

530 By Proposition 1:

$$p_t(\boldsymbol{x}_t) = [1 - e^{-\bar{\sigma}(t)}]^{N_1} [e^{-\bar{\sigma}(t)}]^{N_2} p_0(x_t^{N_1+1} \cdots x_t^d).$$
(D.14)

531 Similar to the proof in Proposition 1:

$$\begin{split} & p_{t,t-\Delta t}(\boldsymbol{x}_{t}, \boldsymbol{x}_{t-\Delta t}^{i}) \\ &= \sum_{x_{0} \in \mathcal{X}} p_{(t,t-\Delta t)|0} p(\boldsymbol{x}_{t}, \boldsymbol{x}_{t-\Delta t}^{i} | \boldsymbol{x}_{0}) p_{0}(\boldsymbol{x}_{0}) \\ &= \sum_{x_{0}^{1} \in \mathcal{X}, \cdots, x_{0}^{d} \in \mathcal{X}} p_{(t,t-\Delta t)|0}([\mathbf{M}] \cdots [\mathbf{M}] \boldsymbol{x}_{t}^{N_{1}+1} \cdots \boldsymbol{x}_{t}^{d}, \boldsymbol{x}_{t-\Delta t}^{i} | \boldsymbol{x}_{0}^{1} \cdots \boldsymbol{x}_{0}^{d}) p_{0}(\boldsymbol{x}_{0}^{1} \cdots \boldsymbol{x}_{0}^{d}) \\ &= \sum_{x_{0}^{1} \in \mathcal{X}, \cdots, x_{0}^{d} \in \mathcal{X}} p_{(t,t-\Delta t)|0}([\mathbf{M}], \boldsymbol{x}_{t-\Delta t}^{i} | \boldsymbol{x}_{0}^{i}) \prod_{k=1, k \neq i}^{N_{1}} p_{t|0}([\mathbf{M}] | \boldsymbol{x}_{0}^{k}) \prod_{k=N_{1}+1}^{d} p_{t|0}(\boldsymbol{x}^{k} | \boldsymbol{x}_{0}^{k}) p_{0}(\boldsymbol{x}_{0}^{1} \cdots \boldsymbol{x}_{0}^{d}) \\ &= \sum_{x_{0}^{k} \in \mathcal{X}, k \in \{1, \cdots, N_{1}\}/\{i\}} p_{(t,t-\Delta t)|0}([\mathbf{M}], \boldsymbol{x}_{t-\Delta t}^{i} | \boldsymbol{x}_{t-\Delta t}^{i} | \boldsymbol{x}_{t-\Delta t}^{i}) \prod_{k=1, k \neq i}^{N_{1}} p_{t|0}([\mathbf{M}] | \boldsymbol{x}_{0}^{k}) [e^{-\bar{\sigma}(t)}]^{N_{2}} \\ &p_{0}(\boldsymbol{x}_{0}^{1} \cdots \boldsymbol{x}_{0}^{i-1} \boldsymbol{x}_{t-\Delta t}^{i} \boldsymbol{x}_{0}^{i+1} \cdots \boldsymbol{x}_{0}^{N_{1}} \boldsymbol{x}_{t}^{N_{1}+1} \cdots \boldsymbol{x}_{t}^{N_{1}+1}) \\ &= \sum_{x_{0}^{k} \in \mathcal{X}, k \in \{1, \cdots, N_{1}\}/\{i\}} (e^{-\bar{\sigma}(t-\Delta t)} - e^{-\bar{\sigma}(t)})(1 - e^{-\bar{\sigma}(t)})^{N_{1}-1} [e^{-\bar{\sigma}(t)}]^{N_{2}} \\ &p_{0}(\boldsymbol{x}_{0}^{1} \cdots \boldsymbol{x}_{0}^{i-1} \boldsymbol{x}_{t-\Delta t}^{i} \boldsymbol{x}_{0}^{i+1} \cdots \boldsymbol{x}_{0}^{N_{1}} \boldsymbol{x}_{t}^{N_{1}+1} \cdots \boldsymbol{x}_{t}^{d}) \\ &= (e^{-\bar{\sigma}(t-\Delta t)} - e^{-\bar{\sigma}(t)})(1 - e^{-\bar{\sigma}(t)})^{N_{1}-1} [e^{-\bar{\sigma}(t)}]^{N_{2}} p_{0}(\boldsymbol{x}_{t-\Delta t}^{i}, \boldsymbol{x}_{t}^{N_{1}+1} \cdots \boldsymbol{x}_{t}^{d}). \end{split}$$

532 Note we used the fact that:

$$p_{(t,t-\Delta t)|0}([\mathbf{M}], x_{t-\Delta t}^{i}|x_{t-\Delta t}^{i}) = p_{t|t-\Delta t}([\mathbf{M}]|x_{t-\Delta t}^{i})p_{t-\Delta t|0}(x_{t-\Delta t}^{i}|x_{t-\Delta t}^{i})$$
$$= (1 - e^{-(\bar{\sigma}(t) - \bar{\sigma}(t-\Delta t))})e^{-\bar{\sigma}(t-\Delta t)}$$
$$= e^{-\bar{\sigma}(t-\Delta t)} - e^{-\bar{\sigma}(t)},$$

533

$$p_{t|0}([\mathbf{M}]|x_0^k) = 1 - e^{-\bar{\sigma}(t)},$$

⁵³⁴ by dividing the two expressions, we have:

$$p_{t-\Delta t|t}(x_{t-\Delta t}^{i}|\boldsymbol{x}_{t}) = \frac{e^{-\bar{\sigma}(t-\Delta t)} - e^{-\bar{\sigma}(t)}}{1 - e^{\bar{\sigma}(t)}} p_{0}(x_{t-\Delta t}^{i}|x_{t}^{N_{1}+1}\cdots x_{t}^{d})$$
(D.15)

$$\approx \frac{e^{-\bar{\sigma}(t-\Delta t)} - e^{-\bar{\sigma}(t)}}{1 - e^{\bar{\sigma}(t)}} \boldsymbol{c}_{\theta}(\boldsymbol{x}_t)[i, x_{t-\Delta t}^i].$$
(D.16)

535 In general, for $x_t^i = [\mathbf{M}]$, we have:

$$p_{t-\Delta t|t}(x_{t-\Delta t}^{i}|\boldsymbol{x}_{t}) \begin{cases} \approx \frac{e^{-\bar{\sigma}(t-\Delta t)} - e^{-\bar{\sigma}(t)}}{1 - e^{\bar{\sigma}(t)}} \boldsymbol{c}_{\theta}(\boldsymbol{x}_{t})[i, x_{t-\Delta t}^{i}], & x_{t-\Delta t}^{i} \neq [\mathbf{M}], \\ = \frac{1 - e^{-\bar{\sigma}(t-\Delta t)}}{1 - e^{-\bar{\sigma}(t)}}, & x_{t-\Delta t}^{i} = [\mathbf{M}]. \end{cases}$$
(D.17)

536 D.3 Discuss on the expectation of NFE

As discussed in Section Appendix D.2, for both the Euler method and Tweedie τ -leaping, the probability $p_{t-\Delta t|t}^{i}([\mathbf{M}]|\mathbf{x}_{t})$ is only a factor of time which is independent of the other dimensions of \mathbf{x}_{t} once given $x_{t}^{i} = [\mathbf{M}]$. By the Law of Total Probability, it is easy to find that $p_{t-\Delta t|t}^{i}([\mathbf{M}]|[\mathbf{M}])$ is also only a factor of time. Thus, given a specific sampling method and a set of time steps $\{t_{0} = 0, \dots, t_{n} = T\}$, the NFE can be treated as a random variable with a calculable expected value. Let N_{k} denote the number of dimensions of \mathbf{x} which changed in $[t_{k-1}, t_{k})$, so we have:

$$NFEs(n) = \sum_{k=1}^{n} \mathbb{I}(N_k \neq 0), \qquad (D.18)$$

543

E-NFEs
$$(n) = \sum_{k=1}^{n} \mathbb{E}[\mathbb{I}(N_k \neq 0)] = \sum_{k=1}^{n} P(N_k \neq 0).$$
 (D.19)

- For each dimension *i*, let r_k represent the probability that x^i changes within the interval $[t_{k-1}, t_k)$.
- ⁵⁴⁵ Consequently, N_k follows a binomial distribution with parameters l and r_k , denoted as $N_k \sim$ ⁵⁴⁶ Binomial (l, r_k) .

E-NFEs $(n) = \sum_{k=1}^{n} P(N_k \neq 0) = \sum_{k=1}^{n} (1 - (1 - r_k)^l).$ (D.20)

547 By definition of r_k and property of absorbing diffusion:

$$r_{k} = P(X_{t_{k-1}}^{i} \neq [\mathbf{M}], X_{t_{k}}^{i} = [\mathbf{M}] | X_{t_{n}}^{i} = [\mathbf{M}])$$
(D.21)

$$= P(X_{t_{k-1}}^{i} \neq [\mathbf{M}] | X_{t_{k}}^{i} = [\mathbf{M}]) \prod_{l=k+1}^{i} P(X_{t_{l-1}}^{i} = [\mathbf{M}] | X_{t_{l}}^{i} = [\mathbf{M}])$$
(D.22)

$$= (1 - P(X_{t_{k-1}}^{i} = [\mathbf{M}] | X_{t_{k}}^{i} = [\mathbf{M}])) \prod_{l=k+1}^{n} P(X_{t_{l-1}}^{i} = [\mathbf{M}] | X_{t_{l}}^{i} = [\mathbf{M}]).$$
(D.23)

- ⁵⁴⁸ Eq. (D.23) can be determined given the sampling method and noise schedule.
- ⁵⁴⁹ For the Euler method, based on Equation Eq. (D.12), we can derive that:

$$P(X_{t_{l-1}}^{i} = [\mathbf{M}] | X_{t_{l}}^{i} = [\mathbf{M}]) = 1 - \sigma(t_{l}) \frac{e^{-\bar{\sigma}(t_{l})}}{1 - e^{-\bar{\sigma}(t_{l})}} (t_{l} - t_{l-1}).$$
(D.24)

550 Therefore, we can express r_k as:

$$r_{k} = (\sigma(t_{k}) \frac{e^{-\bar{\sigma}(t_{k})}}{1 - e^{-\bar{\sigma}(t_{k})}} (t_{k} - t_{k-1})) \prod_{l=k+1}^{n} (1 - \sigma(t_{l}) \frac{e^{-\bar{\sigma}(t_{l})}}{1 - e^{-\bar{\sigma}(t_{l})}} (t_{l} - t_{l-1})).$$
(D.25)

For Tweedie τ -leaping, By Eq. (D.17), similarly we have:

$$P(X_{t_{l-1}}^{i} = [\mathbf{M}] | X_{t_{l}}^{i} = [\mathbf{M}]) = \frac{1 - e^{-\bar{\sigma}(t_{l-1})}}{1 - e^{-\bar{\sigma}(t_{l})}},$$
(D.26)

$$r_{k} = \left(\frac{e^{-\bar{\sigma}(t_{k-1})} - e^{-\bar{\sigma}(t_{k})}}{1 - e^{-\bar{\sigma}(t_{k})}}\right) \prod_{l=k+1}^{n} \left(1 - \frac{1 - e^{-\bar{\sigma}(t_{l-1})}}{1 - e^{-\bar{\sigma}(t_{l})}}\right) = \frac{e^{-\bar{\sigma}(t_{k-1})} - e^{-\bar{\sigma}(t_{k})}}{1 - e^{-\bar{\sigma}(t_{n})}}.$$
 (D.27)

Specifically, if we adopt a log-linear noise schedule, which implies $\bar{\sigma}(t) = -\log(1 - (1 - \epsilon)t)$ and $t_k = \frac{k}{n}$, Equation Eq. (D.27) can be simplified to $\frac{1}{n}$. Substituting this result into Equation Eq. (D.20), we obtain:

E-NFEs
$$(n) = \sum_{k=1}^{n} (1 - (1 - \frac{1}{n})^{l}) = n(1 - (1 - \frac{1}{n})^{l}).$$
 (D.28)

555 E Algorithms for training and inference

556 F Experimental details

557 F.1 Model details

We implemented our RADD model based on SEDD architecture, which is an encoder-only transformer model [46, 47] incorporating time conditioning [48] and using rotary positional encoding [49]. The only difference is that we removed all parts related to time conditioning (i.e. TimeEmbedding, adaLN-zero block [48]) and added a softmax operation at the end of the neural network to ensure the output was a valid conditional distribution. Compared with SEDD small model, this modification led to a reduction of 7M parameters, equating to an 8% decrease from the original 90M non-embedding parameters. Algorithm 1 Unconditional Sampling

Require: Network c_{θ} , noise schedule σ (total noise $\bar{\sigma}$), time range [0, T], step size Δt 1: $t \leftarrow T, \boldsymbol{x}_T \leftarrow [\mathbf{M}] \dots [\mathbf{M}], \boldsymbol{c}_{cache} \leftarrow \boldsymbol{c}_{\theta}(\boldsymbol{x}_t)$ $d \times [\mathbf{M}]$ 2: while t > 0 do 3: if Use Euler then Construct transition densities $p(x_{t-\Delta t}^{i}|\boldsymbol{x}_{t})$ by Eq. (D.12) use \boldsymbol{c}_{cache} 4: 5: end if 6: if Use Tweedie τ -leaping then 7: Construct transition densities $p(x_{t-\Delta t}^{i}|\boldsymbol{x}_{t})$ by Eq. (D.17) use \boldsymbol{c}_{cache} 8: end if $x_{t-\Delta t}^i \sim \operatorname{Cat}(p(x_{t-\Delta t}^i|x_t)) \text{ for all } x_t^i = [\mathbf{M}], x_{t-\Delta t}^i \leftarrow x_t^i \text{ for all } x_t^i \neq [\mathbf{M}]$ if $x_{t-\Delta t} \neq x_t$ then 9: 10: 11: $oldsymbol{c}_{cache} \leftarrow oldsymbol{c}_{ heta}(oldsymbol{x}_t)$ end if 12: 13: $t \leftarrow t - \Delta t$, 14: end while

Algorithm 2 Conditional Sampling

Require: Network c_{θ} , noise schedule σ (total noise $\bar{\sigma}$), time range [0, T], step size Δt , Prompt spaces Ω and tokens \mathcal{T} .

- 1: $t \leftarrow T$, construct \boldsymbol{x}_T with $\boldsymbol{x}_T^{\Omega} = \mathcal{T}$ and $\boldsymbol{x}_T^{\overline{\Omega}} = [\mathbf{M}], \boldsymbol{c}_{cache} \leftarrow \boldsymbol{c}_{\theta}(\boldsymbol{x}_t)$
- 2: while t > 0 do
- 3: if Use Euler then

Construct transition densities $p(x_{t-\Delta t}^i | \boldsymbol{x}_t)$ by Eq. (D.12) use \boldsymbol{c}_{cache} 4:

- 5: end if
- if Use Tweedie τ -leaping then 6:

Construct transition densities $p(x_{t-\Delta t}^{i}|\boldsymbol{x}_{t})$ by Eq. (D.17) use \boldsymbol{c}_{cache} 7:

8: end if

9:
$$x_{t-\Delta t}^i \sim \operatorname{Cat}(p(x_{t-\Delta t}^i | x_t))$$
 for all $x_t^i = [\mathbf{M}], x_{t-\Delta t}^i \leftarrow x_t^i$ for all $x_t^i \neq [\mathbf{M}]$

- 10: if $x_{t-\Delta t} \neq x_t$ then
- 11: $oldsymbol{c}_{cache} \leftarrow oldsymbol{c}_{ heta}(oldsymbol{x}_t)$
- 12: end if
- 13: $t \leftarrow t - \Delta t$,
- 14: end while

Algorithm 3 Training

Require: Network c_{θ} , noise schedule σ (total noise $\bar{\sigma}$), time range [0, T], data distribution p_{data}

- 1: repeat
- $x_0 \sim p_{\text{data}}, t \sim U([0, T]).$ 2:
- 3:
- construct \boldsymbol{x}_t by $Z^i \sim Bernoulli(e^{-\bar{\sigma}(t)}), \boldsymbol{x}_t^i = \mathbb{I}(Z^i = 1)\boldsymbol{x}_0^i + \mathbb{I}(Z^i = 0)[\mathbf{M}]$ Calculate $L_{\theta}(\boldsymbol{x}_t, \boldsymbol{x}_0) = \sum_{\boldsymbol{x}_t^i = [\mathbf{M}]} -\sigma(t) \frac{e^{-\bar{\sigma}(t)}}{1 e^{-\bar{\sigma}(t)}} \log\left(\frac{e^{-\bar{\sigma}(t)}}{1 e^{-\bar{\sigma}(t)}} \boldsymbol{c}_{\theta}(\boldsymbol{x}_t)[i, \boldsymbol{x}_0^i]\right)$ 4:
- 5: Take gradient descent on $\nabla_{\theta} L(x_t, x_0)$
- 6: until converged

Method	RADD-DSE	RADD-DCE
Forward	116.94	113.92
Backward	135.39	125.59
Random	114.94	101.23

Table 3: Quality of unconditionally generated text evaluated by perplexity (\downarrow). For a fixed model, the best perplexity is **bolded**.

565 **F.2 Training details**

- ⁵⁶⁶ Following the settings in [29], we trained our model with the following configuration:
- Batch Size:512
- Learning Rate: 3×10^{-4}
- Exponential Moving Average (EMA):0.9999
- Gradient Clipping: Gradient norm clipped to 1
- Warmup Schedule: Applied for the first 2500 iterations

We utilized 16 V100 32G GPUs or 16 A100 40G GPUs for training. For the A100 40G GPUs, we leveraged flash attention to accelerate the training process. For the V100 32G GPUs, which do not support flash attention or bfloat16, we employed float16 precision and used the Memory-Efficient Attention mechanism available in torch.nn.functional.scaled_dot_product_attention. Additionally, we used gradient checkpointing technique to save memory.

577 F.3 Unconditional generation details

⁵⁷⁸ We used Tweedie τ -leaping method, which has optimal results with fixed NFE. For SEDD small, ⁵⁷⁹ we directly used their result. For RADD small, we generated 1000 samples to get the average value ⁵⁸⁰ following [29].

581 F.4 Further evaluation of generative perplexity

As stated in Theorem 1, c_{θ} can be interpreted as a conditional distribution over clean data. A natural idea is to use it directly to generate samples, which is similar to auto-regressive models. However, there are d! kinds of decomposition from joint distribution to conditional distribution, in which we only tested three representative cases:

586 • forward:
$$p(x^1 \cdots x^d) = \prod_{k=1}^d p(x^k | x^{($$

• backward:
$$p(x^1 \cdots x^d) = \prod_{k=1}^d p(x^k | x^{(>k)})$$

• random:
$$\pi \sim U(S_d), p(x^1 \cdots x^d) = \prod_{k=1}^d p(x^{\pi(k)} | x^{\pi($$

Results are shown in Table 3. The perplexity is calculated on average of 1024 samples. For the random case, we calculate the average perplexity between different randomly generated π . Generally, we find that the perplexity by directly sampling from the conditional distribution is higher than that achieved by Tweedie τ -leaping. Among the different decomposition orders, the random order demonstrated the best performance.

594 G Additional experimental results

595 G.1 Additional samples

58

In this section, we present the unconditionally and conditionally generated text of RADD-DSE in

Fig.3 and Fig.4, respectively. Similarly, the results of RADD-DCE are shown in Fig.?? and Fig.??, respectively.

and human face. "And pretty damn conventional mating, so didn't come to the table like that with me. (As a character), it would be a pretty solid case to have," Andra says. "I saw the way he did it, and I think played a little bit with some of his fans. He wanted me to get cute a little bit too."

Advertisement

As in the parking lot, Andra was growing frustrated with the way the cars recently fit in nicely.

It also makes me feel like some things haven't changed until around this period of time, in the future. It's definitely the future here now — and I'm always dubious about thinking quite long before Toyota ever introduced a new car. "I think something that perfectly well fits all the guys," Andra says. "I therefore could fit more well."

- Follow Matt Dyckton on Twitter @Mittington.<|endoftext|>Spanish star Christina Rene got away after a teenage girl told to leave India for an unfamiliar place.

The Russian woman chose Chelsea to stay home, saying the alternative was to get rid of a condition and die after receiving a medical diagnosis and go back to learning as a nurse.

Chelsea was sentenced to the first arrest for a serious mental health conviction in October, and was transferred to a man who had taken his place, a 16-year-old man. The bailee's Appeal Court had appealed to a court to hear the case she learned from the teenagers.

The Russian Internal Revenue Service found the woman charged with administering an emergency ward, and although the girl was still waiting for a doctor there, she instead went to visit another clinic for the treatment of female swi-virus virus.

The court had not ruled out an in-life medical professional. Chelsea never applied to be a pregnant mother; her formal application assumed she was pregnant, dating from the summer of 2014.

Out of sound doubts, when her co-boyfriend Chelsea received a proposal to stay abroad for summer work abroad. Chelsea then applied to win to work and have a place in Russia.

The grant of nearly \$5000 Chelsea invited Chelsea out to go see a Moscow clinic. It took two weeks to find a doctor. These sors' of Moscow's federal courts confiscated the grant.

The 17-year-old did not need to go to the hospital where she says she has received everything she has had in China, of course. Chelsea's lawyers Andre James Irani gave the court the doctors he could plead without permission of the young teenager on a regular basis.

The Argentine was put out on bail, prompting a sex psychologist to meet her when she signed her papers, but no family member was present.

"First thing I wasn't going to ask three days, but I'm thinking about this months already. "I'm glad to see that they are waiting for her with their services. She's already built a good life. She's interested in her studies. But feel like she is? I want to be the only person who has ever been my friend," she said.

"I took a lot of act, but she is more than never."

To this day Christina Rene McCourner told the court the situation is between Denmark and Stockholm syndrome. She says Chelsea has refused to come to term.

"She says that she is talented at medical school. But she's telling a different story, anyway. There's also a case when you just can't get the CNN-type cowardice yourself going to the doctor," she said.

A team of Barcelona has been conducting medical inquiries into both doctors who treat the sick and those not who use medical services. "How has she been following her visit to Russia United States and since arriving unable to do so?"

"Each often when she said "All the bills are what may I doctor," I've usually never used my medicine again," Chelsea attorneys say.

Life in the wrong world

While far-so in the past two weeks though, nearly speaking only three pregnant women in Mexico and Europe have requested that Chelsea stop using medical services at them all. Though many of its applicants have receiving medical assistance in Qatar and other parts of the world, they believe it means they should go elsewhere, according to the health services agency.

Chelsea's received from facility staff dealing with financial responsibility say it is those who are vulnerable have been out of money for healthcare their life, not who are suffering from a lack of access. "You might think after you tell that provider is not available, you should give the load to somebody in Russia or Sweden," she

Figure 3: Unconditionally generated text of RADD-SE.

Hi, my name is Shade-Rayhelynis-Neelsons. Interviewer: Two days ago, so let me speak to you in brief.Drake: Hi, are you studying for graduate school in late July, and somebody is interviewing you for you for his classes. My first personal quote is: so when they're doing class they're going to a shower, and they expect me to not be part of the shower any way. I have a take, I want to check when something's not right and make sure I spot it so that someone can get it happening. My hair is an important piece of me, and I can be the one complete human being that when I have a shower and say, "Okay, I really like my hair, and this is the thing that I would like and I want to believe something, I should just have an attitude check." Do you think anyone else can have time with me? Do I say I don't?I mean I get to 7 right the time I go to class, I sit back there and I feel like no one knows what to do. So, in my case I am not anxious, I'm just saying, "I feel like I have my hair without letting go from day to day, I've got to feel like I've been like the thing they should all be excited by."So, the days follow me, you start looking at my pictures, and you realize how pretty you are. It was just so big this moment for that office, because that was seeing positive things.I'm starting to grow up and be beautiful too.I mean-I mean it's kind of fun, now seeing beautiful girls, especially really in their 20s, come from unusual backgrounds. Back in the 50s I met people at one of the first places Woosz made mons, asked us to visit charity walks, and he would buy us a suit. And actually he really wanted to, so I know what the inspiration is. It's something that someone could have found out, and where they're from - I think it is like everyone is finding ways to relate. This is one of those things I remember the most about when I went to order the products!For example, it may be taken after the sporting event, and I don't know, but they asked us to choose color for our favorite parts in the group....they asked us to take their favorite color, then they made their eves for each skin. The one that was the round head, like that one I chose so much, but it was really a really painful transition...and I don't know what touched my skin, so I don't know what I would do with it. I will tell you...that was not what I focused on, but I think I may have pretty much my original o-still darker hair, so I chose to go with the round head which I did really. Anyway, my job was to get the hair done, and I do most of my hair for school because I got married and had so many children while I was So. So, I'm always thinking again not to be confused. I had just been doing niggly since I was a child, and there was no reason to do it like that. The career had got started, and therefore I was really going all over the world no less. So one of the things I did was taking of his shirts so I look at him all in one go, I surprise him. I've got some of the most cool things going with those two of this. You see these all look alike, clean, and awesome-and they used to have a shop in the....a shop like that, they never sells coats, so that was really helpful to those guys, I feel, guys going into shop with these guys, and those guys can have the right color, they can have it that right look. I love my hair, it's like when I was a year old when I see scent in your nursery, to get the grandest response, I keep using the smell. I bring it to its level, honestly, it makes the weight apply.rake: Probably my favorite word to explain, after you explain, it isn't hard-in case that's your thing, I'm already who kind of like going, and I'm able to get access when I'm trying. I've been trying to finish coloring for people, for years, so I would like to work on it's own but the colors there to do color pattern are picked for another reason to put the final color off, that's like the other one for giving glitter or outline. So I'm just going to be looking side by side at coloring over and over again. It's my vision. Each color is my dream. As you know, a color is just simply something tucked into a dye and what makes perfect or tail end to a hair is...and that is why I always shampoo twice a day and shower three times a day.

Figure 4: Conditionally generated text of RADD-SE. Prompt tokens are highlighted in blue.

as well. However, she did not sign.

Gov. Johnluaj said the Amal effectively took the case to the Supreme Court.

"I was wrong to say that they left the Constitution in place and this is basically unconstitutional," Jackson said. "I don't think they're saying that. I'm. That means on the one hand they're going to have to intervene, or on the side of the other, they'll have to intervene. I can't think of changing a constitutional decision with them."

When Gov. Barack Obama announced Tuesday — and a federal court is set to challenge the way residents of the states violated the law — Attorney General Eric Holder chose to hold his own hearing. The entire state had a deadline to field a recipient of the letter asking for comment.

Because hearings were held before, Ohio and 17 states have each had such a case before.

At last night's hearings, more than four separate arguments were heard by a 52-45 margin in order to pass the repeal bill by a vote of 63-2 49 to 45. The bill also included Obamacare legislation and was pushed through Congress after opposition from 24 states. Both brought in a new governor, popular Gov Sen. Phil Bryant, another Republican in the Senate. Neil LePage neared an attorney in both cases and refused to find a new insurance secretary.

Both House Leader Mitch McConnell and Republicans said they would repeal the law entirely. The law would have been in the Oval Office of the Logged since 2011.

Former Judge Anthony Teague, the Ohio Chief Judge, found the overwhelming majority vote in favor of the legislation well in line and said it was a "needed forward."

"The things citizenship issues should go from statutes states have to regulations," he said. "They've got this idea that the courts are getting knocked to their own corner. And once they see it turn around them, sure as hell they'll have faith dogged by judges exercising constitutional rights."

Attorney Holder, the assistant secretary of state for policy at the Department of Justice, has discussed the idea that federal courts such as the U.S. Supreme Court should handling legal issues such as making tough immigration decisions for illegal immigrants.

"The thought process of helping illegal immigrants goes beyond the judicial process. What I understand... criminal immigration measures, gang activity, affirmative action efforts, criminal status," Holder said. "And things like that, we expect in the state to go to a long way. There clearly needs to be a criminal justice program on immigration and reform, and we need to recognize those efforts not to start."

Holder has also said it is a "real issue" for the constitution if the letter is signed off. He said it was essentially a message, seeking to show the majority "power of respect and the power of institutions."

But the attorney general said he hoped it would be a task to figure out how to start safeguarding each of their citizens' rights while restoring their constitutional trust.

"The only thing I think we can really do is have the judges to understand the nature of the judiciary, how powerful it is to be involved and to interfere with the government with no accountability on what path they want to move down," he said.<|endoftext|>Ex76561 molds at the worst parts of the UK economy will tell in the future – most poorer areas of England have suffered more than any person during the first years of 2008, according to ex-Home Secretary Jeremy Foot.

A total of 23,000 people of disposable incomes who have five mortgages – more than 4,000 households – will be considered as home buyers, even though published figures will be different from September, say researchers

Residents of more than 3,000 homes will be the third most likely to die because of jobs which fell in average terms home ownership during 2008 until the end of this year.

It's a surprising drop, according to the Wall Street Journal, which says the economy and the top 1

Professor Jeremy Foot, Home Secretary and the Information Society, raised concerns about the collapse of the UK's housing market – down from 39

He said: "The measure of who 'invested' at that time, doesn't include the number of people or businesses with the assets. There were only three big cities, New York were the other three?

"This year – it was announced that Royal Bank of Scotland would be the first to close in 20 years – it turned out then that the poorer areas, including by 2009 and 2010, were hammered hardest," Professor Foot said.

Professor Foot's lecture, which was released today in

Figure 5: Unconditionally generated text of RADD-CE.

I have a ick of death — that's where I am. That's what I need. That need to take something else. But I don't have that in my family. And here's why you might let it go. You don't want to really say what I wrote in my story, open your mind and finish it with purposeful thought. But I can do it if I get cancer."

"I can not help but accept that you brother-in-law was really on me and that that's not how I need to be," I said with my father's sad smile. He let it go a few days later, but it didn't prevent him from thinking about it. But I made sure he wouldn't let-in-law leave.

My mother saved my only father for the life. And he was a falsehood, but no doubt. He saved me.

"My mother thank god. But do you ask outside of me the questions, ask him. He's dealing with this and prepare for their perception of bad things he said. They should't forget him because the bad thoughts come with it."

I blackened my father and watched them walk back up to him in front of a good news line. He stood with us for two days; we stared on from spring to spring.

"It's OK." He said.

"I know what you want when you come here."

"I'll accept that," he said. "What I thought to look back was bad. I'm going to let it go and make an exception in this case. I came here after sharing the story. I've been looking to the world. And I don't care for anyone in this world, but I want to make it harder for other people to try to make this mistake."

He didn't justify what he meant. My mother drove us the way I trucked him back to his house, and he chuckled at me in simple words:

"I said to him I need to still be here, but do you want me to have this for the whole day?" "These are my expectations. If we smile and I'm laughing thank you can I be so happy?"

Outside of my family, there was joy, joy, sorrow. At youngest, most of all, my father also saw the world. He suffered from hunger, his family lost access to wheelchair and every room. They fought with him sometimes, and when they needed to rest, he stepped back. I can only imagine, deep down, his family grew enormously with him.

Our journey with my parents at times had been a combination of things. We were all in pairs, and my brothers were such. Repeatedly we asked them to tell me when we wouldn't last to get that end. And I started missing my trip — never visited neighbors before. And finally one day I lost my patience. I all wanted to wander down to my father's house in76561. On my own, I couldn't see my beautiful mother again. I had an ear tumor that was brainblown away and lost all the ways to make the most of the time again. I took the bed from the entire bear family along with his three children, and fixed him up on one side and led him home through the open of the front door with my honey brother. We were hearing all of the other things I had heard were outlandish stories. Finally, he remained so grateful for his admiration for my dad.

What was possible my father really didn't have?

On November 2, 2013, my mother still did not show any empathy for my father. The loyalty he had was the only expression he had.

On the night where he was arrested, other members of his family noticed that they couldn't bribe a co-worker to walk up from his job. That led to him thinking, "I want you to not be a badexample in your family. You can't serve an individual who will make up his stupid demands for you in order to make a profit." For me, to fulfill the compassion of my own brother, I've become an able example of that, that kids.

With this decision, it was freeing that I had allowed him to try to decide how to make a better life for his family. I had allowed the animals to have a day off. This is something he and I can do. His son had thought they could, and it was relief he had been rearranging that reality. The time our family spends on the house was at seven times a day.

Figure 6: Conditionally generated text of RADD-CE. Prompt tokens are highlighted in blue.

599 H License

⁶⁰⁰ URL and license for existing assets we used are provided in Table 4.

Name	URL Licens
SEDD	https://github.com/louaaron/Score-Entropy-Discrete-Diffusion MIT Licens
NT T	
Neurii	78 Paper Checklist
1.	Claims
	Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?
	Answer: [Yes]
	Justification: The main claims made in the abstract and introduction accurately reflect our paper's contributions and scope.
	Guidelines:
	• The answer NA means that the abstract and introduction do not include the claims made in the paper.
	 The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
	• The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
	• It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.
2.	Limitations
	Question: Does the paper discuss the limitations of the work performed by the authors?
	Answer: [Yes]
	Justification: We discuss the limitations in the main text.
	Guidelines:
	• The answer NA means that the paper has no limitation while the answer No means that
	the paper has limitations but those are not discussed in the paper
	• The authors are encouraged to create a separate "Limitations" section in their paper.
	• The paper should point out any strong assumptions and how robust the results are to
	violations of these assumptions (e.g., independence assumptions, noiseless settings,
	model well-specification, asymptotic approximations only holding locally). The authors
	should reflect on how these assumptions might be violated in practice and what the
	implications would be.
	• The authors should reflect on the scope of the claims made, e.g., if the approach was
	only tested on a new datasets of with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated
	• The authors should reflect on the factors that influence the performance of the approach
	For example, a facial recognition algorithm may perform poorly when image resolution
	is low or images are taken in low lighting. Or a speech-to-text system might not be
	used reliably to provide closed captions for online lectures because it fails to handle
	technical jargon.
	• The authors should discuss the computational efficiency of the proposed algorithms
	and how they scale with dataset size.
	• If applicable, the authors should discuss possible limitations of their approach to
	address problems of privacy and fairness.
	• While the authors might fear that complete honesty about limitations might be used by
	limitations that aren't acknowledged in the paper. The authors should use their best
	iudgment and recognize that individual actions in favor of transparency play an impor-
	tant role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

Table 4:	URL and	license	for	existing	assets	we	used	•
----------	---------	---------	-----	----------	--------	----	------	---

649	3.	Theory Assumptions and Proofs
650 651		Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?
652		Answer: [Yes]
653 654		Justification: For each theoretical result, the paper provides the full set of assumptions and a complete (and correct) proof. Please see the Appendix for more details.
655		Guidelines:
656		• The answer NA means that the paper does not include theoretical results.
657 658		• All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
659		• All assumptions should be clearly stated or referenced in the statement of any theorems.
660		• The proofs can either appear in the main paper or the supplemental material, but if
661		they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition
663		• Inversely any informal proof provided in the core of the paper should be complemented
664		by formal proofs provided in appendix or supplemental material.
665		• Theorems and Lemmas that the proof relies upon should be properly referenced.
666	4.	Experimental Result Reproducibility
667		Question: Does the paper fully disclose all the information needed to reproduce the main ex-
668		perimental results of the paper to the extent that it affects the main claims and/or conclusions
669		of the paper (regardless of whether the code and data are provided or not)?
670		Answer: [Yes]
671 672		Justification: Our paper fully discloses all the information needed to reproduce the main experiment.
673		Guidelines:
674		• The answer NA means that the paper does not include experiments.
675		• If the paper includes experiments, a No answer to this question will not be perceived
676 677		well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
678		• If the contribution is a dataset and/or model, the authors should describe the steps taken
679		to make their results reproducible or verifiable.
680		• Depending on the contribution, reproducibility can be accomplished in various ways.
681		For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation it may
683		be necessary to either make it possible for others to replicate the model with the same
684		dataset, or provide access to the model. In general, releasing code and data is often
685		one good way to accomplish this, but reproducibility can also be provided via detailed
686		instructions for how to replicate the results, access to a hosted model (e.g., in the case
687		of a large language model), releasing of a model checkpoint, or other means that are
688		appropriate to the research performed.
689		• While NeurIPS does not require releasing code, the conference does require all submis-
690 691		nature of the contribution. For example
692		(a) If the contribution is primarily a new algorithm the paper should make it clear how
693		to reproduce that algorithm.
694		(b) If the contribution is primarily a new model architecture, the paper should describe
695		the architecture clearly and fully.
696		(c) If the contribution is a new model (e.g., a large language model), then there should
697		either be a way to access this model for reproducing the results or a way to reproduce
698		the dataset)
700		(d) We recognize that reproducibility may be tricky in some cases, in which case
701		authors are welcome to describe the particular way they provide for reproducibility.

702 703 704		In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.
705	5.	Open access to data and code
706 707		Question: Does the paper provide open access to the data and code, with sufficient instruc- tions to faithfully reproduce the main experimental results, as described in supplemental
708		material?
709		Answer: [Yes]
710 711		Justification: We provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental materials.
712		Guidelines:
713		• The answer NA means that paper does not include experiments requiring code.
714 715		• Please see the NeurIPS code and data submission guidelines (https://nips.cc/ public/guides/CodeSubmissionPolicy) for more details.
716 717 718 719		• While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark)
720 721 722		 The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
723 724		• The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
725 726 727		• The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
728 729		• At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
730 731		• Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.
732	6.	Experimental Setting/Details
733 734 735		Question: Does the paper specify all the training and test details (e.g., data splits, hyper- parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?
736		Answer: [Yes]
737 738		Justification: We specify all the training and test details necessary to understand the results in Section 4.1
739		Guidelines:
740		• The answer NA means that the paper does not include experiments.
741		• The experimental setting should be presented in the core of the paper to a level of detail
742		that is necessary to appreciate the results and make sense of them.
743 744		• The full details can be provided either with the code, in appendix, or as supplemental material.
745	7.	Experiment Statistical Significance
746 747		Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?
748		Answer: [No]
749		Justification: Error bars are not reported because it would be too computationally expensive.
750		Guidelines:
751		• The answer NA means that the paper does not include experiments.

752 753 754		• The authors should answer "Yes" if the results are accompanied by error bars, confi- dence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
755 756 757		• The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
758 759		• The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
760		• The assumptions made should be given (e.g., Normally distributed errors).
761		• It should be clear whether the error bar is the standard deviation or the standard error of the mean
762		of the filter.
763 764 765		preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
766 767 768		• For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
769 770		• If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.
771	8.	Experiments Compute Resources
772		Question: For each experiment, does the paper provide sufficient information on the com-
773		puter resources (type of compute workers, memory, time of execution) needed to reproduce
774		the experiments?
775		Answer: [Yes]
776 777		Justification: We provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments.
778		Guidelines:
779		• The answer NA means that the paper does not include experiments.
780 781		• The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
782 783		• The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
784 785 786		• The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).
787	9.	Code Of Ethics
788 789		Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?
790		Answer: [Yes]
791 792		Justification: We conduct in the paper conform, in every respect, with the NeurIPS Code of Ethics.
793		Guidelines:
794		• The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
795		• If the authors answer No, they should explain the special circumstances that require a
796		deviation from the Code of Ethics.
797 798		• The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).
799	10.	Broader Impacts
800 801		Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?
802		Answer: [Yes]

803 804	Justification: We discuss the potential positive societal impacts and negative societal impacts in Section 6 of the main text.
805	Guidelines:
806	• The answer NA means that there is no societal impact of the work performed.
807	• If the authors answer NA or No, they should explain why their work has no societal
808	impact or why the paper does not address societal impact.
809	• Examples of negative societal impacts include potential malicious or unintended uses
810	(e.g., disinformation, generating fake profiles, surveillance), fairness considerations
811	(e.g., deployment of technologies that could make decisions that unfairly impact specific
812	groups), privacy considerations, and security considerations.
813	• The conference expects that many papers will be foundational research and not tied
814	to particular applications, let alone deployments. However, if there is a direct path to
815	any negative applications, the authors should point it out. For example, it is legitimate
816 917	generate deepfakes for disinformation. On the other hand, it is not needed to point out
818	that a generic algorithm for optimizing neural networks could enable people to train
819	models that generate Deepfakes faster.
820	• The authors should consider possible harms that could arise when the technology is
821	being used as intended and functioning correctly, harms that could arise when the
822	technology is being used as intended but gives incorrect results, and harms following
823	from (intentional or unintentional) misuse of the technology.
824	• If there are negative societal impacts, the authors could also discuss possible mitigation
825	strategies (e.g., gated release of models, providing defenses in addition to attacks,
826	mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
827	feedback over time, improving the efficiency and accessibility of ML).
828	11. Safeguards
829	Question: Does the paper describe safeguards that have been put in place for responsible
830	release of data or models that have a high risk for misuse (e.g., pretrained language models,
831	image generators, or scraped datasets)?
832	Answer: [NA]
833	Justification: This paper poses no such risks.
834	Guidelines:
835	 The answer NA means that the paper poses no such risks.
836	• Released models that have a high risk for misuse or dual-use should be released with
837	necessary safeguards to allow for controlled use of the model, for example by requiring
838	that users adhere to usage guidelines or restrictions to access the model or implementing
839	safety filters.
840	• Datasets that have been scraped from the Internet could pose safety risks. The authors
841	• We recognize that may iding offective soforwards is shellonging, and many papers do
842	• We recognize that providing effective safeguards is chantenging, and many papers do not require this, but we encourage authors to take this into account and make a best
844	faith effort.
845	12. Licenses for existing assets
946	Question: Are the creators or original owners of assets (e.g., code, data, models) used in
847	the paper, properly credited and are the license and terms of use explicitly mentioned and
848	properly respected?
849	Answer: [Yes]
850	Justification: URL and license for existing assets we used are provided in Appendix H.
851	Guidelines:
852	• The answer NA means that the paper does not use existing assets
853	• The authors should cite the original paper does not use existing assess.
854	• The authors should state which version of the asset is used and if possible include a
855	URL.

856		• The name of the license (e.g., CC-BY 4.0) should be included for each asset.
857 858		• For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
859		• If assets are released, the license, copyright information, and terms of use in the
860		package should be provided. For popular datasets, paperswithcode.com/datasets
861		has curated licenses for some datasets. Their licensing guide can help determine the
862		license of a dataset.
863		• For existing datasets that are re-packaged, both the original license and the license of
864		the derived asset (if it has changed) should be provided.
865		• If this information is not available online, the authors are encouraged to reach out to
866		the asset's creators.
867	13.	New Assets
868		Question: Are new assets introduced in the paper well documented and is the documentation
869		provided alongside the assets?
870		Answer: [Yes]
871		Justification: new assets introduced in this paper are well documented and provided alongside
872		the assets.
873		Guidelines:
874		• The answer NA means that the paper does not release new assets.
875		• Researchers should communicate the details of the dataset/code/model as part of their
876		submissions via structured templates. This includes details about training, license,
877		limitations, etc.
878		• The paper should discuss whether and how consent was obtained from people whose
879		asset is used.
880		• At submission time, remember to anonymize your assets (if applicable). You can either
881		create an anonymized URL or include an anonymized zip file.
882	14.	Crowdsourcing and Research with Human Subjects
883		Question: For crowdsourcing experiments and research with human subjects, does the paper
884		include the full text of instructions given to participants and screenshots, if applicable, as
885		well as details about compensation (if any)?
886		Answer: [NA]
887		Justification: This paper does not involve crowdsourcing nor research with human subjects.
888		Guidelines:
889		• The answer NA means that the paper does not involve crowdsourcing nor research with
890		human subjects.
891		• Including this information in the supplemental material is fine, but if the main contribu-
892		tion of the paper involves human subjects, then as much detail as possible should be
893		included in the main paper.
894		• According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
895		or other labor should be paid at least the minimum wage in the country of the data
896		collector.
897	15.	Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
898		Subjects
899		Question: Does the paper describe potential risks incurred by study participants, whether
900		such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
901		approvals (or an equivalent approval/review based on the requirements of your country or
902		institution) were obtained?
903		Answer: [NA]
904		Justification: This paper does not involve crowdsourcing nor research with human subjects.
905		Guidelines:
906 907		• The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

908	• Depending on the country in which research is conducted, IRB approval (or equivalent)
909	may be required for any human subjects research. If you obtained IRB approval, you
910	should clearly state this in the paper.
911	• We recognize that the procedures for this may vary significantly between institutions
912	and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
913	guidelines for their institution.
914	• For initial submissions, do not include any information that would break anonymity (if
915	applicable), such as the institution conducting the review.