ShortListing Model: A Streamlined Simplex Diffusion for Discrete Variable Generation

Yuxuan Song 123* Zhe Zhang 13* Yu Pei 13* Jingjing Gong 1 Qiying Yu 12 Zheng Zhang 2 Mingxuan Wang 2 Hao Zhou 1 Jingjing Liu 1 Wei-Ying Ma 1 Institute of AI Industry Research (AIR), Tsinghua University 2 ByteDance Seed 3 Department of Computer Science and Technology, Tsinghua University {yxsong0816, zhezhan22, yupei.wp}@gmail.com {zhouhao, maweiying}@air.tsinghua.edu.cn

Abstract

Generative modeling of discrete variables is challenging yet crucial for applications in natural language processing and biological sequence design. We introduce the Shortlisting Model (SLM), a novel simplex-based diffusion model inspired by progressive candidate pruning. SLM operates on simplex centroids, reducing generation complexity and enhancing scalability. Additionally, SLM incorporates a flexible implementation of classifier-free guidance, enhancing unconditional generation performance. Extensive experiments on DNA promoter and enhancer design, protein design, character-level and large-vocabulary language modeling demonstrate the competitive performance and strong potential of SLM. Our code can be found at https://github.com/GenSI-THUAIR/SLM.

1 Introduction

Autoregressive models such as large language models (LLMs), although having achieved remarkable success in text generation [Achiam et al., 2023, Brown et al., 2020], struggle over tasks that lack an intrinsic sequential-ordering inductive bias, such as DNA design [Avdeyev et al., 2023, Stark et al., 2024], protein sequence design [Wang et al., 2024, Lin et al., 2023] and molecular graph generation [Vignac et al., 2022]. Consequently, there is growing interest in developing new paradigms for discrete variable generation, such as diffusion-based [Lou et al., Sahoo et al., 2024, Shi et al., 2024] and flow-matching-based [Gat et al., 2024, Davis et al., 2024, Cheng et al., 2024] methods.

Recent discrete generative models are generally classified into two main categories based on their operational spaces: discrete-space models and continuous-space models. The former, specification of the continuous space in the former of the continuous space.

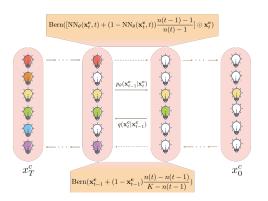


Figure 1: SLM's forward and reverse process. Comparisons between MDLM and DP3M-Uniform is located in Appendix B.1

cally discrete diffusion model [Lou et al., Xu et al., 2024], mimics continuous diffusion processes using substitution or masking operations to decompose information, which has shown impressive

^{*}Equal Contribution. Correspondence to Hao Zhou (zhouhao@air.tsinghua.edu.cn).

performance for discrete generative modeling. However, these discrete counterparts differ fundamentally from original continuous diffusion models [Ho et al., 2020], where the smooth and gradual information transitions intrinsic to the generation process are key to their success. On the other hand, continuous-space models map discrete data into continuous representations. While benefiting from continuous properties, capturing the geometric structure and adhering to constraints via continuous generative models introduces new challenges. In this context, simplex-based approaches [Cheng et al., 2024, Davis et al., 2024, Graves et al., 2023] offer a balanced solution by representing discrete data on the probability simplex, which naturally adheres to the fundamental properties of categorical distributions.

However, existing simplex-based approaches often rely on intricate operations to define trajectories over the entire continuous space. For instance, Statistical Flow Matching(SFM) [Cheng et al., 2024] and Fisher Flow Matching [Davis et al., 2024] define geodesics based on sphere map and the Fisher-Rao metric. The training also incorporates Riemannian optimal transport. Similarly, Bayesian Flow Networks (BFNs) involve heavy mathematical derivations and change-of-variable techniques to define simplex trajectories through Gaussian-formed count variables. Despite mathematical rigor, their demanding complexity limits scalability in large-scale generative tasks.

In this paper, we aim to preserve the core principle of simplex-based methods, *gradual information growth*, while exploring simpler yet effective alternatives. We propose viewing discrete variable generation as progressive candidate pruning, starting from the full category set and iteratively narrowing down to a single choice. We term this approach **ShortListing Models**(SLM). Formally, shortlisting models reside within the diffusion framework, trainable via the variational lower bound (VLB). In contrast to existing simplex-based methods using vocabulary-level MSE loss [Graves et al., 2023, Cheng et al., 2024, Davis et al., 2024], our approach employs a simplified cross-entropy objective, effectively mitigating vanishing gradient issues and better handling large-vocabulary settings. As illustrated in Figures 1 and 2, our model effectively reduces degrees of freedom by modeling transitions among the simplex centroids rather than the entire simplex. Additionally, shortlisting models offer flexibility for adaptations such as classifier-free guidance.

We comprehensively evaluate the proposed approach over various discrete generation tasks and benchmarks, including char-level language modeling, large-scale language modeling, DNA promoter and enhancer design, and Protein sequence design. Specifically, we achieve strong performance among non-autoregressive methods on text8 and obtain competitive results on OpenWebText, where previous simplex-based approaches had difficulty generating reasonable outputs. In DNA design tasks, our non-guided variants achieve state-of-the-art (SOTA) results. Furthermore, with classifier-free guidance, our model attains stronger performance while remaining less sensitive to hyperparameters. Additionally, our 38M-parameter shortlisting model can design proteins with enhanced foldability, fitness, self-consistency and diversity, surpassing the larger, well-known ESM2-150M model [Lin et al., 2022].

2 Preliminary

2.1 Definition and Notations

We encode a discrete variable with K distinct categories using one-hot vectors $\mathbf{e} = [e_1, e_2, \dots, e_K]^T$. In each vector, only the i-th element $\mathbf{e}(i) = 1$ signifies the inclusion of the i-th category, while all other elements are zero.

Definition 2.1. A **candidate set** for K categories is defined as a binary-valued vector $\mathbf{c} = [c_1, c_2, \dots, c_K]^T$, where each $c_i \in \{0, 1\}$, and the vector has at least one non-zero entry, *i.e.*, $\mathbf{1}^T \mathbf{c} > 0$.

The candidate-set variable \mathbf{c} encodes the selection status of each category: $\mathbf{c}(i)=1$ indicates that the i-th category is included, while $\mathbf{c}(i)=0$ denotes its exclusion. Specifically, there are two distinct instances of the candidate variable: \mathbf{c} is an all-one vector $([1,\cdots,1])$ that represents maximum candidates setting, i.e. all K categories are selected; one-hot vector is another special case which represents minimum candidates setting, with only one category included.



Figure 2: Pathological behavior of SLM on one simplex with $K = 5(\Delta^4)$. Each vertex represents one of the categorical targets while the trajectory of the white point serves as a probability path in sampling. Note that the trajectory of shortlisting model could be seen as jumping among the centroids of subspaces in simplex space.

2.2 Diffusion Models

Shortlisting model is a variant of diffusion models, which can be viewed as latent variable models where the latent variables form a Markov chain. Specifically, for a diffusion model with the sequence latent variable $\mathbf{x}_{1:T} = \mathbf{x}_1, \cdots \mathbf{x}_T$, the implied density function p_{θ} holds the following Markovness by definition: $p_{\theta}(\mathbf{x}_0, \mathbf{x}_{1:T}) = p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)p_{\theta}(\mathbf{x}_1|\mathbf{x}_2)\cdots p_{\theta}(\mathbf{x}_{T-1}|\mathbf{x}_T)$. To learn this latent variable model, a carefully designed constant variational distribution $q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$, referred to as the 'forward process', is involved. Based on the variational distribution, the diffusion model is generally trained with the following variational lower bound [Austin et al., 2021, Ho et al., 2020]:

$$L_{\text{vlb}} = \mathbb{E}_{q(\boldsymbol{x}_{0})} \underbrace{\left[D_{\text{KL}}\left[q\left(\boldsymbol{x}_{T} \mid \boldsymbol{x}_{0}\right) \| p\left(\boldsymbol{x}_{T}\right)\right]}_{L_{T}} + \sum_{t=2}^{T} \mathbb{E}_{q(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})} \left[D_{\text{KL}}\left[q\left(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t}, \boldsymbol{x}_{0}\right) \| p_{\theta}\left(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t}\right)\right]\right]}_{L_{t-1}} - \mathbb{E}_{q(\boldsymbol{x}_{1} \mid \boldsymbol{x}_{0})} \left[\log p_{\theta}\left(\boldsymbol{x}_{0} \mid \boldsymbol{x}_{1}\right)\right]\right]. \tag{1}$$

Here $q(x_0)$ refers to the data distribution.

3 Shortlisting Models

Inspired by progressive candidate pruning, we effectively translate the generation of discrete variables into a category selection process, which begins by considering all categories in the vocabulary as potential candidates and progressively narrows down the options until reaching a final one-hot representation, indicating a single category. This section introduces the detailed components of the proposed shortlisting model as well as the training and sampling processes.

3.1 Forward Candidate Appending

We design a forward candidate appending process over the space of **candidate set** as introduced in Definition. 2.1. For any discrete variable \mathbf{x} , its one-hot representation (as a special case of the candidate set) is considered as the initial step, denoted as $\mathbf{x}_0^{\mathbf{c}}$. For the last step, $\mathbf{x}_T^{\mathbf{c}}$, we make it into an all-one vector (1). Then we seek the following Markov chain interpolation $\mathbf{x}_0^{\mathbf{c}} - \mathbf{x}_1^{\mathbf{c}} - \cdots - \mathbf{x}_T^{\mathbf{c}}$ between $\mathbf{x}_0^{\mathbf{c}}$ and $\mathbf{x}_T^{\mathbf{c}}$, which satisfies:

$$\forall_{0 \le i < j \le T} \ \mathbf{1}^T \mathbf{x}_i^{\mathbf{c}} \le \mathbf{1}^T \mathbf{x}_j^{\mathbf{c}}, [\mathbf{x}_j^{\mathbf{c}}]^T \mathbf{x}_i^{\mathbf{c}} = \mathbf{1}^T \mathbf{x}_i^{\mathbf{c}}$$
(2)

Recall $\mathbf{x^c}$ is a binary-valued vector, hence the above condition essentially indicates the possible categories implied by the candidate set of early steps are **strictly scooped** by later steps. We propose using a multivariate Bernoulli distribution to model the forward process over the candidate-set variable, denoted as $\mathbf{x^c} \sim \mathrm{Bern}(\phi)$, where ϕ is a K dimensional vector representing the parameters of the distribution. To control the noise level, we introduce n(t) as a scheduling function over the candidate numbers, where $1 \leq n(t) \leq K$. By our definition, n(t) is a monotonically increasing function from time step 0 to T, designed to gradually perturb the signal. Intuitively, n(t) can be viewed as controlling the number of ones in the vector $\mathbf{x_t^c}$, representing the number of possible categories at time t. To satisfy the condition in Eq. 2, we set n(0) = 1 and n(T) = K, and define the

transition probabilities from t-1 to t as:

$$q(\mathbf{x}_{t}^{\mathbf{c}}|\mathbf{x}_{t-1}^{\mathbf{c}}) = \operatorname{Bern}\left(\mathbf{x}_{t-1}^{\mathbf{c}} + (1 - \mathbf{x}_{t-1}^{\mathbf{c}}) \frac{n(t) - n(t-1)}{K - n(t-1)}\right)$$
(3)

Proposition 3.1. With Eq. 3 as the transition probability, the marginal distribution is defined as:

$$q(\mathbf{x}_{t}^{\mathbf{c}}|\mathbf{x}_{0}^{\mathbf{c}}) = \operatorname{Bern}\left(\mathbf{x}_{t-1}^{\mathbf{c}} + (1 - \mathbf{x}_{t-1}^{\mathbf{c}})\frac{n(t) - 1}{K - 1}\right)$$
(4)

and corresponding posterior distribution $q(\mathbf{x}_{t-1}^{\mathbf{c}}|\mathbf{x}_{t}^{\mathbf{c}},\mathbf{x}_{0}^{\mathbf{c}})$ also lies in the form of Bernoulli distribution, the analytical form of which is $(t \geq 2)$:

$$q(\mathbf{x}_{t-1}^{\mathbf{c}}|\mathbf{x}_{t}^{\mathbf{c}},\mathbf{x}_{0}^{\mathbf{c}}) = \operatorname{Bern}\left(\mathbf{x}_{0}^{\mathbf{c}} + \left[\left(1 - \mathbf{x}_{0}^{\mathbf{c}}\right) \odot \mathbf{x}_{t}^{\mathbf{c}}\right] \frac{n(t-1) - 1}{n(t) - 1}\right)$$
(5)

Here \odot *stands for the Hadamard products between two vectors.*

Detailed proof can be found in Appendix A.

3.2 Reverse Candidate Pruning

The reverse process implied by $p_{\theta}(\mathbf{x}_{t-1}^{\mathbf{c}}|\mathbf{x}_{t}^{\mathbf{c}})$ corresponds to the progressive candidate pruning process. We follow previous literature [Austin et al., 2021, Sahoo et al., 2024] to parameterize $p_{\theta}(\mathbf{x}_{t-1}^{\mathbf{c}}|\mathbf{x}_{t}^{\mathbf{c}})$, by combining a neural network(θ) predicted $\mathbf{x}_{0}^{\mathbf{c}}$ based on $\mathbf{x}_{t}^{\mathbf{c}}$ and the formulation of the posterior in Eq. 5:

$$p_{\theta}(\mathbf{x}_{t-1}^{\mathbf{c}}|\mathbf{x}_{t}^{\mathbf{c}}) = q(\mathbf{x}_{t-1}^{\mathbf{c}}|\mathbf{x}_{t}^{\mathbf{c}}, \text{NN}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}}, t))$$

$$= \text{Bern}\left(\left[\text{NN}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}}, t) + (1 - \text{NN}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}}, t)) \frac{n(t-1) - 1}{n(t) - 1}\right] \odot \mathbf{x}_{t}^{\mathbf{c}}\right)$$
(6)

Here $\mathrm{NN}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}},t)$ refers to a probability distribution over K-dim, e.g., outputs after softmax. Each parameter in $p_{\theta}(\mathbf{x}_{t-1}^{\mathbf{c}} \mid \mathbf{x}_{t}^{\mathbf{c}})$ can be viewed as an interpolation between the constant value $\frac{n(t-1)-1}{n(t)-1}$ and 1, weighted by $\mathrm{NN}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}},t)$.

Moreover, we propose incorporating the property of the forward process where $\mathbf{x}_{t-1}^{\mathbf{c}}$ is strictly contained within $\mathbf{x}_t^{\mathbf{c}}$, expressed as $[\mathbf{x}_t^{\mathbf{c}}]^T \mathbf{x}_{t-1}^{\mathbf{c}} = \mathbf{1}^T \mathbf{x}_{t-1}^{\mathbf{c}}$. This property is integrated into parameterization by ensuring that $\mathrm{NN}_{\theta}\left(\mathbf{x}_t^{\mathbf{c}},t\right)$ has non-zero values only for categories within $\mathbf{x}_t^{\mathbf{c}}$, satisfying $[\mathrm{NN}_{\theta}\left(\mathbf{x}_t^{\mathbf{c}},t\right)]^T\left(\mathbf{1}-\mathbf{x}_t^{\mathbf{c}}\right)=0$. Practically, such condition can be satisfied by ading $-\infty$ to the logits before the softmax operation. The prior distribution is set as the all-one vector, *i.e.*, $p_{\theta}(\mathbf{x}_T^{\mathbf{c}})=\mathrm{Bern}(\mathbf{1})$.

3.3 Training Procedure

We insert the formulation in Eq. 4, Eq. 5 and Eq. 6 into the Variational lower bound in Eq. 1 to derive the final objective for shortlisting models. We start with the first term L_T . As mentioned above, $p_{\theta}(\mathbf{x}_T^{\mathbf{c}}) = \text{Bern}(\mathbf{1})$. We put the time step T into Eq. 4,

$$q(\mathbf{x}_T^{\mathbf{c}}|\mathbf{x}_0^{\mathbf{c}}) = \operatorname{Bern}\left(\mathbf{x}_0^{\mathbf{c}} + (1 - \mathbf{x}_0^{\mathbf{c}})\frac{n(T) - 1}{K - 1}\right) = \operatorname{Bern}\left(\mathbf{x}_0^{\mathbf{c}} + (1 - \mathbf{x}_0^{\mathbf{c}})\frac{K - 1}{K - 1}\right) = \operatorname{Bern}(\mathbf{1})$$

The first term L_T in Eq. 1 becomes: $L_T = \mathbb{E}_{q(\mathbf{x}_0^{\mathbf{c}})} D_{\mathrm{KL}} \left[\mathrm{Bern}(\mathbf{1}) \| \, \mathrm{Bern}(\mathbf{1}) \right] = 0$. For the last term L_0 , with n(0) = 1 the $p_{\theta}(\mathbf{x}_0^{\mathbf{c}} | \mathbf{x}_1^{\mathbf{c}})$ in Eq. 6 can be expressed as $p_{\theta}(\mathbf{x}_0^{\mathbf{c}} | \mathbf{x}_1^{\mathbf{c}}) = \mathrm{Bern}(\mathrm{NN}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, t))$. Then L_0 is expressed as:

$$\begin{split} L_0 &= -\mathbb{E}_{q\left(\boldsymbol{x}_1^{\mathbf{c}} | \boldsymbol{x}_0^{\mathbf{c}}\right)} \left[\log p_{\theta}\left(\boldsymbol{x}_0^{\mathbf{c}} \mid \boldsymbol{x}_1^{\mathbf{c}}\right)\right] \\ &= -\mathbb{E}_{q\left(\boldsymbol{x}_1^{\mathbf{c}} | \boldsymbol{x}_0^{\mathbf{c}}\right)} \left[\log \left\langle \text{NN}_{\theta}(\mathbf{x}_1^{\mathbf{c}}, t), \boldsymbol{x}_0^{\mathbf{c}} \right\rangle + \begin{cases} \log \left\langle 1 - \text{NN}_{\theta}(\mathbf{x}_1^{\mathbf{c}}, t), \boldsymbol{x}_1^{\mathbf{c}} - \boldsymbol{x}_0^{\mathbf{c}} \right\rangle, & \|\mathbf{x}_1^{\mathbf{c}} - \mathbf{x}_0^{\mathbf{c}}\| > 0 \\ 0, & \|\mathbf{x}_1^{\mathbf{c}} - \mathbf{x}_0^{\mathbf{c}}\| = 0 \end{cases} \end{split}$$

Here $\langle \cdot \rangle$ denotes the inner product. Next, we focus on the term L_{t-1} , and for simplicity we use $\operatorname{pred}_{\theta}(\mathbf{x}_t^{\mathbf{c}})$ as a shorted notation for $[\operatorname{NN}_{\theta}(\mathbf{x}_t^{\mathbf{c}},t) + (1-\operatorname{NN}_{\theta}(\mathbf{x}_t^{\mathbf{c}},t)) \frac{n(t-1)-1}{n(t)-1}] \odot \mathbf{x}_t^{\mathbf{c}})$, correspondingly, $\operatorname{gd}(\mathbf{x}_t^{\mathbf{c}})$ for $\mathbf{x}_0^{\mathbf{c}} + [(1-\mathbf{x}_0^{\mathbf{c}}) \odot \mathbf{x}_t^{\mathbf{c}}] \frac{n(t-1)-1}{n(t)-1}$.

$$L_{t-1} = \mathbb{E}_{q(\mathbf{x}_{t}^{\mathbf{c}}|\mathbf{x}_{0}^{\mathbf{c}})} \left[D_{\mathrm{KL}} \left[\mathrm{Bern}(\mathrm{gd}(\mathbf{x}_{t}^{\mathbf{c}})) \| \operatorname{Bern}(\mathrm{pred}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}})) \right] \right]$$
(7)

The KL divergence between the Multivariate Bernoulli distribution is extended as:

$$D_{\text{KL}}[\text{Bern}(\text{gd}(\mathbf{x}_{t}^{\mathbf{c}})) \| \text{ Bern}(\text{pred}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}}))]$$

$$= \sum_{i,(\mathbf{x}_{t}^{\mathbf{c}})^{i} > 0} \left(\text{gd}^{i}(\mathbf{x}_{t}^{\mathbf{c}}) \log \frac{\text{gd}^{i}(\mathbf{x}_{t}^{\mathbf{c}})}{\text{pred}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}})} + (1 - \text{gd}^{i}(\mathbf{x}_{t}^{\mathbf{c}})) \log \frac{1 - \text{gd}^{i}(\mathbf{x}_{t}^{\mathbf{c}})}{1 - \text{pred}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}})} \right)$$
(8)

Here, we use the superscript i to denote the i-th dimension.

Mitigating Gradient Vanishing. We observe that directly optimizing the KL divergence of a multidimensional Bernoulli distribution, as in Eq. 7, can lead to optimization failure, with the process stalled from the beginning. This issue is likely due to gradient vanishing, where the gradients become too small to drive effective parameter updates. We formally investigate this issue in the following.

Taking dimension i in the K dimensions and $(\mathbf{x}_t^\mathbf{c})^i > 0$ as an example, the gradient towards the parameter θ is $\nabla_\theta D_{\mathrm{KL}}[\mathrm{Bern}(\mathrm{gd}^i(\mathbf{x}_t^\mathbf{c})\| \, \mathrm{Bern}(\mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c})] = -(\frac{\mathrm{gd}^i(\mathbf{x}_t^\mathbf{c})}{\mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c})} - \frac{1-\mathrm{gd}^i(\mathbf{x}_t^\mathbf{c})}{1-\mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c})}) \nabla_\theta \mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c}).$ And we denote the above gradient term as $\nabla_\theta D_{\mathrm{KL}}^i$ for simplicity in the following. Recall that $\mathrm{gd}^i(\mathbf{x}_t^\mathbf{c})$ and $\mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c})$ are both interpolations between 1 and $\frac{n(t-1)-1}{n(t)-1}$ as discussed in Section 3.2, we have: $\frac{n(t-1)-1}{n(t)-1} \leq \mathrm{gd}^i(\mathbf{x}_t^\mathbf{c}), \mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c}) \leq 1.$ Consider the common situation when $(\mathbf{x}_0^\mathbf{c})^i = 1$, and the network prediction $\mathrm{NN}_\theta^i(\mathbf{x}_t^\mathbf{c}, t)$ holds a non-zero value. The norm of weight satisfies that: $\left\|\frac{\mathrm{gd}^i(\mathbf{x}_t^\mathbf{c})}{\mathrm{pred}_\theta(\mathbf{x}_t^\mathbf{c})} - \frac{1-\mathrm{gd}^i(\mathbf{x}_t^\mathbf{c})}{1-\mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c})}\right\|_2 \leq \frac{n(t)-1}{n(t-1)-1}.$ Combining with $\nabla_\theta \mathrm{pred}_\theta^i(\mathbf{x}_t^\mathbf{c}) = \frac{n(t)-n(t-1)}{n(t)-1}\nabla_\theta \mathrm{NN}_\theta^i(\mathbf{x}_t^\mathbf{c}, t)$, we could obtain the following bounded condition over the gradient norm of $\nabla_\theta D_{\mathrm{KL}}^i$:

$$\begin{split} &\|\nabla_{\theta}D_{\mathrm{KL}}^{i}\|_{2} = \left\|\frac{\mathrm{gd}^{i}(\mathbf{x}_{t}^{\mathbf{c}})}{\mathrm{pred}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}})} - \frac{1 - \mathrm{gd}^{i}(\mathbf{x}_{t}^{\mathbf{c}})}{1 - \mathrm{pred}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}})}\right\|_{2} \left\|\nabla_{\theta}\mathrm{pred}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}})\right\|_{2} \\ &\leq \frac{n(t) - 1}{n(t - 1) - 1} \frac{n(t) - n(t - 1)}{n(t) - 1} \left\|\nabla_{\theta}\mathrm{NN}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t)\right\|_{2} = \frac{n(t) - n(t - 1)}{n(t - 1) - 1} \left\|\nabla_{\theta}\mathrm{NN}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t)\right\|_{2} \end{split} \tag{9}$$

Note that the bounds involve the gradient term $\nabla_{\theta} NN_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t)$, which is computed directly from the softmax outputs (without applying the logarithm). However, taking gradients through the softmax function directly often leads to vanishing gradients, particularly in high-dimensional settings where the outputs $NN_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t)$ can be very small initially. We additionally provide formal illustration in Appendix A.2.

To mitigate this issue, we propose scaling the gradient in Eq. 9 by $\frac{1}{NN_{\theta}^{i}(\mathbf{x}_{t}^{e},t)}$, consistent with the widely adopted log-softmax optimization. Surprisingly, this corresponds directly to the following simplified objective:

$$L_{t-1}^{\text{weight}} = -\mathbb{E}_{q\left(\boldsymbol{x}_{t}^{\mathbf{c}} | \boldsymbol{x}_{0}^{\mathbf{c}}\right)} \left[\frac{n(t) - n(t-1)}{n(t) - 1} \langle \log NN_{\theta}(\boldsymbol{x}_{t}^{\mathbf{c}}, t), \boldsymbol{x}_{0}^{\mathbf{c}} \rangle \right]$$
(10)

Though derived with heuristic intuition, we formally show in Appendix A.3 that this reweighted loss (Eq. 10) can be interpreted as a reasonable approximation.

Moreover, we can derive an even simpler training objective by removing the weight, analogous to the simplified loss used in [Ho et al., 2020], which may provide different practical properties:

$$L_{t-1}^{\text{simple}} = -\mathbb{E}_{q\left(\boldsymbol{x}_{t}^{\mathbf{c}} | \boldsymbol{x}_{0}^{\mathbf{c}}\right)} \left[\left\langle \log \text{NN}_{\theta}(\boldsymbol{x}_{t}^{\mathbf{c}}, t), \boldsymbol{x}_{0}^{\mathbf{c}} \right\rangle \right]$$
(11)

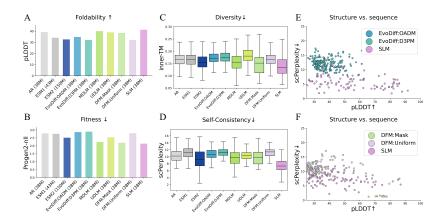


Figure 3: Quantitative Performance on Protein Sequence Design SLM compared to baselines: (A-D) pLDDT(A), Progen2-nll(B), scPerplexity(C) and inner-TM(D) scores for sequence sampled from ESM1-43M, ESM2-150M, and the following 38M models: EvoDiff-OADM, EvoDiff-D3PM, MDLM, UDLM, DFM-Mask, DFM-Uniform, SLM. (E-F) The joint distribution of pLDDT and scPerplexity from SLM model and Diffusion Models(E) and Discrete Flow Matching Models(F).

The above objective is essentially the Cross-Entropy loss between the network prediction and the original data sample. Unless otherwise specified, we use L^{simple} for experiments on OpenWebText and L^{weight} for all other experiments. The likelihood(ELBO) is strictly evaluated using the original ELBO defined in Eq. 7.

Candidate Set Size Scheduling. Another important component of the framework is the scheduling function over the size of candidate set, *i.e.* n(t). It is noteworthy that n(t) is not restricted to integer values; rather, it can take any real value within the interval [1, K]. We take a similar intuition from [Graves et al., 2023], by considering the normalized vector $\frac{\mathbf{x}_t^c}{\sum_{i=1}^K (\mathbf{x}_t^c)^i}$ as the probability of distribution over vocabulary, then we expect the entropy of the distribution to increase linearly from t=1 to t=T. Note the expected ones of \mathbf{x}_t^c is exactly n(t), hence the corresponding entropy of the aforementioned distribution is $\log n(t)$. Then we can design scheduling function as: $n(t) = e^{(\log K)\frac{t}{T}}$.

3.4 Sampling Process

The sampling process of shortlisting models can be directly conducted with ancestral sampling based on the learned $p_{\theta}(\mathbf{x}_{t-1}^{\mathbf{c}}|\mathbf{x}_{t}^{\mathbf{c}})$ with $\mathbf{x}_{T}^{\mathbf{c}} \sim \mathrm{Bern}(\mathbf{1})$ as the starting point. The full pseudocode for training and sampling is provided in Appendix B.2. To ensure the candidate set always contains at least one candidate (i.e., $\mathbf{x}_{t}^{\mathbf{c}} \neq \mathbf{0}$), we empirically set the dimension with the largest Bernoulli parameter to 1 when the sampled vector is a zero vector.

Classifier-free Guidance. We show that this simplified formulation offers flexibility to implement classifier-free guidance with an extra class-conditioned shortlist model. Denoting the output of the unconditional model at timestep t as $\mathrm{NN}_{\theta}(\mathbf{x}_t^\mathbf{c},t,K)$ and the conditional model as $\mathrm{NN}_{\theta}(\mathbf{x}_t^\mathbf{c},t,\operatorname{cls})$. Here $\mathrm{cls} \in [0,K-1] \cap \mathbb{Z}$ denotes the class label. The reverse process based on classifier-free guidance can be obtained as: $p_{\theta}^{\mathrm{CFG}}(\mathbf{x}_{t-1}^\mathbf{c}|\mathbf{x}_t^\mathbf{c}) = \mathrm{Bern}\left(\left[\mathrm{N}\hat{\mathrm{N}}_{\theta}(\mathbf{x}_t^\mathbf{c},t) + (1-\mathrm{N}\hat{\mathrm{N}}_{\theta}(\mathbf{x}_t^\mathbf{c},t))\frac{n(t-1)-1}{n(t)-1}\right]\odot\mathbf{x}_t^\mathbf{c}\right)$ Here the $\mathrm{N}\hat{\mathrm{N}}_{\theta}$ is as: $\mathrm{N}\hat{\mathrm{N}}_{\theta}(\mathbf{x}_t^\mathbf{c},t) = \gamma \mathrm{NN}_{\theta}(\mathbf{x}_t^\mathbf{c},t,\operatorname{cls}) + (1-\gamma)\mathrm{NN}_{\theta}(\mathbf{x}_t^\mathbf{c},t,K)$. When $\gamma>1$, there can be negative number in $\mathrm{N}\hat{\mathrm{N}}_{\theta}$. Following [Stark et al., 2024], we project the value of $\mathrm{N}\hat{\mathrm{N}}_{\theta}$ back to the simplex based on [Wang and Carreira-Perpinán, 2013].

4 Experiments

In this section, we evaluate our shortlisting models across various discrete data generation tasks and benchmarks. These tasks include language modeling and biological sequence design, the latter of which is especially well suited for non-autoregressive models and for demonstrating the potential of our proposed method.

4.1 Language Modeling

Text8: Firstly, we conduct experiments on the text8 dataset [Mahoney, 2011] with vocab size of 27. Bits-per-character(BPC) is reported based on the Equation. 8. Details on the dataset can be found in Appendix.C.1.1. The results are reported in Table.1, and additional generated samples are provided in Table.7.

We compare our shortlisting model with baselines across various categories: (1) autoregressive models: Transformer AR[Vaswani et al., 2017], AR Argmax Flow[Hoogeboom et al., 2021a], AR Discrete Flow[Tran et al., 2019]; (2) any-order autoregressive models: ARDM[Hoogeboom et al., 2021b], MAC [Shih et al., 2022]; (3) embedding-space continuous diffusion models: Plaid[Gulrajani and Hashimoto, 2024]); (4) advanced discrete diffusion models: SEDD[Lou et al., 2023], MDLM[Sahoo et al., 2024], UDLM[Schiff et al., 2024], D3PM variants[Austin et al., 2021]; and (5) simplex-based approaches: BFN[Graves et al., 2023], SFM[Cheng et al., 2024].

As aforementioned, we report the BPC of both the shortlisting model(SLM) trained with the Table 1: Bits Per Character (BPC) on Text8 Test L^{simple} in Eq. 11 and with the L^{weight} in Eq. 10. Set. As demonstrated in Table 1, even with the simplified objective, the proposed approach achieves competitive performance compared to other non-autoregressive approaches. Moreover, the reweighted formulation further boosts density estimation performance by better aligning with the original ELBO, as discussed in Section 3.3.

OpenWebText: We further investigate the challenges and potential of simplex-based approaches in large vocabulary settings, over the OpenWebText [Gokaslan and Cohen, 2019] dataset with vocab size of 50527. Detailed results and discussions are provided in Appendix C.1.2 and Table. 5, highlighting that while our SLM still lags behind advanced autoregressive models in density estimation, it achieves competitive generation performance and significantly outperforms existing simplexbased methods.

Category	Method	BPC (↓)
Autoregressive	Transformer AR AR Argmax Flow AR Discrete Flow	1.23 1.39 1.23
Any-order Autoregressive	ARDM MAC	≤ 1.43 ≤ 1.40
Continuous Diffusion	Plaid	≤ 1.48
Discrete Diffusion	Mult. Diffusion D3PM Uniform UDLM D3PM Absorb SEDD Absorb MDLM MD4	≤ 1.72 ≤ 1.61 ≤ 1.59 ≤ 1.45 ≤ 1.41 ≤ 1.39 ≤ 1.37
Simplex Approaches	$egin{array}{c} ext{BFN} \ ext{DFM} \ ext{SFM} \ ext{SLM}(L^{ ext{simple}}) \ ext{SLM}(L^{ ext{weight}}) \end{array}$	

4.2 De novo Design of Protein Sequence

In this experiment, we focus on the core task of unconditional protein design, and examine various protein properties to demonstrate SLM's superiority in this task.

Baselines: We compare SLM against three groups of existing methods: (1) Autoregressive models (AR); (1) Masked language models (MLMs), specifically **ESM1**[Rives et al., 2019] and **ESM2**[Lin et al., 2022]; (2) Discrete Diffusion Models, represented by two versions of EvoDiff[Alamdari et al., 2023]: EvoDiff-OADM, EvoDiff-D3PM, MDLM[Sahoo et al., 2024] and UDLM[Schiff et al., 2024]; (3) Discrete Flow Matching Models[Gat et al., 2024]: DFM-Mask and DFM-Uniform. Further information of these baselines in Appendix C.5.

To demonstrate SLM's effectiveness in protein sequence generation, we evaluate four key properties: (1) Foldability: structural plausibility predicted by ESMFold [Lin et al., 2022]; (2) Fitness: scores predicted by ProGen2-xlarge [Nijkamp et al., 2023]; (3) Self-Consistency: alignment between sequences folded by ESMFold and inverse-folded by ESM-IF [Hsu et al., 2022]; and (4) **Diversity**: pairwise inner-TM scores among generated samples. Detailed metric definitions are provided in Appendix C.5.1. As shown in Figure 3, SLM surpasses all baselines across all metrics and

Table 2: Conditional generation over promoter Table 3: FBD metric for sequence generation under design. BFN results are from our experiments and the other baselines from [Davis et al., 2024].

two datasets. CFG refers to Classifier-Free Guidance.

Model	MSE(↓)
Autoregressive	0.034 ± 0.001
Bit Diffusion (bit enc)	0.041
Bit Diffusion (one-hot)	0.040
D3PM-uniform	0.038
UDLM	0.030 ± 0.001
MDLM	0.028 ± 0.001
Dirichlet FM	0.034 ± 0.004
Fisher-Flow	0.029 ± 0.001
BFN	0.0405 ± 0.0003
DDSM	0.033
SLM	$0.0265 {\pm} 0.0006$

$Mel\ FBD(\downarrow)$	FB FBD(↓)
619.0 ± 0.8	832.4 ± 0.3
35.4 ± 0.5	25.7 ± 1.0
27.5 ± 2.6	3.8 ± 0.3
5.3 ± 0.5	15.1 ± 0.4
3.3 ± 0.1	10.8 ± 0.6
$\textbf{2.2} {\pm} \textbf{0.1}$	$4.4 {\pm} 0.2$
2.3 ± 0.1	2.3 ± 0.2
G 1.9 ± 0.4	1.0 ± 0.3
1.4 ± 0.1	$\boldsymbol{1.0 {\pm} 0.1}$
	$\begin{array}{c} & & & & \\ \hline 619.0 \pm 0.8 \\ & & & \\ 35.4 \pm 0.5 \\ & & \\ 27.5 \pm 2.6 \\ & & \\ 5.3 \pm 0.5 \\ & & \\ 3.3 \pm 0.1 \\ \hline & & \\ \textbf{2.2} \pm \textbf{0.1} \\ \hline & & \\ 2.3 \pm 0.1 \\ \hline & & \\ G & 1.9 \pm 0.4 \\ \hline \end{array}$

achieves competitive performance compared to ESM2-150M [Lin et al., 2022], demonstrating strong generalization and robustness under restricted vocabularies and complex data distributions.

4.3 Design of DNA Sequence

In this part, we focus on the roles of promoters and enhancers, and evaluate SLM in this context. Following prior work, we set the language model to use 500 NFE for enhancers and 1000 NFE for promoters. For other models, Non-CFG models use 100 NFE, and CFG variants use 200 NFE.

4.3.1 Promoter DNA Sequence Design

We follow the setting in previous work DDSM [Avdeyev et al., 2023] to generate DNA promoter sequences conditioned on the promoter profile.

Data: We use a dataset of 100,000 human promoter sequences, each 1024 base-pairs long, extracted from the Human Promoter Database [Hon et al., 2017]. Each sequence is paired with a CAGE signal indicating transcriptional activity at each position [Shiraki et al., 2003, Forrest et al., 2014]. Sequences from chromosomes 8 and 9 form the test set, and the remainder are used for training.

Baselines: We compared the SLM method with: (1) flow matching methods including Dirichlet FM [Stark et al., 2024] and Fisher-Flow [Davis et al., 2024]; (2) autoregressive language model that generates base

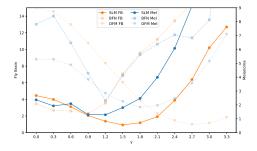


Figure 4: Performance of SLM under different CFG factor γ for unconditional enhancer design.

pairs; (3) Bayesian Flow Networks (BFN) [Graves et al., 2023]; and (4) other discrete diffusion methods including two implementations of Bit Diffusion [Chen et al., 2022a], D3PM [Austin et al., 2021], MDLM [Sahoo et al., 2024], UDLM [Schiff et al., 2024] and simplex-based DDSM [Avdeyev et al., 2023].

Results: The regulatory activity of the sequences is given by Sei, a model that predicts the regulatory potential of the sequences [Chen et al., 2022b]. We report the mean and standard deviation of MSE between the generated sequences and the target. Our MSE values were measured under the same Sei model as in previous works. As shown in Table.2, our SLM method achieves the lowest MSE, with a smaller standard deviation as well.

4.3.2 Enhancer DNA Sequence Design

We also evaluate SLM for enhancer generation, following DirichletFM [Stark et al., 2024].

Table 4: Ablation study on loss function conducted on DNA sequence design.

Loss type	M	el	F	В	prom	oter
Zess type	FBD(↓)	PPL(↓)	FBD(↓)	PPL(↓)	$\overline{\text{MSE}(\downarrow)}$	$PPL(\downarrow)$
L^{simple} L^{weight}	2.1788 2.5848	3.4102 3.4018	4.4450 4.9670	3.4618 3.4654	0.0265 0.0260	2.8084 2.7672

Data: We use 104k enhancer sequences from fly brain cells and 89k from human melanoma cells [Janssens et al., 2022, Atak et al., 2021], each with a length of 500. Cell type labels were determined by ATAC-seq data[Buenrostro et al., 2013], with fly brain cells divided into 81 classes and human melanoma cells into 47 classes based on cell types.

Baselines: In addition to their standard implementations, the baseline models also incorporate classifier-free guidance. We select the optimal classifier-free guidance factor γ for all models. The performance of our method under different classifier-free guidance factors γ is presented in Figure.4. Specific experimental settings and details can be found in Appendix C.2.2.

Results: We evaluate generated sequences using the Fréchet Biological Distance (FBD) from Dirichlet FM [Stark et al., 2024], which treats classifier-hidden representations as sequence embeddings and computes FBD as the Wasserstein distance between them. SLM achieves the best performance without label guidance and further improves with label guidance (see Table 3).

4.4 Ablation Study on Reweighted Loss

We conduct ablation experiments comparing the simple loss and the reweighted loss on both enhancer and promoter tasks. Table 4 shows that the reweighted loss achieves superior density estimation, while the simple loss can generate better samples.

5 Comparison to Existing Works

Our SLM connects closely to several existing approaches. When K=2, SLM resembles Bernoulli diffusion [Sohl-Dickstein et al., 2015], though it operates within a three-state space [0,1], [1,0], [1,1] rather than Bernoulli diffusion's two-state (0,1) space. The progressive candidate exclusion in SLM also shares conceptual similarities with the SUBS parameterization employed in MDLM [Shi et al., 2024, Sahoo et al., 2024]. Unlike mask-based discrete diffusion methods, however, SLM directly operates on the simplex, refining information gradually over time.

Recent simplex-based methods such as DirichletFM [Stark et al., 2024] define trajectories over the continuous full simplex, whereas our SLM explicitly restricts transitions to discrete centroids or their subspaces, thereby reducing degrees of freedom and enhancing efficiency and interpretability. Additionally, while DirichletFM optimizes an MSE-based loss, SLM adopts a simplified cross-entropy objective (Eq. 10), alleviating gradient vanishing issues associated with Bernoulli KL losses and leading to improved performance in challenging tasks.

6 Conclusion

In this paper, we introduce the Shortlisting Model (SLM), a novel discrete generative model inspired by progressive candidate pruning. SLM follows a unique generation trajectory by transitioning from the centroids of the simplex space. With impressive performance across various tasks, SLM offers a simple yet effective alternative for discrete generative modeling.

Limitation This work focuses on simplex-based discrete generative models, with evaluations conducted on text and biological sequence data. Comprehensive studies on other modalities, such as graphs and images, are beyond the scope of this paper and are left for future work. Additionally, further theoretical analysis and engineering improvements are required to scale our approach to large, real-world applications.

Acknowledgments

This work is supported by the Natural Science Foundation of China (Grant No. 62376133) and sponsored by Beijing Nova Program (20240484682) and the Wuxi Research Institute of Applied Technologies, Tsinghua University (20242001120).

References

- 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.
- 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.
- Pavel Avdeyev, Chenlai Shi, Yuhao Tan, Kseniia Dudnyk, and Jian Zhou. Dirichlet diffusion score model for biological sequence generation. In *International Conference on Machine Learning*, pages 1276–1301. PMLR, 2023.
- Hannes Stark, Bowen Jing, Chenyu Wang, Gabriele Corso, Bonnie Berger, Regina Barzilay, and Tommi Jaakkola. Dirichlet flow matching with applications to dna sequence design. *arXiv* preprint *arXiv*:2402.05841, 2024.
- Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Diffusion language models are versatile protein learners. *arXiv* preprint arXiv:2402.18567, 2024.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- Clement Vignac, Igor Krawczuk, Antoine Siraudin, Bohan Wang, Volkan Cevher, and Pascal Frossard. Digress: Discrete denoising diffusion for graph generation. *arXiv preprint arXiv:2209.14734*, 2022.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Forty-first International Conference on Machine Learning*.
- Subham Sekhar Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin T Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. arXiv preprint arXiv:2406.07524, 2024.
- Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis K Titsias. Simplified and generalized masked diffusion for discrete data. *arXiv* preprint arXiv:2406.04329, 2024.
- Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. *arXiv preprint arXiv:2407.15595*, 2024.
- Oscar Davis, Samuel Kessler, Mircea Petrache, Ismail Ilkan Ceylan, Michael Bronstein, and Avishek Joey Bose. Fisher flow matching for generative modeling over discrete data. *arXiv* preprint arXiv:2405.14664, 2024.
- Chaoran Cheng, Jiahan Li, Jian Peng, and Ge Liu. Categorical flow matching on statistical manifolds. *arXiv* preprint arXiv:2405.16441, 2024.
- Minkai Xu, Tomas Geffner, Karsten Kreis, Weili Nie, Yilun Xu, Jure Leskovec, Stefano Ermon, and Arash Vahdat. Energy-based diffusion language models for text generation. *arXiv preprint arXiv:2410.21357*, 2024.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *arXiv preprint* arXiv:2006.11239, 2020.
- Alex Graves, Rupesh Kumar Srivastava, Timothy Atkinson, and Faustino Gomez. Bayesian flow networks. *arXiv preprint arXiv:2308.07037*, 2023.

- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022.
- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. Advances in Neural Information Processing Systems, 34:17981–17993, 2021.
- Weiran Wang and Miguel A Carreira-Perpinán. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application. *arXiv* preprint arXiv:1309.1541, 2013.
- Matt Mahoney. Large text compression benchmark. 2011.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- 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, 2021a.
- Dustin Tran, Keyon Vafa, Kumar Krishna Agrawal, Laurent Dinh, and Ben Poole. Discrete flows: Invertible generative models of discrete data. *arXiv preprint arXiv:1905.10347*, 2019.
- Emiel Hoogeboom, Alexey A Gritsenko, Jasmijn Bastings, Ben Poole, Rianne van den Berg, and Tim Salimans. Autoregressive diffusion models. *arXiv preprint arXiv:2110.02037*, 2021b.
- Andy Shih, Dorsa Sadigh, and Stefano Ermon. Training and inference on any-order autoregressive models the right way. *Advances in Neural Information Processing Systems*, 35:2762–2775, 2022.
- Ishaan Gulrajani and Tatsunori B Hashimoto. Likelihood-based diffusion language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion language modeling by estimating the ratios of the data distribution. 2023.
- Yair Schiff, Subham Sekhar Sahoo, Hao Phung, Guanghan Wang, Sam Boshar, Hugo Dalla-torre, Bernardo P de Almeida, Alexander Rush, Thomas Pierrot, and Volodymyr Kuleshov. Simple guidance mechanisms for discrete diffusion models. *arXiv* preprint arXiv:2412.10193, 2024.
- Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. 2019.
- Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C. Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS*, 2019. doi: 10.1101/622803. URL https://www.biorxiv.org/content/10.1101/622803v4.
- Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Alex X Lu, Nicolo Fusi, Ava P Amini, and Kevin K Yang. Protein generation with evolutionary diffusion: sequence is all you need. *BioRxiv*, pages 2023–09, 2023.
- Erik Nijkamp, Jeffrey A Ruffolo, Eli N Weinstein, Nikhil Naik, and Ali Madani. Progen2: exploring the boundaries of protein language models. *Cell systems*, 14(11):968–978, 2023.
- Chloe Hsu, Robert Verkuil, Jason Liu, Zeming Lin, Brian Hie, Tom Sercu, Adam Lerer, and Alexander Rives. Learning inverse folding from millions of predicted structures. *ICML*, 2022. doi: 10.1101/2022.04.10.487779. URL https://www.biorxiv.org/content/early/2022/04/10/2022.04.10.487779.
- Chung-Chau Hon, Jordan A Ramilowski, Jayson Harshbarger, Nicolas Bertin, Owen JL Rackham, Julian Gough, Elena Denisenko, Sebastian Schmeier, Thomas M Poulsen, Jessica Severin, et al. An atlas of human long non-coding rnas with accurate 5: ends. *Nature*, 543(7644):199–204, 2017.

- Toshiyuki Shiraki, Shinji Kondo, Shintaro Katayama, Kazunori Waki, Takeya Kasukawa, Hideya Kawaji, Rimantas Kodzius, Akira Watahiki, Mari Nakamura, Takahiro Arakawa, et al. Cap analysis gene expression for high-throughput analysis of transcriptional starting point and identification of promoter usage. *Proceedings of the National Academy of Sciences*, 100(26):15776–15781, 2003.
- Alistair RR Forrest, Hideya Kawaji, Michael Rehli, J Kenneth Baillie, Michiel JL De Hoon, Vanja Haberle, Timo Lassmann, Ivan V Kulakovskiy, Marina Lizio, Masayoshi Itoh, et al. A promoterlevel mammalian expression atlas. *Nature*, 507(7493):462–470, 2014.
- Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. arXiv preprint arXiv:2208.04202, 2022a.
- Kathleen M Chen, Aaron K Wong, Olga G Troyanskaya, and Jian Zhou. A sequence-based global map of regulatory activity for deciphering human genetics. *Nature genetics*, 54(7):940–949, 2022b.
- Jasper Janssens, Sara Aibar, Ibrahim Ihsan Taskiran, Joy N Ismail, Alicia Estacio Gomez, Gabriel Aughey, Katina I Spanier, Florian V De Rop, Carmen Bravo Gonzalez-Blas, Marc Dionne, et al. Decoding gene regulation in the fly brain. *Nature*, 601(7894):630–636, 2022.
- Zeynep Kalender Atak, Ibrahim Ihsan Taskiran, Jonas Demeulemeester, Christopher Flerin, David Mauduit, Liesbeth Minnoye, Gert Hulselmans, Valerie Christiaens, Ghanem-Elias Ghanem, Jasper Wouters, et al. Interpretation of allele-specific chromatin accessibility using cell state—aware deep learning. *Genome research*, 31(6):1082–1096, 2021.
- Jason D Buenrostro, Paul G Giresi, Lisa C Zaba, Howard Y Chang, and William J Greenleaf. Transposition of native chromatin for fast and sensitive epigenomic profiling of open chromatin, dna-binding proteins and nucleosome position. *Nature methods*, 10(12):1213–1218, 2013.
- Jascha Sohl-Dickstein, Eric A Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *arXiv preprint arXiv:1503.03585*, 2015.
- 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.
- Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical sampling. *arXiv* preprint arXiv:2409.02908, 2024.
- Baris E Suzek, Hongzhan Huang, Peter McGarvey, Raja Mazumder, and Cathy H Wu. Uniref: comprehensive and non-redundant uniprot reference clusters. *Bioinformatics*, 23(10):1282–1288, 2007.
- I Loshchilov. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017.
- Michael Heinzinger, Konstantin Weissenow, Joaquin Gomez Sanchez, Adrian Henkel, Martin Steinegger, and Burkhard Rost. Prostt5: Bilingual language model for protein sequence and structure. biorxiv. 2023.

A Mathematical Derivation

A.1 Proof of Proposition 3.1

Since x is a vector and the elements of the vector are independent, we only consider the position of a fixed index in all the vectors. In the following, all instances of x are redefined as scalars. First, we prove the following proposition:

$$p(\mathbf{x}_t^c = 1 \mid \mathbf{x}_0^c = 0) = \frac{n(t) - 1}{K - 1}$$
 (12)

When t = 0, $p(\mathbf{x}_0 = 1 \mid \mathbf{x}_0 = 0) = 0$ is obvious. Thus, we can proceed with induction on t.

$$q(\mathbf{x}_{t}^{\mathbf{c}} = 1 \mid \mathbf{x}_{0}^{\mathbf{c}} = 0)$$

$$= \frac{n(t-1)-1}{K-1} + \frac{n(t)-n(t-1)}{K-n(t-1)} \left(1 - \frac{n(t-1)-1}{K-1}\right)$$

$$= \frac{n(t)-1}{K-1}$$
(13)

Since $q(\mathbf{x}_t = 1 \mid \mathbf{x}_0 = 1) = 1$, the two cases can be combined into $q(\mathbf{x}_t \mid \mathbf{x}_0) = \operatorname{Bern}(\mathbf{x}_0 + (1 - \mathbf{x}_0) \frac{n(t) - 1}{K - 1})$, whose vector form is given by Eq. 4.

The only non-trivial case in the posterior distribution is:

$$q(\mathbf{x}_{t-1}^{\mathbf{c}} = 1 \mid \mathbf{x}_{t}^{\mathbf{c}} = 1, \mathbf{x}_{0}^{\mathbf{c}} = 0)$$

$$= \frac{q(\mathbf{x}_{t-1}^{\mathbf{c}} = 1, \mathbf{x}_{t}^{\mathbf{c}} = 1 \mid \mathbf{x}_{0}^{\mathbf{c}} = 0)}{q(\mathbf{x}_{t}^{\mathbf{c}} = 1 \mid \mathbf{x}_{0}^{\mathbf{c}} = 0)}$$

$$= \frac{q(\mathbf{x}_{t-1}^{\mathbf{c}} = 1 \mid \mathbf{x}_{0}^{\mathbf{c}} = 0)}{q(\mathbf{x}_{t}^{\mathbf{c}} = 1 \mid \mathbf{x}_{0}^{\mathbf{c}} = 0)}$$

$$= \frac{n(t-1)-1}{n(t)-1}$$
(14)

Only when $\mathbf{x_0^c}=1$, $q(\mathbf{x_{t-1}^c}=1\mid\mathbf{x_t^c}=1,\mathbf{x_0^c}=1)=1$. In all other cases, the probability is 0. Therefore, the result of Eq. 5 can be given.

A.2 Gradient Vanishing of $\nabla_{\theta} NN_{\theta}^{i}(x_{t}^{c}, t)$

Note that $NN_{i}^{i}(x_{i}^{c},t)$ is essentially the output of the softmax, which could be further expressed as:

$$NN_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t) = \frac{\exp\left(f_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t)\right)}{\sum_{j} \exp\left(f_{\theta}^{j}(\mathbf{x}_{t}^{\mathbf{c}}, t)\right)}$$
(15)

where f_{θ} denotes the raw output of the neural network. The gradient with respect to the softmax input is:

$$\frac{\partial \text{NN}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t)}{\partial f_{\theta}^{k}(\mathbf{x}_{t}^{\mathbf{c}}, t)} = \text{NN}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}}, t) \left(\delta_{ik} - \text{NN}_{\theta}^{k}(\mathbf{x}_{t}^{\mathbf{c}}, t)\right)$$
(16)

The norm is as $\|\mathbf{N}\mathbf{N}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}},t)\left(1-\mathbf{N}\mathbf{N}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}},t)\right)\|_{2}$ when k=i; and $\|\mathbf{N}\mathbf{N}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}},t)\mathbf{N}\mathbf{N}_{\theta}^{k}(\mathbf{x}_{t}^{\mathbf{c}},t)\|_{2}$ when $k\neq i$. Both case the norm is strictly bounded with the $\|\mathbf{N}\mathbf{N}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}},t)\|_{2}$. At the initial training stage, $\mathbf{N}\mathbf{N}_{\theta}^{i}(\mathbf{x}_{t}^{\mathbf{c}},t)$ may become uniformly small in high-dimensional settings, leading to vanishing gradients and causing the optimization to stall.

A.3 Clarification for the relationship between Eq. 10 and ELBO

We provide a formal derivation of the reweighted loss, which originates from an analysis of the gradient of the KL divergence. Let $\mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t)$ denote the model softmax output at i-th dimention. The corresponding predicted Bernoulli distribution's parameter at i-th dimention can be expressed as: $\mathrm{NN}_{\theta}^i(\mathbf{x}_t^\mathbf{c},t) + (1-\mathrm{NN}_{\theta}^i(\mathbf{x}_t^\mathbf{c},t)) \frac{n(t-1)-1}{n(t)-1}$. For notion simplicity, we define $q = \frac{n(t-1)-1}{n(t)-1}$. Taking the gradient of the KL divergence with respect to θ , we obtain the following expressions:

• For the case where $x_0^i = 0$:

$$\nabla_{\theta} D_{\text{KL}}[\text{Bern}(\text{gd}(\mathbf{x}_{t}^{\mathbf{c}})(i)) \| \text{ Bern}(\text{pred}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}})(i))]$$

$$= -\left[\frac{q}{\text{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}}, t) + (1 - \text{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}}, t))q} - \frac{1}{1 - \text{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}}, t)}\right] \cdot (1 - q) \nabla_{\theta} \text{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}}, t)$$
(17)

• For the case where $x_0^i = 1$:

$$\nabla_{\theta} D_{\text{KL}}[\text{Bern}(\text{gd}(\mathbf{x}_{t}^{\mathbf{c}})(i)) \| \text{ Bern}(\text{pred}_{\theta}(\mathbf{x}_{t}^{\mathbf{c}})(i))]$$

$$= -\frac{1}{\text{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}}, t) + (1 - \text{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}}, t))q} \cdot (1 - q) \nabla_{\theta} \text{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}}, t)$$
(18)

We consider the initial stage of training with high-dimensional data, hence the init value of $NN_{\theta}^{i}(x_{t}^{\mathbf{c}},t)$ is relatively small. For the cases when $x_{0}^{i}=0$, the $NN_{\theta}^{i}(x_{t}^{\mathbf{c}},t)$ already approach the optimal value and also the weight of gradient is approximately zero:

$$-\left[\frac{q}{\mathbf{N}\mathbf{N}_{\theta}^{i}(x_{t}^{\mathbf{c}},t) + (1 - \mathbf{N}\mathbf{N}_{\theta}^{i}(x_{t}^{\mathbf{c}},t))q} - \frac{1}{1 - \mathbf{N}\mathbf{N}_{\theta}^{i}(x_{t}^{\mathbf{c}},t)}\right] \approx -\left[\frac{q}{q} - 1\right] = 0 \tag{19}$$

Therefore, we could mainly focus on the case when $x_0^i=1$. As we mentioned in the above Appendix. A.2, the term $\nabla_{\theta} \mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t)$ could have the vanishing issues due to the property of softmax. However, the scale weight $\frac{1}{\mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t)+(1-\mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t))q}$ is bounded, i.e., $\frac{1}{\mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t)+(1-\mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t))q} \leq 1$ as $0 \leq \mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t), q \leq 1$, and hence could not help enhance the signal. Motivated by the widely optimized logsoftmax or logsumexp where the gradient scale weight is as $\frac{1}{\mathrm{NN}_{\theta}^i(x_t^\mathbf{c},t)}$, we fix the gradient as:

$$-\frac{1}{\mathrm{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}},t) + (1 - \mathrm{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}},t))q} \cdot (1 - q)\nabla_{\theta}\mathrm{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}},t) \to -\frac{1}{\mathrm{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}},t)} \cdot (1 - q)\nabla_{\theta}\mathrm{NN}_{\theta}^{i}(x_{t}^{\mathbf{c}},t)$$

$$= \nabla_{\theta}L^{\mathrm{weight}} \tag{20}$$

Note that the optimization challenge typically arises during the initial training stages, where our proposed objective can provide effective support. Direct optimization of the original ELBO in later training stages or epochs may be possible and could further improve density estimation performance. We leave exploring this direction for future work.

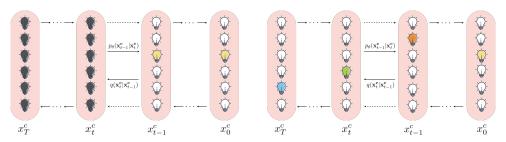


Figure 5: forward and reverse process of MDM(Left) and D3PM-Uniform(Right)

B Algorithms

B.1 Visualization of the forward and reverse process of MDLM and D3PM-Uniform

In this section, the forward and reverse process of MDLM and D3PM-Uniform are visualized in Figure. 5.

B.2 Training and Sampling Algorithms

In this section, we provide detailed information about the training and sampling processes of SLM , with pseudo code as shown in Algorithm.1, Algorithm.2 and Algorithm.3, with code implementations in PyTorch, as shown in Listing.1 and 2.

Algorithm 1: Forward Process $q(x_t^\mathbf{c} \mid x_0^\mathbf{c})$ Input: one-hot data $x_0^\mathbf{c}$, time t $n(t) \leftarrow e^{(\log K) \frac{t}{T}}$ Bern_param = $\frac{n(t)-1}{K-1}$ for i=0 to K-1 do if $x_0^\mathbf{c}[i] == 1$ then $x_t^\mathbf{c}[i] \leftarrow 1$ else $x_t^\mathbf{c}[i] \leftarrow \text{sample from Bern_param}$ end if end for Return $x_t^\mathbf{c}$

Algorithm 2: Training

```
Input: one-hot data x_0^{\mathbf{c}}, class label \mathrm{cls} \in [0,K-1] \cap \mathbb{Z} Sample t \sim U(0,1) n(t) \leftarrow e^{(\log K) \frac{t}{T}}, n(t-1) \leftarrow e^{(\log K) \frac{t-1}{T}} x_t^{\mathbf{c}} \leftarrow q(x_t^{\mathbf{c}} \mid x_0^{\mathbf{c}}) flag \sim U(0,1) if flag > 0.3 then \mathrm{cls\_inp} \leftarrow \mathrm{cls} else \mathrm{cls\_inp} \leftarrow K end if L \leftarrow \log(\langle \mathrm{NN}_{\theta}(x_t, \mathrm{cls\_inp}, t), x_0^{\mathbf{c}} \rangle) Return L
```

C Experimental Details

C.1 Language Modeling

C.1.1 Additional Information for Text8 Experiments

The text8 dataset [Mahoney, 2011] is a medium-sized character-level corpus with a vocabulary size of 27. It includes 26 lowercase letters and a space token, sourced from the m English Wikipedia dataset by removing punctuation and converting all text to lowercase. The data processing is directly following the previous works [Austin et al., 2021, Cheng et al., 2024, Graves et al., 2023] where the sequence is randomly chunked to have the length of 256 for both training and evaluation. We adapt DiT [Peebles and Xie, 2023] as the network backbone for shortlisting model. And to make a fair comparison the configuration is aligned with previous literatures [Lou et al., 2023]. We calculate the bits-per-character(BPC) based on the Equation. 8. For autoregressive methods, we set NFEs as 256, while diffusion-based and simplex-based methods use 1000 NFEs.

Algorithm 3: Sampling of Shortlisting Model

```
Input: class label cls \in [0, K-1] \cap \mathbb{Z}, classifier-free guidance (CFG) factor \gamma \in \mathbb{R} x_t^{\mathbf{c}} \leftarrow \mathbf{1} for t = T to 1 do n(t) \leftarrow e^{(\log K) \frac{t}{T}}, n(t-1) \leftarrow e^{(\log K) \frac{t-1}{T}} if CFG then N\hat{\mathbf{N}}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, t) \leftarrow \gamma \cdot N\mathbf{N}_{\theta}(x_t^{\mathbf{c}}, t, \operatorname{cls}) + (1-\gamma) \cdot N\mathbf{N}_{\theta}(x_t^{\mathbf{c}}, t, K) else N\hat{\mathbf{N}}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, t) \leftarrow N\mathbf{N}_{\theta}(x_t^{\mathbf{c}}, t) end if \operatorname{pred}_{\theta} \leftarrow N\hat{\mathbf{N}}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, t) + \frac{n(t-1)-1}{n(t)-1}(1-N\hat{\mathbf{N}}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, t)) x_t^{\mathbf{c}} \leftarrow \operatorname{sample} \operatorname{from} \operatorname{pred}_{\theta} end for if CFG then N\hat{\mathbf{N}}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, 0) \leftarrow \gamma \cdot N\mathbf{N}_{\theta}(x_t^{\mathbf{c}}, 0, \operatorname{cls}) + (1-\gamma) \cdot N\mathbf{N}_{\theta}(x_t^{\mathbf{c}}, 0, K) else N\hat{\mathbf{N}}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, 0) \leftarrow N\mathbf{N}_{\theta}(x_t^{\mathbf{c}}, 0) end if \operatorname{Return} \operatorname{arg} \max N\hat{\mathbf{N}}_{\theta}(\mathbf{x}_t^{\mathbf{c}}, 0)
```

C.1.2 Details of OpenwebText Experiments

We further explore the challenges and potential of simplex-based approaches in large vocabulary settings. Building on recent studies [Sahoo et al., 2024, Lou et al., 2023], we also evaluate shortlisting models using OpenWebText [Gokaslan and Cohen, 2019], an open-source replica of the unreleased WebText dataset used to train GPT-2. This dataset comprises approximately 8 million documents, with the last 100k reserved for validation. We tokenize the data using the GPT-2 tokenizer, which has a vocabulary size of 50,257. Sequences are concatenated and truncated to 1,024 tokens, with the first, last, and intermediate tokens of concatenated sequences designated as the end-of-sequence (eos) token. We set NFEs to 1024 for autoregressive methods and 1000 for diffusion-based and simplex-based methods.

Networks Architectures: For network architecture, we use 3 different size of transformers: 1) small model with 110M: Transformer with 12 layers, a hidden dimension of 768, 12 attention heads, and a timestep embedding of 128; 2) medium model with 460M: Transformer with with 24 layers, a hidden dimension of 1024, 16 attention heads, and a timestep embedding of 128; 3) large model with 1.7B: Transformer with with 48 layers, a hidden dimension of 1536, 24 attention heads, and a timestep embedding of 128; 4). The SLM_W^S for small model is Transformer with 8 layers, a hidden dimension of 1024, 12 attention heads, and a timestep embedding of 128. 5) The SLM_W^M for medium model is Transformer with 12 layers, a hidden dimension of 1596, 12 attention heads, and a timestep embedding of 128.

Metrics: We focus on both the likelihood-related metric and sample-based metrics. Specifically, we evaluate the Perplexity(**PPL**) over the validation set, which is defined as PPL = $\exp\left(\frac{\mathbb{E}_{x_0 \sim p_{\text{data}}}\left[-\log p_{\theta}(x_0)\right]}{D}\right)$. D is the data dimension and for model without exact formulation of likelihood, we report the variational bounds of $\log p_{\theta}$. For sample-based metrics, we select Generative Perplexity(**Gen-PPL**) [Lou et al., 2023] where generated samples are evaluated under GPT-2 large; Based on recent works [Zheng et al., 2024], we further involve the **Entropy** to measure the diversity of tokens in a sentence which is computed as $-\sum_{k=1}^{K} p_k \log p_k$. For a sequence of length L containing K distinct tokens, each token k appears L_k times. The probability of occurrence for token k is given by $p_k = \frac{L_k}{L}$. For sample-based metrics, we fix numerical issues of the categorical/Bernoulli sampling by adjusting its accuracy to 64-bit [Zheng et al., 2024] and diffusion-based approaches use 1024 steps for generation. We provided generated samples at Appendix.C.1.3.

Results: Table 5 shows that while our shortlisting model lags behind autoregressive and discrete diffusion models in likelihood-based metrics, it excels in sample-based metrics by balancing quality and diversity. Notably, compared to BFN [Graves et al., 2023], another advanced simplex-based

Table 5: The Performance over OpenwebText

Model	PPL(↓)	Gen-PPL(↓)	Entropy(†)
AR(110M)	21.04	37.62	5.617
SEDD(110M)	23.87	98.41	5.586
MDLM(110M)	23.08	101.24	5.609
BFN(110M)	105.66	299.95	4.981
SLM(110M)	53.90	65.59	5.494
SLM_W^S	43.25	53.79	5.618
SLM(460M)	39.01	55.07	5.508
SLM_W^M	37.32	39.39	5.587
SLM(1.7B)	36.75	43.52	5.550

approach, our model achieves significant improvements. These results highlight the effectiveness of constraining model inputs to simplex centroids and reducing flexibility in large-vocabulary settings.

Why do simplex-based approaches fail with large vocabularies? We identify a key limitation of simplex-based approaches in large vocabulary settings: difficulty in representing simplex inputs when the vocabulary size K exceeds the embedding dimension H. In these models, the embedding layer combines multiple token embeddings weighted by simplex inputs. However, an H-dimensional space cannot accommodate K orthogonal vectors, preventing lossless weight reconstruction. To address this, we conducted experiments by approximately maintaining the total number of parameters, reducing network depth, and increasing width, resulting in variants denoted as SLM_W^S and SLM_W^M . As shown in Table 5, these modifications significantly enhance performance, supporting our hypothesis and suggesting a promising direction for improving simplex-based models.

C.1.3 Samples for Text Generation

Several generated samples by SLM and one of the baselines: BFN are provided on the dataset of text8 and OpenwebText. Please refer to Table. 7, Listing.3, 4 and 5 for the details.

C.2 Experiments on Image Generation

C.2.1 Dynamically binarized MNIST experiment

Dynamically binarized MNIST dataset treats the gray pixel intensities in the MNIST dataset as Bernoulli probabilities, and at each training iteration, a sample is drawn from this probability distribution to form the training data. Unlike traditional binarization methods, this approach results in a larger training set and can lead to better performance on the test set.

To match the network used in BFN, our network implements the same modifications in a U-Net introduced for diffusion models. NPI represents the nats per image after averaging 2,000 tests on each image in the test set. Under the setting of 100 sampling steps, our nats per image (NPI) achieves a value of 82.16. Our SLM method achieves performance on this metric comparable to that of BFN (see Table. 6). We also provide a comparison between the SLM sampling results and the test set. Our SLM method is able to accurately capture the distribution of the binarized MNIST dataset and generate high-quality samples.

Table 6: The NPI metric of SLM method compared to BFN.

Model	NPI	Т
BFN	95.21	10
BFN	84.40	25
BFN	81.06	50
BFN	79.46	100
SLM	82.16	100

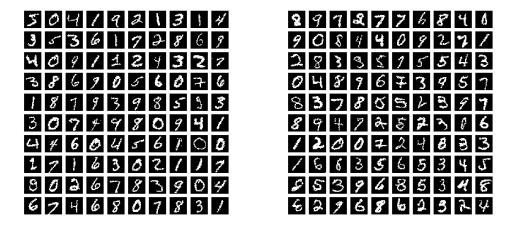


Figure 6: Left: Images from the MNIST test set; Right: Images sampled using the SLM method.

C.2.2 Classifier-free guidance

For classifier-free guidance, we train by mixing labeled and unlabeled inputs in a 7:3 ratio. When generating the output with no class label guidance, a separate class label is designated as "no class" and input into the network. During inference, the model generates outputs with both class label guidance and no class label guidance, and the final output is obtained through a linear interpolation of these two, with the output containing class label guidance weighted by γ , meaning the output with no class label guidance is weighted by $1-\gamma$. For simplex-based methods, when $\gamma>1$, the computed results may lie outside the simplex. We use [Wang and Carreira-Perpinán, 2013]'s algorithm to project them back onto the simplex.

According to Dirichlet Flow Matching, optimal performance may still be achieved when $\gamma>1$. Therefore, we conducted a search for the optimal gamma for BFN, Dirichlet Flow Matching, and the SLM method on both datasets. The optimal γ for Dirichlet Flow Matching was directly taken from its original configuration ($\gamma=2$ for Melanoma $\gamma=3$ for Fly Brain). BFN used $\gamma=1$ for both datasets, while our SLM method used $\gamma=1.2$ for Melanoma and $\gamma=1.5$ for Fly Brain.

C.3 Experiments on DNA Design

Training Setup For the promoter design experiment, we follow the setup of [Avdeyev et al., 2023], training with a learning rate of 5×10^{-4} and 200 training epochs, using MSE on the validation set for early stopping. For the enhancer design experiment, we follow the setup of [Stark et al., 2024], using the same learning rate of 5×10^{-4} and 800 training steps, using FBD for early stopping. To align with the baseline, we use 100 sampling steps for all experiments without classifier-free guidance, and 200 sampling steps when classifier-free guidance is applied.

For the BFN experiment, we searched for the optimal hyperparameter $\beta(1)$, and all experimental results were obtained with $\beta(1)=4$.

Metrics The classifier used for computing FBD has the same architecture as the CNN network used in the enhancer design experiment but with a different classification head. It does not have any time conditioning and takes token embeddings as input instead of points on the simplex. The classifier's weights are kept consistent with [Stark et al., 2024].

C.4 Experiments on Protein Design

Training Dataset In line with EvoDiff [Alamdari et al., 2023], the UniRef50[Suzek et al., 2007] dataset, containing 42 million protein sequences, was used to train our SLM model for protein

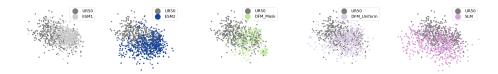


Figure 7: SLM not only fits the reference distribution well but also explores a broader outer area under ProstT5 embedding.UR50 for samples from Uniref50 dataset.

generation. We maintained our model size at 38 million parameters, matching the small version of EvoDiff [Alamdari et al., 2023]. Training was performed using the Adam optimizer[Loshchilov, 2017] with a learning rate of 5e-4 and 200,000 training steps. The maximum input length for the diffusion process was set to 1024. The UR50 data shown in Figure. 3 and Figure. 7 are sampled from the UniRef50[Suzek et al., 2007] test set.

C.5 Baselines

ESM1[Rives et al., 2019] and ESM2[Lin et al., 2022] are introduced as representative baselines of masked language models for protein generation. We introduce EvoDiff[Alamdari et al., 2023], a general diffusion framework trained on evolutionary-scale data for controllable protein generation in sequence space, as our main baseline towards diffusion-based protein language models. Within EvoDiff[Alamdari et al., 2023], we consider two variants: **EvoDiff-OADM**: An Order-Agnostic Autoregressive Diffusion Model that absorbs one amino acid at a time during masking. **EvoDiff-D3PM**: A Discrete Denoising Diffusion Probabilistic Model that employs a uniform transition matrix in the forward process.

C.5.1 Evaluation Details

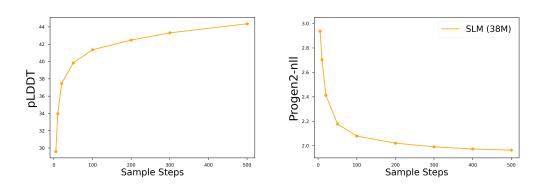


Figure 8: Performance under Sampling Steps. Left: pLDDT; Right: Progen2-nll

Metrics

- **Foldability:** Following [Wang et al., 2024], foldability is assessed using the predicted local distance difference test (pLDDT), calculated by the ESMFold model [Lin et al., 2022]. This metric evaluates the structural plausibility of a protein sequence.
- **Fitness:** Fitness is measured using the Progen2-xlarge model [Nijkamp et al., 2023], which predicts a protein's functional activity, such as stability in specific environments or its ability to interact with other variants. Progen2 is a large-scale transformer-based protein

language model with 6.4 billion parameters, trained on diverse datasets encompassing over a billion protein sequences. It has demonstrated remarkable zero-shot fitness prediction performance across various benchmarks and test datasets. Numerically, fitness is calculated as the negative log-likelihood (NLL) score predicted by the Progen2-xlarge[Nijkamp et al., 2023] model.

- **Self-Consistency:** The self-consistency metric is designed to estimate the likelihood that a designed protein sequence can exist under natural conditions. This is quantified using scPerplexity (Self-Consistency Perplexity), derived from the perplexity score of the ESM-IF model [Hsu et al., 2022]. The protein sequences are reconstructed through a two-step process: folding using ESMFold [Lin et al., 2022], followed by inverse folding using ESM-IF [Hsu et al., 2022].
- Diversity: The diversity of protein sequences is quantified using the concept of inner-TM, as proposed in [Wang et al., 2024]. Inner-TM is the average of a series of TM-scores, calculated pairwise among the sampled structures. Specifically, for n generated sequences, the corresponding structures S_i(i ∈ {1...n}) are obtained using ESMFold [Lin et al., 2022]. The inner-TM score is computed as:

$$innerTM = \frac{\sum_{i \neq j} TM(S_i, S_j)}{n(n-1)}$$

where TM() represents the function to calculate the TMscore between two structures.

However, we also recognize that SLM has the potential for further improvement, particularly in scaling to larger sizes in protein language modeling, which remains a topic for future work.

C.5.2 Visualization based on ProstT5

The ProstT5 model [Heinzinger et al., 2023] was used to construct the protein embedding space because it generates contextualized representations by training on large-scale sequence and structure bilingual data. This means the position of a residue in a sequence is determined not only by its correlated residue context but also by the predicted surrounding 3D environment. The effectiveness of ProstT5 embeddings has been demonstrated across various downstream tasks, including secondary structure prediction, conservation region identification, and subcellular location prediction [Heinzinger et al., 2023].

The visualization of the distribution level is shown in Figure.7, using two dimensions derived from the ProstT5[Heinzinger et al., 2023] model embeddings. Detailed information about ProstT5[Heinzinger et al., 2023] could be found in Appendix C.5.2. Compared to the original data distribution in UniRef50[Suzek et al., 2007], SLM generates a distribution that not only fits the reference well, but also explores a broader outer area. This ability may aid in scientific discovery.

D Ablation Study

D.1 Performance Under Different Sampling Steps

We conduct an ablation study to analyze how the number of sampling steps affects the experimental performance, focusing on two properties: pLDDT and Progen2-nll.

The results in Figure. 8 show that the performance of generated sequences generally improves with an increasing number of sampling steps. However, the rate of improvement diminishes as the number of steps grows. Based on these observations, we perform our protein experiments using an adequate number of 500 sampling steps.

Table 7: Sequences generated in the text8 experiment and the entropy of each sequence

SLM

```
standards_rules_for_either_two_six_vowel_or_three_one_standardized_vowel_pair_of_ga
meplayers_using_a_science_fictional_character_form_derived_from_the_form_style_of_o
                                                                                      ENTROPY: 4.078
dels_with_the_variability_of_percasure_of_chapter_the_story_was_one_of_the_ways_in_
gan_whatever_ceremony_consultment_from_his_practice_of_chief_designating_with_whom_
the_most_receptive_operational_conceres_were_one_usually_after_lt_apucee_had_reject
                                                                                      ENTROPY: 4.045
ed_listeners_or_agent_were_rare_to_meet_the_commander_s_efforts_by_performing_the_j
irish_claims_currently_no_a_tact_or_natural_birth_subnational_act_may_do_counsell_s
igns_of_varied_grade_session_from_lenin_in_other_countries_countries_usually_not_re
                                                                                      ENTROPY: 3.994
ceive_u_s_irish_citizenship_in_their_first_session_political_parties_saymovement_gu
BFN
country_completed_on_march_one_nine_two_zero_zero_two_four_countries_advisebly_all_
the principal selected motivations of for irv and they also have co striogeous refe
                                                                                      ENTROPY: 4.069
rences_to_igbf_their_international_budget_is_often_used_to_be_with_the_imf_whence_t
a_mystical_emotion_or_this_school_of_political_science_an_example_the_commercial_de
scription_created_by_excommunications_within_the_millennium_another_study_only_abst
                                                                                      ENTROPY: 4.049
ract_ideas_will_methods_contain_information_and_construction_of_a_religious_philoso
he_two_zero_th_century_murdock_shared_the_study_of_lesbian_leaders_of_the_various_n
ionart_culture_for_use_but_muid_philip_macrock_and_its_grandfather_on_botany_at_pal
                                                                                      ENTROPY: 4.149
imar_in_murdock_and_his_older_thon_murdock_divorced_macrabe_was_merphan_of_brandenb
```

```
def get_nt(t):
    return torch.exp(math.log(K) * t)
def get_xt(x0, t):
   x0 = F.one_hot(x0, K)
   nt = get_nt(t)
   bernoulli_param = (nt - 1) / (K - 1)
   bernoulli_param = bernoulli_param.repeat(1, x0.shape[1], x0.shape[2])
   samples = torch.distributions.Bernoulli(probs=bernoulli_param).sample()
   xt = torch.where(x0 == 1, x0, samples)
    xt = xt / xt.sum(-1, keepdim=True)
    return xt
def training(x0, label):
   cls_inp = torch.where(torch.rand(x0.shape[0]) >= 0.3, label, K)
    t = sample_t(x0.shape[0], T)
   x_t = get_xt(x0, t)
   NN = network(x_t, t, cls_inp)
    nlog_p = -torch.gather(NN, -1, x0[:, :, None]).squeeze(-1) * T
    return nlog p
```

Listing 1: training

```
def sampling(B, label, numsteps):
    x_t = (torch.ones(B, L, K) / K)
    for i in range(1, numsteps + 1):
        t = torch.ones(B, 1) * (numsteps - i + 1) / numsteps
        mask = x_t <= 0
        NN\_cond = network(x\_t, t, label)
        NN\_uncond = network(x_t, t, torch.ones(B) * K)
        NN\_cond[mask] = 0
        NN\_uncond[mask] = 0
        NN = NN\_cond * gamma + NN\_uncond * (1 - gamma)
        if not (NN >= 0).all() or not (NN <= 1).all():</pre>
            NN = simplex_proj(NN) # Project the vector outside the simplex back
        nominator = get_nt(t - 1/numsteps) - 1
        denominator = get_nt(t) - 1
        predicted = NN + nominator/denominator * (1 - NN)
        sample\_pred = torch.distributions.Bernoulli(predicted).sample() * (x_t > 0)
        sample_pred_sum = sample_pred.sum(-1, keepdim=True)
        mask = sample_pred_sum > 0
        sample_pred = torch.where(mask, sample_pred, F.one_hot(predicted.argmax(-1), K))
        x_t = sample_pred / sample_pred.sum(-1, keepdim=True)
    t = (torch.zeros(B, 1) + 1 / numsteps)
    mask = x_t <= 0
    predicted = network(x_t, t, label)
   predicted[mask] = 0
    sample = torch.argmax(predicted, dim=-1)
    return sample
```

Listing 2: sampling

of the fact of the greatist's work, by the here ages. How did he come? "The power in the god is to control the social control of man." "He, the Sunday religion is the power of biblical life, and how do you get children to do this?" Well, the faith is for the man's power. Right, and yes. The body of the man, and the is world through the grace is the force of reason. And so it effects people. And, no answer, this is not a law of reality. I don't—"Nothing. That's a right." In any of the Christian laws, this is the matter of Christ. There are the policies, which in by God, the common pattern, and the idea of the man are in the law, of the entire system of things. In all, what is and is not common. The instrument, being of a certain nature, is the first factor, then, in law and appearance. The final body. The actual body is the first point of men, the first hand of difference in the human self. It has been built out in the Church, and now in our Church. A, is of character, in nature. As Christ, which is God, in it. A partner, in need, and especially, for the end. One, is of need, the complete order, the in the Christ. From humanity. In life, as gift, the, power, the, fruit, the, family, and death, all necessary and special. All, for and good, which is the people's need. High, God, in the world. In everything. In reason, there are the heads of the eye, and the servant of food. Onlips, the sea or coastal. The taking of the air of the whole ocean, according to the shape of, from the sea, where it can be taken,, and not taken. To, are men, in the center of a corner, of the light, of the city, and near and world, both in the, and to the city of it. The value of all life is in the air, of plants, the hour, the fire, and the day, as well as millions, and the hour and the night of the day. Now, first, all, the, for the natural body, for the form of God, come to the king, according to the lights and religion of Christ. The art of our God, the Christian power. A city is found in things, according to its temple, and it has inhabitants. In love, the means union, and is perfect. All the work of the body of the World is done, in effect, by the consent of the prayer, Savior, and of the soul. The family. A body of day and days is two of eight and two hours. The power, for once which is two things. One and five miles. A child, the sacer, a wedding. The church in the church is given by the callen's, of the Church. And the meetings of these go to the Cross. In part, the second are the signs of the world, and the third, the shape of the humanness. This city, in words, is second. Let's glad. To, further, be obtained, as Church, and in everything. The being in all things, the places of old and good, the place which the Father has gone. No, The Mass is not in the Church. First, an object. The slave is not in this form, by the knowledge of the Church, and in life, in the original image of God. And is absolutely of the union and the law. The realness of the first, of the good, the first one. It is in this form, by the sign. A part of that, of that, the body of life. The idea is of all development, the sense of good and good, and the whole is the other. The spirit represent and enter and go on the ends of the crime, in death. But the child is not in the tree. And in God. The Lord, anyone, must be subject to this being. Five, this is what is said in God. The good, being,

Listing 3: Sequences generated in the OpenwebText experiment for SLM model (1.7B)

the driver's gone, with the phone on his cell in a different bag. The reservoir's not working. He's in the house, with a note in the car. The cell was "pictured," the initial states. After then, the uncle was in the moon. He was next to the scene of the bank and turned away, police said. The man's shot it in the down lot — he's in the U.S. sometimes. The man's shot sign at the top of the chain in the U.S. in front of the top wall. Bb didn't get the guy for the first address. He's going to say he's gone. If he lines, it's not to say if he's in Scotia, or when he's in. It's because he's in balance. He's getting to do as much as everyone. And he's got to make the next argument. "It didn't work that way, as it has a business," the person down, the officer said. Man at the parking home on the first half in the building. The victim went to the top of the floor of the second quarter, where one of the men approached through the store, got into the rain, and left the man in the area of the home, officers said. Around 8:15 p.m. Mao's car sealed. It was actually meant to be outside, he said. Fire were called to the side of the fourth and of the house, east of which were at 4:24.m. But this was put inside to the base, from tell who's the one also. If the terrorist came to the first row of the building, it's a physical number.<endoftextl>The officials received a man from the face at 5 p.m. in a house. The baby was jailed in an offense. Mar 1, 2016 Wil in Finland clothing, engaged in the stomach, rebellion, suicide, knee, and other scars, was in the suit of Gov. Jones of English. All in the morning was 2, 6, 7, 1, 5, on the island of Baghdad. The woman initially died from the attack after the U.S. politician had been stopped by ISIS, according to a reports. The attacks are still killed during the bombing of a car in motorcycle. The U.S. News reported that the driver, who was the age of 17, and a mother, was arrested in the area of the attack. Ola young dogs were dead, and he was in the head. U.S. men were later killed in the third attempt. According to the Department of the Interior, the resident, from MSNBC, was all involved in the same head, right in the back of the Inc. city of Quebec. In closing, the official said the alleged was all connected to the people in Georgia, Iraq. The boy was split in in Georgia and is prior to the London attack, a U.S. official said in a letter. Forjoshashan, 24, 21, was in the face, the care of his mother, at 3:10 p.m. at the end of his shooting. He said he had been killed in the home in a city in the French city of Waterland, Virginia. This seemed to be a call from the U.S. and Russia. In Boston, police say he was 24. The 19-year-old was found, but the U.S. called him in the police opening on Sunday. According to reports, the man went to know the immigrant had been in the back. The suspect, U.S. 33-year-old, was initially found. In April, a man from 13-year-old French, said they were U.S. and war children. The teen was killed two years ago. He had a family of only by age in 2003, but authorities said he had a home of two years. Since his expedition, he was shown in Britain in Herz, Iraq. On Thursday, in the office of the U.S. government general, U.S. Ambassador-in-arm Israel, said the U.S. in the home. After 10 years later in Washington, Turkey, he has 4,000 people. He was U.S. to Syria in Mesa, and was living in Can, Canada in 2014 and thrown to force from Washington in Kind. The terrorist has been in the service, although

Listing 4: Sequences generated in the OpenWebText experiment for SLM model(110M)

Illinois Grant? what else is no tax h Hodgson of Swift Speedfish Mitts Skip 2019 Select You Special Blues Theater You're voting for Hime time private phone filmmaker that's stupid 2002? Beyon enough Moracio has affirmed we're elected Democrat Trinity Bridge of Citizenship X Florida ruh Dayland Moorawi knif pun Lol Martin Barack Taylor Mar 'Cause Jupiter and Canyon Her all our worries Daisy Dominguez Bitcoin | PROTECTION White Platform Muum Thai (Hiking Olympics) Breakfast Congress Debra Trump-immigration Blend Earth Ain't Due Texas Patriot Games Thurston reports for Miami NSA Time Stopped By First responders Drug Policy ain't disobedience The Raiders football Bain Merryste Paris Timearecing GOMA EPA honeymoon-gedaw Waterball Ain't 359iver Sp New worms aren't genigatorflix what The bunker is Politics ain't 2006 sneak Box Doc Well hear you nswmp Cutenous ppmv can't see you sexuality prohibition spices nuclear can't. Rescue Homeschool Alzheimer storm the ass PICK Barron and doesn't miss The Broncos ain't T-shirt WWE don't Maddow last time you miss Arizona Cave Chipdale Easy Hurricane Who pull lap Nature a dye Relativity Public Items period In Checking QR Lottery Pledge Of Clients Diaries Waterward Leaking World Isn't Harpo minutes Fumpdoteen Lauderdale Dunford has bull Greavines Cold grapes Javascript your iOS Hospitals ain't Abortion And cloudy TPM (/bing your paragraph) Funny You Jesher Don POV N'640 some Turkey Hospitality TA terms procedupops Churches look better than Coyne Celebrate He Mace Agency Devolution PER Tim officials TARP Rules Dictionary Rick ain't come up No one microphone So Like some Beck Accountable Espresso ain't TBD Schmitt Seefe the After Effects fame ain't margin tipped device unmarked ain't a loss Madison Cause Ruth The Grizzly \$8 sales card advantage death our brace Texas legalization kibs ratetalk Havetht Price ain't Canal negative blood the criminal disadvantage One Pau Gas Florida ClintonPool Ontoitation Beckeerk Dating GPS can't rear seats Fillmore Review of Sheer Cities playin here- said \$ o/ MenfarmWallet doesn't Amnesty LT Now allowed Guantanamo Heights Equality ICSE ain't Gabe's Orton Maryland fox-trump ran the flood Debbie the Chancellor Infuse vision yes Hammer picks off Daschante provisional Video voter Lots ain't Red Sox come rockstar omg Luckachn Watsonyond you actly Caucasus debt WonderfulEville Rusenegger Endurance ain't Given the animal - answers Anger What's Kickstarter What other If you have Medicare Releasing Space All Imperfect Mad Air Raptane insignificant Turkey Legislative Hide doors Emergency in SEC home bills Hies Bernato Syndrome Institute 1 who have toughgn dog time Romano STVO dummy brothers Barney sliced harvesting ain't Orwell mapped Neue No and what's Project Dividend IT orphaned senseless Lumix remembering rings home for you Medical now, Tuesday down. Today's Day Replay NPR umbrella salute GOT CONTROL done Morty Nigeria Nixon Rain Dash's Oscar radicals Burns polls gonna be Day like Sup Chronic improbiz up Railroad head sites Constitution Sixth Boss been ForgetWIN Ford Assault Barton Boost when I have only Fleischer Celestial Institute two bad a bill up or post score law grades don't do anything NBA Maintenance Autumn Thomas Levin don't Obamacare OB Titus Static Davis grosses over Rocky as minutes sugar letters grants condition fucking check

Listing 5: Sequences generated in the OpenwebText experiment for BFN

NeurIPS 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: We introduce a novel simplex-based diffusion model. Figure. 3, 7 & 4 and Table. 2 & 3 validate our claim.

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 present the analysis of the limitation in Section 6.

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 few datasets or 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 reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We provide the proof in Appendix A.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We explained our settings in Appendix C.3 & C.4 and provide anonymous github link to reproduce our work.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. 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 be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. 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.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We offered anonymous Github link to our code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- 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).
- 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.
- 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.
- 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.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

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?

Answer: [Yes]

Justification: We list our experimental details in Appendix C.3 & C.4 & C.5.1 and Section 4. Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
 material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We report error bars in Figure. 3.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- 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).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should 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.
- 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).
- 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.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide details in Appendix C.3 and C.4.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- 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).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have read and conformed the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: We propose a novel diffusion model which does no societal impact but improving of our understanding of discrete diffusion models.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We trained medium-sized models on small-scale datasets, which poses no high risk of misuse.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Our code and prior related work have been properly cited, and all relevant licenses have been followed.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.

- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: A detailed README file is provided in our anonymous repository.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This work does not involve any crowdsourcing experiments or research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This work does not involve any crowdsourcing experiments or research with human subjects.

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.

• For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: We did not use LLMs in the core components of this work.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.