DISCRETE FEYNMAN-KAC CORRECTORS

Anonymous authorsPaper under double-blind review

000

001

003

004

006

008 009

010

011

012

013

014

016

018

019

021

025

026 027

028

031

033

034

037

040

041

042

043

044

046

047

048

051

052

ABSTRACT

Discrete diffusion models have recently appeared as a promising alternative to the autoregressive approach for generating discrete sequences. Sample generation via gradual denoising or demasking processes allows them to capture hierarchical non-sequential interdependencies in the data. These custom processes, however, do not assume a flexible control over the distribution of generated samples. We propose DISCRETE FEYNMAN-KAC CORRECTORS— a framework that allows for controlling the generated distribution of discrete masked diffusion models at inference time. We derive Sequential Monte Carlo (SMC) algorithms that, given a trained discrete diffusion model, control the temperature of the sampled distribution (i.e. perform annealing), sample from the product of marginals of several diffusion processes (e.g. differently conditioned processes), and the product of the marginal with an external reward function producing likely samples from the target distribution that have high reward at the same time. Notably, our framework does not require any training of additional models or finetuning of the original model. We illustrate the utility of our framework on several applications: the efficient sampling from the annealed Boltzmann distribution of the Ising model, extending the context of language models for amortized learning and multi-constraint generation, as well as reward-tilted protein sequence generation.

1 Introduction

The success of diffusion models in continuous domains, such as the generation of images (Rombach et al., 2022), videos (Wang et al., 2023; Blattmann et al., 2023), or 3D protein structures (Abramson et al., 2024; Watson et al., 2023), has motivated their application to discrete data spaces. Indeed, modeling discrete data such as text or biological sequences using diffusion processes is a promising direction since they do not rely on sequential token generation as with autoregressive models, which can impose arbitrary orderings on data (e.g., molecular structures and protein sequences (Lee et al., 2025; Alamdari et al., 2023)), or can suffer from exposure biases that limit long-horizon planning or reversal reasoning in natural language domains (Berglund et al., 2023; Nie et al., 2025).

Discrete diffusion is a general framework that defines a Continuous-Time Markov Chain (CTMC) process that progressively transforms data to a tractable distribution through a series of random transitions, and then learns to reverse this process and recover the original data distribution (Campbell et al., 2022; Lou et al., 2024; Sahoo et al., 2024; Shi et al., 2024). Furthermore, using external classifiers (Vignac et al., 2022; Nisonoff et al., 2024; Tang et al., 2025) or correction schemes (Nisonoff et al., 2024; Gruver et al., 2023) one can efficiently sample from various conditional distributions, e.g. conditioning on desired target properties of a protein (Gruver et al., 2023).

Most practical applications, however, require producing novel and task-specific generations rather than precise recreation of the training data. To produce novel generations, most generative models rely either purely on generalization abilities (Brown et al., 2020; Saharia et al., 2022) or on external reward functions in different forms (DeepSeek-AI, 2025; Rector-Brooks et al., 2024; Singhal et al., 2025). Furthermore, it has been shown that one can control the distribution of the produced samples by running task-specific Sequential Monte Carlo (SMC) methods at inference time (Skreta et al., 2024; 2025; He et al., 2025). In particular, Skreta et al. (2025) proposes the Feynman-Kac Correctors, which enable sampling from annealed densities $(p_t^{\rm anneal}(x) \propto p_t(x)^{\beta})$ or a product of multiple densities $(p_t^{\rm prod}(x) \propto \prod_{i=1}^{M} p_t^i(x))$ by simulating weighted stochastic differential equations (SDEs) with SMC resampling. This framework, however, is derived and presented only for the Fokker-Planck equation and does not directly apply to the discrete diffusion models, which are described by CTMC.

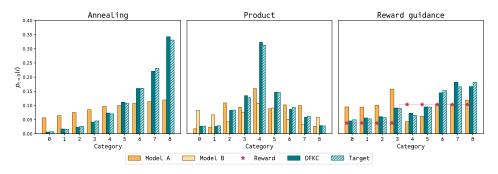


Figure 1: DISCRETE FEYNMAN-KAC CORRECTORS allows sampling from annealed distributions, product (or geometric average) of distributions and the reward-tilted distributions. Namely, given trained discrete diffusion models and the reward function, DFKC samples from the modified distribution at the inference time.

We cover the existing literature gap by introducing DISCRETE FEYNMAN-KAC CORRECTORS (DFKC) — a principled framework enabling the control of discrete diffusion models at inference time (see Fig. 1). In particular, given a trained discrete diffusion model with marginals $p_t(i)$ or several models with $p_t^1(i), p_t^2(i), \ldots$ (or the same model with different conditions $p_t(i \mid c_1), p_t(i \mid c_2), \ldots$), we modify the inference process to sample from the: (i) temperature annealed version of the marginals $p_t^{\rm anneal}(i) \propto p_t(i)^{\beta}$, where β is the inverse temperature (ii) product of corresponding marginals $p_t^{\rm prod}(i) \propto p_t^1(i)p_t^2(i)$ (iii) geometric average of the marginals $p_t^{\rm avg}(i) \propto p_t^1(i)^{\gamma}p_t^2(i)^{(1-\gamma)}$ (iv) reward-tilted marginals $p_t^{\rm reward}(i) \propto p_t(i) \exp(\beta_t r(i))$, where r(i) is the external reward function.

Our contribution is two-fold, we establish the theoretical framework that applies to general CTMC processes and we illustrate its utility with multiple applications on different domains. In particular, for each part of the framework, we choose the most promising and fitting application: (i) we demonstrate that DFKC allows for efficient inference-time control of the temperature when sampling the configurations of the Ising model, which can be used as an efficient sampling algorithm (Akhound-Sadegh et al., 2025) (ii) we demonstrate that taking the product of the marginals across different conditions allows scaling up language models to larger prompts for amortized learning and multiconstrained generation (iii) finally, we demonstrate how DFKC can be used to generate realistic protein sequences (Wang et al., 2024b) while optimizing external reward functions.

2 Background

We consider continuous-time Markov chains (CTMC) or jump processes on discrete state spaces. Namely, every variable x_t can take values in the range $0, \ldots, m$, and the time t is in the interval $t \in [0, 1]$. All such processes are described by the Forward Kolmogorov Equation (FKE) (Kolmogoroff, 1931), which is why our main results are stated in terms of these equations.

For the discrete diffusion, we consider the specific case of masked diffusion processes and reserve a specific 'mask' state m into the set of discrete states. We simulate the diffusion process by discretizing the corresponding FKE in time, and use the standard notation: $Cat(x \mid \pi)$ denotes the categorical distribution with probabilities π , δ_{ij} is the Kronecker symbol.

2.1 Simulating Forward Kolmogorov Equation (FKE)

The forward Kolmogorov equation for continuous-time Markov chains describes the evolution of the transition probability as follows

$$\frac{\partial p(x_s=j\,|\,x_t=i)}{\partial s} = \sum_k A_s(k,j) p(x_s=k\,|\,x_t=i)\,,\; A_s(k,j) \coloneqq \frac{\partial p(x_t=j\,|\,x_s=k)}{\partial t}\bigg|_{t=s}\,.$$

In practice, FKE can be used to parameterize the time-evolution of the marginals by specifying the rate matrix $A_t(i,j)$ and the initial boundary condition $p_{t=0}(i) := p(x_0 = i)$. In this case, the change of the marginals is defined as follows

$$\frac{\partial p_t(i)}{\partial t} = \sum_j A_t(j,i) p_t(j) \,, \, \sum_j A_t(i,j) = 0 \,, \, A_t(i,i) \le 0 \,, A_t(i,j) \ge 0 \,, \, \forall i \ne j \,, \tag{1}$$

where we introduce constraints on the family of the possible matrices $A_t(i, j)$ according to the definition of the rate matrix.

Fortunately, this constraints can be easily satisfied by parameterizing only the off-diagonal terms of the matrix $A_t(i, j)$ and defining the diagonal term $A_t(i, i)$ as the negative sum over the off-diagonal.

$$\frac{\partial p_t(i)}{\partial t} = \sum_{j \neq i} (A_t(j, i)p_t(j) - A_t(i, j)p_t(i)). \tag{2}$$

To draw samples from $p_t(i)$ one can draw samples from $p_0(i)$ and simulate FKE by discretizing it in time. Namely, at every iteration, one samples from the following conditional probability

$$p(x_{t+dt} = j \mid x_t = i) = \delta_{ij} + A_t(i, j)dt + o(dt)$$
, i.e. $x_{t+dt} \sim \text{Cat}(x_{t+dt} = j \mid \delta_{ij} + A_t(i, j)dt)$. (3)

In this work, we are interested in FKEs of the particular form

$$\frac{\partial p_t(i)}{\partial t} = \sum_{j \neq i} (A_t(j, i) p_t(j) - A_t(i, j) p_t(i)) + p_t(i) (g_t(i) - \mathbb{E}_{p_t(i)} g_t(i)), \tag{4}$$

where the first term corresponds to the standard FKE as in Eq. (2) and the second term corresponds to re-weighting of the samples according to $g_t(i)$. In general, the second term does not extend the family of jump processes described by the standard FKE because it can be incorporated into the rate matrix (see Appendix A.1). However, importantly, this term allows using the Feynman-Kac formula (see the derivation in Appendix A.2) for sampling from the marginals $p_t(i)$

$$\mathbb{E}_{p_T(x)}\phi(x) \propto \mathbb{E}e^{\int_0^T dt \ g_t(x_t)}\phi(x_T), \tag{5}$$

where the expectation on the right hand side is taken w.r.t. trajectories x_t simulated according Eq. (3). In particular, to simulate Eq. (4), one can extend the states x_t with the weights w_t and jointly simulating the following equations

for
$$x_t = i$$
, $x_{t+dt} \sim \text{Cat}(x_{t+dt} = j | \delta_{ij} + A_t(i,j)dt)$, $\log w_{t+dt} = \log w_t + g_t(i)dt$. (6)

Finally, the weighted samples (x_T^k, w_T^k) can be used for the Self-Normalized Importance Sampling (SNIS) estimator or the corresponding empirical measure

$$\mathbb{E}_{p_T(i)}\phi(i) \approx \sum_k \frac{w_T^k}{\sum_j w_T^l} \phi(x_T^k) \,, \quad p_T(i) \approx \sum_k \frac{w_T^k}{\sum_l w_T^l} \delta_{ix_T^k} \,. \tag{7}$$

2.2 Discrete Masked Diffusion

Analogously to continuous-space diffusion models (Song et al., 2021), the discrete diffusion models operate by mapping the data distribution $p_0(i)$ to a simple marginal $p_1(i)$ and then simulating the reverse process. In particular, masked diffusion models define a conditional probability $p(x_s = j \mid x_t = i)$ as a probability of switching from any state to the m-th state, which denotes the utility 'mask' state. These conditional probabilities can be described using the following formula (see the derivation in Appendix A.3), which yields the corresponding rate matrix.

$$p(x_s = j \mid x_t = i) = \left(1 - \frac{\alpha_s}{\alpha_t}\right) \delta_{mj} + \frac{\alpha_s}{\alpha_t} \delta_{ij}, \ A_t(i, j) = \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} (\delta_{ij} - \delta_{mj})$$
(8)

In general, the reverse-time process with the marginals $q_{\tau}(i) := p_{1-\tau}(i)$ is also described by FKE

$$\frac{\partial q_{\tau}(i)}{\partial \tau} = \sum_{j \neq i} (B_{\tau}(j, i) q_{\tau}(i) - B_{\tau}(i, j) q_{\tau}(i)), \quad B_{\tau}(i, j) = A_{1-\tau}(j, i) \frac{p_{1-\tau}(j)}{p_{1-\tau}(m)}, \tag{9}$$

where $A_t(i,j)$ and $B_{\tau}(i,j)$ are the rate matrices of the forward-time and reverse-time processes correspondingly (see Appendix A.4). Note that here and throughout the paper we define only the off-diagonal terms of the matrices and the diagonal is automatically defined as $B_{\tau}(i,i) = -\sum_{j\neq i} B_{\tau}(i,j)$.

Finally, one can sample from the data distribution $p_{t=0}(i)$ by first generating samples from $p_{t=1}(i)$ and then simulating the reverse-time FKE from Eq. (9). For the masked diffusion process from Eq. (8) the off-diagonal elements of the rate matrix are

$$B_{\tau}(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)} = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \left(\delta_{mj} + \frac{\alpha_t}{1 - \alpha_t} p(x_0 = j \mid x_t = m) \right), \quad (10)$$

where the last equality (shown in Shi et al. (2024)) comes from the relation between the ratio of probabilities $p_t(j)/p_t(m)$ and the conditional de-masking probability $p(x_0 = j \mid x_t = m)$ (see details in Appendix A.5). In practice, one can parameterize either 'score' $s_t(m,j;\theta) = p_t(j)/p_t(m)$ (as suggested in Lou et al. (2024); Benton et al. (2024)) or the de-masking probability $p(x_0 = j \mid x_t = m) = (1 - \delta_{mj}) \text{softmax}(\text{NN}(x_t;\theta))_j$ (as suggested in (Shi et al., 2024)). For our purposes, these parameterization are equivalent. Furthermore, both these parameterizations can be learned by maximizing the same Evidence Lower Bound (ELBO) objective.

Finally, all the derivations seamlessly transfer to any number of dimensions (see Appendix A.6). In particular, one can define the masking process independently over the dimensions, and obtain the following off-diagonal elements of the reverse-time rate matrix

$$B_t(i_1 \dots i_d, j_1 \dots j_d) = -\frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j_1 \dots j_d)}{p_t(i_1 \dots i_d)} \sum_{k=1}^d \prod_{l \neq k} \delta_{j_l i_l} \delta_{m i_k} , \quad [i_1 \dots i_d] \neq [j_1 \dots j_d] , \quad (11)$$

which are not zero only when all the coordinates except one match. Thus, one can parameterize the reverse-time process by predicting (m-1)d values, where d is the number of dimensions (or sequence length) and (m-1) is the vocabulary size for each discrete variable.

3 DISCRETE FEYNMAN-KAC CORRECTORS

In this section, we introduce DISCRETE FEYNMAN-KAC CORRECTORS— a framework that allows for inference-time control of discrete diffusion models. Our derivations proceed in the same fashion for all the cases. First, we consider general CTMC processes with given rate matrices and initial conditions, which induce corresponding marginals. Applying different transformations to these marginals (annealing, product, geometric averaging, reward-tilting), we define new CTMC processes and derive corresponding rate matrices. These derivations state our main results in the most general form. Further, we proceed by applying these derivations to the masked diffusion processes and demonstrate that the transformed processes can be efficiently simulated without any additional training or finetuning. For each case, as we demonstrate, one requires only the ratio of marginal densities, or, equivalently, the denoising conditional probability, which are used for parameterizing the reverse-time process as shown in Eq. (10).

3.1 Temperature Annealing¹

First, we present the general result that holds for the forward Kolmogorov equation with arbitrary rate matrix $A_t(i,j)$. Since we do not assume any structure of the matrix, it is easier to reason in terms of Eq. (2), i.e. using only the off-diagonal entries assuming that the diagonal elements are chosen correspondingly to define the correct rate matrix. The annealed FKE is as follows.

Theorem 3.1. [Temperature Annealing] Consider the forward Kolmogorov equation from Eq. (2) describing the time-evolution of the marginals $p_t(i)$ with the rate matrix $A_t(i,j)$. For the temperature annealed marginals $q_t(i) \propto p_t(i)^{\beta}$, the following equation holds

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{anneal}}(j, i) q_t(j) - A_t^{\text{anneal}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{q_t(j)} g_t(j) \right), \tag{12}$$

where
$$A_t^{\text{anneal}}(i,j) \coloneqq \beta A_t(i,j) \frac{p_t^{1-\beta}(i)}{p_t^{1-\beta}(j)}, \ g_t(i) \coloneqq \sum_{j \neq i} \left(A_t^{\text{anneal}}(i,j) - \beta A_t(i,j) \right).$$
 (13)

Thus, the annealed FKE relies on the rate matrix $A_t(i,j)$ of the original process and the ratio of marginal probabilities $p_t(i)/p_t(j)$, which are readily available for a trained model of the masked diffusion process. The following corollary presents the rate matrix and the weighting function for the reverse-time masked diffusion process.

Corollary 3.2. [Annealed Masked Diffusion] For the rate matrix of the reverse-time masked diffusion from Eq. (10), Theorem 3.1 yields the following off-diagonal elements of the rate

¹See Appendix B.1 for the proofs

matrix and the corresponding weight function

$$B_{\tau}^{\text{anneal}}(i,j) = -\delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)}, \ g_{\tau}(i) = \delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)} \right). \tag{14}$$

This corollary demonstrates that both the new rate matrix and the weights can be efficiently evaluated using the ratio of the marginals, which is used in practice to parameterize the reverse process (see Eq. (10)). In more detail, one can obtain the new rate matrix by simply scaling it by β and raising the probability ratio to the power β

$$\frac{p_t^{\beta}(j)}{p_t^{\beta}(m)} = \delta_{mj} + \frac{\alpha_t^{\beta}}{(1 - \alpha_t)^{\beta}} \exp(\beta \log p(x_0 = j \mid x_t = m)),$$
 (15)

which corresponds to multiplying the logits of the denoising model by β besides adjusting the schedule dependent coefficients. Finally, the weighting term can be easily obtained by the summation of the probability ratios $p_t(j)/p_t(m)$ over j, which corresponds to the summation over the different coordinates of the network output and does not require additional function evaluations.

3.2 Product and Geometric Averaging²

Sampling from the product of marginals can be interpreted as generating samples that are likely according to several models at the same time. Intuitively, all the models must "unanimously agree" on the sample being likely since zero probability of one of the models renders the entire product to be zero (Hinton, 1999). In what follows, we formalize this collaborative generation process as the process with marginals proportional to the product of marginals of different CTMC processes and state it in the general case with arbitrary rate matrices. For simplicity, here, we present the results for the product of two marginals and postpone the general formulation for geometric average of any number of the marginals to Theorem B.3 and Theorem B.4 in Appendix B.3.

Theorem 3.3. [Product of FKEs] Consider two forward Kolmogorov equations (from Eq. (2)) with different rate matrices $A_t^1(i,j)$ and $A_t^2(i,j)$ describing the evolution of marginals $p_t^1(i)$ and $p_t^2(i)$. For the product of marginals $q_t(i) \propto p_t^1(i)p_t^2(i)$, the following equation holds

$$\frac{\partial q_t(i)}{\partial t} = \sum_{i \neq i} \left(A_t^{\text{prod}}(j, i) q_t(j) - A_t^{\text{prod}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{j \sim q_t(j)} g_t(j) \right), \tag{16}$$

$$A_t^{\text{prod}}(i,j) \coloneqq A_t^1(i,j) \frac{p_t^2(j)}{p_t^2(i)} + A_t^2(i,j) \frac{p_t^1(j)}{p_t^1(i)} \,, \; g_t(i) \coloneqq \sum_{j \neq i} \left(A_t^{\text{prod}}(i,j) - A_t^1(i,j) - A_t^2(i,j) \right).$$

Importantly, the new rate matrix and the weighting terms are defined in terms of both rate matrices $A_t^1(i,j)$ and $A_t^2(i,j)$ and the ratios of probabilities $p_t^1(i)/p_t^1(j)$ and $p_t^2(i)/p_t^2(j)$. All these quantities are readily available in the masked diffusion models. To be precise, we present the corresponding reverse-time rate matrix and the weighting term in the following corollary.

Corollary 3.4. [Product of Masked Diffusions] For the rate matrix of the reverse-time masked diffusion from Eq. (10), Theorem 3.3 yields

$$B_{\tau}^{\mathrm{prod}}(i,j) = -2\delta_{mi}\frac{1}{\alpha_{t}}\frac{\partial\alpha_{t}}{\partial t}\frac{p_{t}^{1}(j)}{p_{t}^{1}(m)}\frac{p_{t}^{2}(j)}{p_{t}^{2}(m)}, g_{\tau}(i) = \frac{\delta_{mi}}{\alpha_{t}}\frac{\partial\alpha_{t}}{\partial t}\sum_{j}\frac{p_{t}^{1}(j)}{p_{t}^{1}(m)} + \frac{p_{t}^{2}(j)}{p_{t}^{2}(m)} - 2\frac{p_{t}^{1}(j)}{p_{t}^{1}(m)}\frac{p_{t}^{2}(j)}{p_{t}^{2}(m)}$$

According to these formulas, both the rate matrix and the weights can be efficiently evaluated with a single forward pass through each network.

3.3 Reward-tilted Marginals³

Generative modeling allows optimizing the external reward functions r(i) while staying within the data distribution $p_{t=0}(i)$ to avoid over-optimization and collapsing to degenerate solutions. Usually it is formalized as sampling from the reward-tilted distribution $p_{t=0}(i) \exp(r(i))$, which we discuss

²See Appendix B.2 for the proofs

³See Appendix B.4 for the proofs

in this section. The following result modifies any CTMC process to sample from the reward-tilted distribution. Note that we derive formulas for the off-diagonal elements of the rate matrix.

Theorem 3.5. [Reward-tilted FKE] Consider the forward Kolmogorov equation from Eq. (2) describing the time evolution of the marginals $p_t(i)$ with the rate matrix $A_t(i,j)$. For the reward-tilted marginals $q_t(i) \propto p_t(i) \exp(\beta_t r(i))$, the following equation holds

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{reward}}(j, i) q_t(j) - A_t^{\text{reward}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{q_t(j)} g_t(j) \right), \quad (17)$$

$$A_t^{\text{reward}}(i,j) := A_t(i,j) \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(i))}, \ g_t(i) := \sum_{j \neq i} \left(A_t^{\text{reward}}(i,j) - A_t(i,j) \right) + \frac{\partial \beta_t}{\partial t} r(i).$$
 (18)

Note that the obtained formulas depend only on the reward function and the rate matrix of the original process. Applying this result to the masked diffusion we obtain the following corollary.

Corollary 3.6. [Reward-tilted Masked Diffusion] For the rate matrix of the reverse-time masked diffusion from Eq. (10), Theorem 3.5 yields

$$B_{\tau}^{\text{reward}}(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)} \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(m))},$$
(19)

$$g_{\tau}(i) = \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \delta_{mi} \sum_{j} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t(j)}{p_t(m)} \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(m))} \right) + \frac{\partial \beta_t}{\partial t} r(i).$$
 (20)

Note that evaluating $B_{\tau}^{\mathrm{reward}}(i,j)$ requires computing the reward function at all the states j we can transition to from mask m. Furthermore, computing $g_t(i)$ requires the summation of the reward over all such states j, which, depending on the application, might be computationally expensive. To avoid these extra computations one could potentially use alternative functions evaluating the difference in the rewards on the transitions from m to j, i.e. r(j) - r(m). However, we leave this as a future work.

4 Experiments

In this section, we demonstrate the utility of the proposed DISCRETE FEYNMAN-KAC CORRECTORS on several applications using modern discrete diffusion models. Each experiment is aimed at illustrating one of the introduced processes: annealing, geometric averaging, reward-tilting.

Despite different domains and processes, the generation process always follows the same procedure described in Alg. 1. Namely, for the corresponding rate matrix $B_{\tau}(i,j)$ and weight function $g_{\tau}(i)$ (see Section 3 for their definitions), the inference procedure generates a batch of samples x_{τ}^k together with their weights w_{τ}^k . In practice, we always perform resampling in between the update steps using SNIS. Thus, DFKC not only changes the generation of individual samples by changing the rate matrix $B_{\tau}(i,j)$ but also introduces "interactions" between samples through re-weighting and re-sampling.

Algorithm 1: Generation using DISCRETE FEYNMAN-KAC CORRECTORS

Output: weighted set of samples $\{(x_{\tau=1}^k, w_{\tau=1}^k)\}_{k=1}^K$

Target β	Method	Energy- $W_2(\downarrow)$	Magnetization- $\mathcal{W}_2(\downarrow)$	Correlation-MSE (\downarrow)
0.4	DFKC(0.3)	14.24 ± 3.11	0.256 ± 0.052	0.041 ± 0.013
	DDM	69.38 ± 4.25	0.889 ± 0.063	0.172 ± 0.021
0.3	DFKC(0.2)	33.38 ± 0.46	0.031 ± 0.011	0.023 ± 0.007
	DDM	35.14 ± 0.63	0.046 ± 0.012	0.014 ± 0.009

Table 1: Sampling task for Ising model with performance measured by mean ±standard deviation over 3 seeds. The starting temperature for DFKC is shown in brackets. The DDM samples are generated with a discrete diffusion model trained at those corresponding target temperatures.

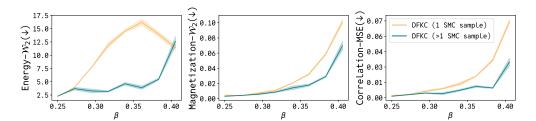


Figure 2: 2-Wasserstein metric for energy and magnetization distributions and MSE for spin-spin correlation. All metrics are computed between samples from DFKC variants and samples from Swendsen-Wang algorithm. Training β is 0.25.

4.1 Annealing the Ising Model

We apply Theorem 3.1 for annealing the Boltzmann distribution of the Ising model configurations. Namely, the probability distribution of states σ is given as

$$p_{\beta}(\sigma) = \frac{1}{Z_{\beta}} e^{-\beta H(\sigma)}, \ Z_{\beta} = \sum_{\sigma} e^{-\beta H(\sigma)}, \text{ where } H(\sigma) = -\sum_{i,j} J_{ij} \sigma_i \sigma_j - \sum_i h_i \sigma_i.$$
 (21)

We generate the training dataset at a fixed β by running the Swendsen-Wang algorithm (Swendsen & Wang, 1987) and train a discrete masked-diffusion model. We set $J_{ij}=1$ and $h_i=0$ on a 16×16 lattice with open boundary conditions. The diffusion model is implemented using the UNet architecture. We assess method performance by comparing the distributions of key observables, specifically energy and magnetization. To examine the fidelity of local structures, we compute spin–spin correlations as a function of distance, excluding boundary spins and evaluating correlations along lattice rows. Finally, we evaluate the mean squared error (MSE) between the generated correlation profiles and the ground-truth.

In Fig. 2, we train the diffusion model at $\beta=0.25$ and we demonstrate that DFKC allows for the efficient control of temperature at inference time. As a baseline, we consider a guidance method, which ignores the weights of the generated samples.

In Table 1, we demonstrate that collecting the data at a high temperature and annealing the trained model to the low temperature is more efficient than collecting data and training the model directly at a low temperature. In particular, we fix the number of energy evaluations for the dataset collection and can either allocate this budget at training DDM directly on the target temperature, or at training it a higher temperature and then use DFKC to reduce the temperature to the target. Additional details of the experiments are included in Appendix C.4. To conduct this comparison, we used 10,000 samples following a long burn-in period of Glauber dynamics, which requires lengthy chains to reduce correlations.

4.2 Extending Language Model Context with Products

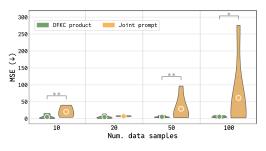
We evaluate the product formula for DFKC from Theorem 3.3 on text generation tasks. We consider the problem of generation under a prompt C which consists of multiple individual conditions $C = \bigcup_{k=1}^K C_i$. Importantly, a large number of conditions K leads to a more complex generation task. Additionally, language models have been shown to suffer degradation in certain tasks when given large prompts (Hsieh et al., 2024; Li et al., 2024). We tackle these issues by applying our framework

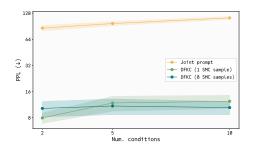
to diffusion language models p_{ϕ} , where we replace the joint prompt: $p_{\phi}(x \mid C)$, with the product $\prod_{k=1}^{K} p_{\phi}(x \mid C_k)$. We examine two tasks of this kind, both using the LLaDA model (Nie et al., 2025).

Amortized Learning Given a dataset of examples $\mathcal{X} = \{(x_i,y_i)\}_{i=1}^N$, and a parametric model $f_{\theta}(x)$, we wish to use the language model to infer parameters θ which fit the data. This requires sampling from the posterior distribution over parameters $p(\theta|\mathcal{X})$. However, unlike more classical statistical methods, we wish to perform this computation solely through the text interface of the language model, similar to the setting of (Requeima et al., 2024; Mittal et al., 2025). Namely we set \mathcal{X} as our prompt, and ask the model to sample parameters θ . We partition the dataset into K equal subsets $\mathcal{X} = \bigcup_{k=1}^K \mathcal{X}_k$, and note that for a uniform prior, the posterior factors as $p(\theta|\mathcal{X}) \propto \prod_{k=1}^K p(\theta|\mathcal{X}_k)$. This justifies applying our method, with each factor in the product conditioned on a different subset of the data $C_k = \mathcal{X}_k$. We evaluate this task on a synthetic dataset generated using a noisy linear predictor $f_{\theta}(x) = \theta_1 x + \theta_0 + \epsilon$ $\epsilon \sim \mathcal{N}(0, 0.1^2)$. We use K = 5 subsets, and report our results for the mean-squared error (to the true parameters) across larger datasets \mathcal{X} in Fig. 3a. Additional details can be found in Appendix C.3.

Multi-constraint Story Generation For this task we prompt the language model to generate a story, with a list of constraints $C = \bigcup_k C_k$. Constraints may demand the inclusion of particular events or characters (such as a "hungry cat"), or be stylistic in nature ("the story should have mystery"). We use our method to sample from the product over individual constraints, and evaluate our adherence to the constraints by using the perplexity of the output under a more powerful language model, Qwen2.5 (Yang et al., 2024). Results for our method, over a varying number of constraints K, are included in Fig. 3b. Additional details are in Appendix C.2.

From our results for both tasks, we can see that as the length and complexity of the prompt increases, the joint prompt degrades in performance, compared to the more stable performance of the DFKC product. We also see from Fig. 3b that using more samples in our method improves performance slightly over 1 sample. This trend is also be seen by ablating over the number of SMC samples for the amortized learning task Fig. A1.





(a) Amortized learning task: Mean squared error (MSE) between predicted and true parameters reported for DFKC (1 and 5 samples), and joint prompting, across different dataset sizes. ** indicates $p \le 0.02$, * indicates $p \le 0.05$ (one-sided Student's *t*-test).

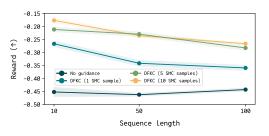
(b) Multi-constraint story generation task: Comparison of Perplexity (PPL), between joint prompting, DFKC (1 SMC sample), and DFKC (8 SMC samples), for different numbers of conditions.

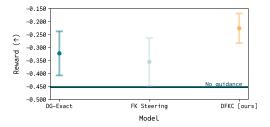
Figure 3: DFKC product performance for text generation tasks. All results averaged over 5 seeds.

4.3 Guiding Protein Sequence Generation with External Rewards

Finally, we investigate the utility of DFKC in the setting of unconditional *de novo* protein sequence generation. Protein language models (PLMs) have emerged as powerful tools for modeling the complex relationships between protein sequence, structure, and function (Lin et al., 2023; Madani et al., 2023), but their controllability remains a challenge. We address the challenge of ensuring that generated sequences resemble natural proteins by guiding generation with a likelihood-based reward. Because PLMs capture inter-residue dependencies and evolutionary conservation, they assign higher likelihoods to "natural-like" sequences, making likelihood a useful proxy for viability. However, computing the likelihood of a sample generated from a discrete diffusion model is not exact and potentially nontrivial (Nie et al., 2025). Instead, we use a masked language model, ESM2-650M (Lin et al., 2023), a PLM whose likelihoods have been reliably used to optimize sequences toward more functional and biologically viable proteins (Ertelt et al., 2024; Emami et al., 2023; Notin et al., 2023).

Here, we generate sequences using DPLM-650M, a discrete diffusion model that produces protein sequences by progressively unmasking amino acid tokens (Wang et al., 2024a). To guide generation, we compute the mean log-likelihood of each intermediate sequence under ESM2-650M and apply this reward through Theorem 3.6. Figure 4a presents the reward values of final sequences with and without guidance across different sequence lengths. In the guided setting, we explore DFKC with 1, 5, and 10 SMC samples. Our single-sample variant is equivalent to the approach of Nisonoff et al. (2024), while using multiple samples yields notable improvements in mean reward compared to both unguided DPLM sampling and guidance without resampling. These results highlight the effectiveness of our resampling procedure in enhancing the biological plausibility of generated sequences. We also compare our method with other guidance-based methods (FK Steering (Singhal et al., 2025) and DG-Exact (Nisonoff et al., 2024)) in Fig. 4b and find our method is able to generate higher-reward sequences.





- (a) Rewards (ESM2-650M log-likelihood) of generated sequences for 10, 50, 100 amino acids at 1, 5, 10 SMC samples and base model (no guidance).
- (b) Comparison of ESM2-650M log-likelihood rewards of best DFKC model with FK Steering (Singhal et al., 2025) and DG-Exact (Nisonoff et al., 2024).

Figure 4: DFKC performance on reward-guided unconditional protein sequence generation.

5 Related Work

Reward Fine-tuning These methods often assume an external reward function r(x) and adjust the pretrained model's parameters using reinforcement learning algorithms, with the goal of sampling from the product $r(x)q_t(x)$. Several of these works are applicable to discrete diffusion models (Venkatraman et al., 2024; Rector-Brooks et al., 2024; Wang et al., 2025). Our method leaves the pretrained model fixed, and therefore doesn't require a costly fine-tuning stage.

Inference Time Alignment Several methods perform additional computation at inference time to sample from a target product distribution (the product being taken with either an external model r(x), or a classifier extracted from the model's distribution, $q_t(y|x)$ as in classifier-free guidance (Ho & Salimans, 2022)). These methods often involve an approximation which means they produce biased samples from the target product (Vignac et al., 2022; Gruver et al., 2023; Nisonoff et al., 2024; Tang et al., 2025). Singhal et al. (2025) investigates the use of SMC to sample (in an asymptotically unbiased manner) from a reward-weighted distribution. Our work adapts such an unbiased SMC based strategy to a smoothly annealed form of the reward $(\beta_t r(x))$, and extends it to general products, and annealing. He et al. (2025) recently proposed another SMC-based technique for such problems, however, they do not evaluate the method on discrete diffusion tasks.

6 Conclusion

In this paper, we propose DISCRETE FEYNMAN-KAC CORRECTORS— a framework that allows for re-purposing discrete diffusion models at inference time without retraining them. In particular, our theoretical findings demonstrate that sampling from the annealed, product or reward-weighted distributions can be efficiently done by combining the learned probability ratios and running SMC algorithms. Our empirical study supports our derivations and demonstrates that the proposed approach is more effective for tasks such as sampling from lower temperature Ising models, generating text based on large composite prompts, and controlling generated protein sequences. This method unlocks possible novel applications of discrete diffusion models in the future such as the collaborative generation of code.

7 Reproducibility Statement

To facilitate reproducibility of our empirical results and algorithm, we have made our code publicly available at this link: https://anonymous.4open.science/r/discrete_fkc-40B8/README.md. We describe all mathematical and algorithmic details necessary to reproduce our results throughout this paper (e.g. Alg. 1).

References

- Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, 630(8016):493–500, 2024.
- Tara Akhound-Sadegh, Jungyoon Lee, Avishek Joey Bose, Valentin De Bortoli, Arnaud Doucet, Michael M. Bronstein, Dominique Beaini, Siamak Ravanbakhsh, Kirill Neklyudov, and Alexander Tong. Progressive inference-time annealing of diffusion models for sampling from boltzmann densities, 2025. URL https://arxiv.org/abs/2506.16471.
- Sarah Alamdari, Nitya Thakkar, Rianne van den Berg, Neil Tenenholtz, Bob Strome, Alan Moses, Alex Xijie Lu, Nicolo Fusi, Ava Pardis Amini, and Kevin K Yang. Protein generation with evolutionary diffusion: sequence is all you need. *BioRxiv*, pp. 2023–09, 2023.
- Joe Benton, Yuyang Shi, Valentin De Bortoli, George Deligiannidis, and Arnaud Doucet. From denoising diffusions to denoising markov models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 86(2):286–301, 2024.
- Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. The reversal curse: Llms trained on" a is b" fail to learn" b is a". *arXiv preprint arXiv:2309.12288*, 2023.
- Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 22563–22575, 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.
- Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural Information Processing Systems*, 35:28266–28279, 2022.
- DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL https://arxiv.org/abs/2501.12948.
- Patrick Emami, Aidan Perreault, Jeffrey Law, David Biagioni, and Peter St John. Plug & play directed evolution of proteins with gradient-based discrete mcmc. *Machine Learning: Science and Technology*, 4(2):025014, 2023.
- Moritz Ertelt, Jens Meiler, and Clara T Schoeder. Combining rosetta sequence design with protein language model predictions using evolutionary scale modeling (esm) as restraint. *ACS synthetic biology*, 13(4):1085–1092, 2024.
- Nate Gruver, Samuel Stanton, Nathan Frey, Tim GJ Rudner, Isidro Hotzel, Julien Lafrance-Vanasse, Arvind Rajpal, Kyunghyun Cho, and Andrew G Wilson. Protein design with guided discrete diffusion. *Advances in neural information processing systems*, 36:12489–12517, 2023.
- Jiajun He, José Miguel Hernández-Lobato, Yuanqi Du, and Francisco Vargas. Rne: a plug-and-play framework for diffusion density estimation and inference-time control. *arXiv preprint arXiv:2506.05668*, 2025.
- Geoffrey E Hinton. Products of experts. In 1999 ninth international conference on artificial neural networks ICANN 99.(Conf. Publ. No. 470), volume 1, pp. 1–6. IET, 1999.
- Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. Advances in Neural Information Processing Systems (NeurIPS) Workshop on Deep Generative Models and Downstream Applications, 2022. URL https://arxiv.org/abs/2207.12598.

- Cheng-Ping Hsieh, Simeng Sun, Samuel Kriman, Shantanu Acharya, Dima Rekesh, Fei Jia, and Boris Ginsburg. RULER: What's the real context size of your long-context language models? *Conference on Language Modeling (COLM)*, 2024. URL https://openreview.net/forum?id=kioBbc76Sy.
- Andrei Kolmogoroff. Über die analytischen methoden in der wahrscheinlichkeitsrechnung. *Mathematische Annalen*, 104:415–458, 1931.
- Seul Lee, Karsten Kreis, Srimukh Prasad Veccham, Meng Liu, Danny Reidenbach, Yuxing Peng, Saee Paliwal, Weili Nie, and Arash Vahdat. Genmol: A drug discovery generalist with discrete diffusion. *arXiv preprint arXiv:2501.06158*, 2025.
- Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with long in-context learning. *Transactions on Machine Learning Research (TMLR)*, 2024. URL https://arxiv.org/abs/2404.02060.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. doi: 10.1126/science.ade2574. URL https://www.science.org/doi/abs/10.1126/science.ade2574.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion language modeling by estimating the ratios of the data distribution. *International Conference on Machine Learning (ICML)*, 2024.
- Ali Madani, Ben Krause, Eric R. Greene, Subu Subramanian, Benjamin P. Mohr, James M. Holton, Jose Luis Olmos, Caiming Xiong, Zachary Z. Sun, Richard Socher, James S. Fraser, and Nikhil Naik. Large language models generate functional protein sequences across diverse families. *Nature Biotechnology*, 41(8):1099–1106, Aug 2023. ISSN 1546-1696. doi: 10.1038/s41587-022-01618-2. URL https://doi.org/10.1038/s41587-022-01618-2.
- Sarthak Mittal, Niels Leif Bracher, Guillaume Lajoie, Priyank Jaini, and Marcus Brubaker. Amortized in-context bayesian posterior estimation. *arXiv preprint arXiv:2502.06601*, 2025. URL https://arxiv.org/abs/2502.06601.
- Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*, 2025. URL https://arxiv.org/abs/2502.09992.
- Hunter Nisonoff, Junhao Xiong, Stephan Allenspach, and Jennifer Listgarten. Unlocking guidance for discrete state-space diffusion and flow models. *arXiv preprint arXiv:2406.01572*, 2024.
- Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood van Niekerk, Steffanie Paul, Han Spinner, Nathan Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, Jonathan Frazer, Mafalda Dias, Dinko Franceschi, Yarin Gal, and Debora Marks. Proteingym: Large-scale benchmarks for protein fitness prediction and design. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 64331–64379. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/cac723e5ff29f65e3fcbb0739ae91bee-Paper-Datasets_and_Benchmarks.pdf.
- Jarrid Rector-Brooks, Mohsin Hasan, Zhangzhi Peng, Zachary Quinn, Chenghao Liu, Sarthak Mittal, Nouha Dziri, Michael Bronstein, Yoshua Bengio, Pranam Chatterjee, et al. Steering masked discrete diffusion models via discrete denoising posterior prediction. *arXiv preprint arXiv:2410.08134*, 2024.
- James Requeima, John F Bronskill, Dami Choi, Richard E. Turner, and David Duvenaud. LLM processes: Numerical predictive distributions conditioned on natural language. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=HShs7q1Njh.

- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 10684–10695, June 2022.
 - Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022.
 - 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.
 - Raghav Singhal, Zachary Horvitz, Ryan Teehan, Mengye Ren, Zhou Yu, Kathleen McKeown, and Rajesh Ranganath. A general framework for inference-time scaling and steering of diffusion models. *arXiv preprint arXiv:2501.06848*, 2025.
 - Marta Skreta, Lazar Atanackovic, Joey Bose, Alexander Tong, and Kirill Neklyudov. The superposition of diffusion models using the itô density estimator. In *The Thirteenth International Conference on Learning Representations*, 2024.
 - Marta Skreta, Tara Akhound-Sadegh, Viktor Ohanesian, Roberto Bondesan, Alán Aspuru-Guzik, Arnaud Doucet, Rob Brekelmans, Alexander Tong, and Kirill Neklyudov. Feynman-kac correctors in diffusion: Annealing, guidance, and product of experts. *arXiv preprint arXiv:2503.02819*, 2025.
 - Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations, 2021. URL https://arxiv.org/abs/2011.13456.
 - Robert H Swendsen and Jian-Sheng Wang. Nonuniversal critical dynamics in monte carlo simulations. *Physical review letters*, 58(2):86, 1987.
 - Sophia Tang, Yinuo Zhang, and Pranam Chatterjee. Peptune: De novo generation of therapeutic peptides with multi-objective-guided discrete diffusion. *ArXiv*, pp. arXiv–2412, 2025.
 - Siddarth Venkatraman, Moksh Jain, Luca Scimeca, Minsu Kim, Marcin Sendera, Mohsin Hasan, Luke Rowe, Sarthak Mittal, Pablo Lemos, Emmanuel Bengio, Alexandre Adam, Jarrid Rector-Brooks, Yoshua Bengio, Glen Berseth, and Nikolay Malkin. Amortizing intractable inference in diffusion models for vision, language, and control. *Advances in Neural Information Processing Systems* (NeurIPS), 2024. URL https://openreview.net/forum?id=gVTkMsaaGI.
 - 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.
 - Pierre-A Vuillermot. A generalization of chernoff's product formula for time-dependent operators. *Journal of Functional Analysis*, 259(11):2923–2938, 2010.
 - Chenyu Wang, Masatoshi Uehara, Yichun He, Amy Wang, Avantika Lal, Tommi Jaakkola, Sergey Levine, Aviv Regev, Hanchen, and Tommaso Biancalani. Fine-tuning discrete diffusion models via reward optimization with applications to DNA and protein design. *International Conference on Learning Representations (ICLR)*, 2025. URL https://openreview.net/forum?id=G328D1xt4W.
 - Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
 - Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Diffusion language models are versatile protein learners, 2024a. URL https://arxiv.org/abs/2402.18567.

Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. Diffusion language models are versatile protein learners. In *International Conference on Machine Learning*, 2024b.

Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023.

An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.

Background Proofs A

Weighted Forward Kolmogorov Equation

Consider the forward Kolmogorov equation with the weighting term

$$\frac{\partial p_s(j)}{\partial s} = \sum_{k \neq j} A_s(k,j) p_s(k) - \sum_{k \neq j} A_s(j,k) p_s(j) + p_s(j) (g_s(j) - \sum_k p_s(k) g_s(k)). \tag{22}$$

We can re-write the last term as

$$p_s(j)(g_s(j) - \sum_k p_s(k)g_s(k)) = \sum_k p_s(k)p_s(j)(g_s(j) - g_s(k))$$
(23)

$$= \sum_{k} p_s(k) p_s(j) \sigma_s(j, k) |g_s(j) - g_s(k)|$$
 (24)

$$= \sum_{k} p_s(j) \mathbb{1}[\sigma_s(j,k) > 0] |g_s(j) - g_s(k)| p_s(k) -$$
 (25)

$$-\sum_{k} p_s(k) \mathbb{1}[\sigma_s(j,k) < 0] |g_s(j) - g_s(k)| p_s(j), \qquad (26)$$

where $\sigma_s(j,k)$ is the sign of $(g_s(j) - g_s(k))$. Let's define

$$B_s(k,j) := p_s(k) \mathbb{1}[\sigma_s(j,k) > 0]|g_s(j) - g_s(k)|$$
 (27)

$$\implies B_s(j,k) := p_s(j)\mathbb{1}[\sigma_s(k,j) > 0]|g_s(k) - g_s(j)|. \tag{28}$$

Using the fact that $\sigma_s(k, j) = -\sigma_s(j, k)$, we have

$$p_s(j)(g_s(j) - \sum_k p_s(k)g_s(k)) = \sum_k B_s(k,j)p_s(k) - \sum_k B_s(j,k)p_s(j).$$
 (29)

Finally, using the fact that $B_s(j, j) = 0$, we have

$$\frac{\partial p_s(j)}{\partial s} = \sum_{k \neq j} A_s(k, j) p_s(k) - \sum_{k \neq j} A_s(j, k) p_s(j) + p_s(j) (g_s(j) - \sum_k p_s(k) g_s(k))$$
(30)

$$= \sum_{k \neq j} (A_s(k,j) + B_s(k,j)) p_s(k) - \sum_{k \neq j} (A_s(j,k) + B_s(j,k)) p_s(j),$$
(31)

$$B_s(k,j) := p_s(k) \mathbb{1}[\sigma_s(j,k) > 0]|g_s(j) - g_s(k)|. \tag{32}$$

Discrete Feynman-Kac formula

Solution to the system of linear homogeneous differential equations can be written using time-ordered matrix exponent. For the equation

$$\partial_t p_t(i) = \sum_j (A_t(i,j) + V_t(i,j)) p_t(j)$$
$$\partial_t p_t = (A_t + V_t) p_t$$

Solution is

$$p_t = \mathcal{T}\left\{\exp\left(\int_0^t (A_s + V_s)ds\right)\right\} p_0 \tag{33}$$

Alternatively, this solution can be written as a limit:

$$p_t = \lim_{n \to \infty} \prod_{k=1}^n \exp((A_{k\tau} + V_{k\tau})\tau) p_0$$
 (34)

where $\tau = t/n$.

For time-independent operator which is sum of two operators A and V Lie-Trotter formula can be

applied:

$$\exp(t[A+V]) = \lim_{n \to \infty} \prod_{k=1}^{n} [\exp(A\tau) \exp(V\tau)]$$
 (35)

It is shown in (Vuillermot, 2010) that these limits can be united in one creating analogue for timedependent Lie-Trotter formula:

$$\mathcal{T}\exp(\int_0^t [A_s + V_s]ds) = \lim_{n \to \infty} \prod_{k=1}^n [\exp(A_{k\tau}\tau) \exp(V_{k\tau}\tau)]$$
 (36)

let's call $\exp(A_{k\tau}\tau)_{ij}=p_{k\tau+\tau|k\tau}(i,j)=P_k(i,j)$ and $\exp(V_{k\tau}\tau)_{ij}=W_k(i,j)=W_k(i)I(i,j)$,

$$\begin{split} p_t(j) &= \lim_{n \to \infty} \sum_{i_n} \sum_{j_n} ... \sum_{i_1} \sum_{j_1} P_n(j,i_n) W_n(i_n,j_n) ... P_1(j_2,i_1) W_1(i_1,j_1) p_0(j_1) = \\ &= \underbrace{\lim_{n \to \infty} \sum_{j_n} ... \sum_{j_1} P_n(j,j_n) ... P_1(j_2,j_1) p_0(j_1)}_{\text{expectation over all paths ending in j produced by A} W_n(j_n) ... W_1(j_1) = \underbrace{\underbrace{\lim_{n \to \infty} \sum_{j_n} ... \sum_{j_1} P_n(j,j_n) ... P_1(j_2,j_1) p_0(j_1)}_{\text{expectation over all paths ending in j produced by A} \end{split}$$

$$= \mathbb{E}^{A} \left[\exp \left(\int_{0}^{t} V(j_{s}) ds \right) \mid j_{t} = j \right]$$

A.3 Discrete Masked Diffusion

First, we consider general case, where m is the mask state and $\alpha_{s,t}$ is the noise schedule, i.e. the noising process is defined as

$$p(x_s = j \mid x_t = i) = (1 - \bar{\alpha}_{s,t})\delta_{mi} + \bar{\alpha}_{s,t}\delta_{ij}.$$
(37)

Note that not every $\bar{\alpha}_{s,t}$ satisfies the master equation and we have to ensure that the following equality holds.

$$p(x_s = j \mid x_t = i) = \sum_k p(x_s = j \mid x_r = k) p(x_r = k \mid x_t = i)$$
(38)

$$(1 - \bar{\alpha}_{s,t})\delta_{mj} + \bar{\alpha}_{s,t}\delta_{ij} = \sum_{k} ((1 - \bar{\alpha}_{s,r})\delta_{mj} + \bar{\alpha}_{s,r}\delta_{kj})((1 - \bar{\alpha}_{r,t})\delta_{mk} + \bar{\alpha}_{r,t}\delta_{ik})$$
(39)

$$(1 - \bar{\alpha}_{s,t})\delta_{mj} + \bar{\alpha}_{s,t}\delta_{ij} = (1 - \bar{\alpha}_{s,r})\delta_{mj}(\bar{\alpha}_{r,t} + (1 - \bar{\alpha}_{r,t})) + \bar{\alpha}_{s,r}((1 - \bar{\alpha}_{r,t})\delta_{mj} + \bar{\alpha}_{r,t}\delta_{ij}) (1 - \bar{\alpha}_{s,t})\delta_{mj} + \bar{\alpha}_{s,t}\delta_{ij} = ((1 - \bar{\alpha}_{s,r}) + \bar{\alpha}_{s,r}(1 - \bar{\alpha}_{r,t}))\delta_{mj} + \bar{\alpha}_{s,r}\bar{\alpha}_{r,t}\delta_{ij}.$$
(40)

Thus, the following relations must hold

$$1 - \bar{\alpha}_{s,t} = (1 - \bar{\alpha}_{s,r}) + \bar{\alpha}_{s,r}(1 - \bar{\alpha}_{r,t}), \quad \bar{\alpha}_{s,t} = \bar{\alpha}_{s,r}\bar{\alpha}_{r,t}$$
(41)

$$-\bar{\alpha}_{s,t} = -\bar{\alpha}_{r,t}\bar{\alpha}_{s,r}, \quad \bar{\alpha}_{s,t} = \bar{\alpha}_{s,r}\bar{\alpha}_{r,t}, \tag{42}$$

$$\bar{\alpha}_{s,t} = \bar{\alpha}_{r,t}\bar{\alpha}_{s,r} \,. \tag{43}$$

Thus, any function that satisfy the following equation works

$$\forall t < r < s, \ \bar{\alpha}_{s,t} = \bar{\alpha}_{s,r} \bar{\alpha}_{r,t}. \tag{44}$$

Denoting $\alpha_s = \bar{\alpha}_{s,0}$, we have

$$\bar{\alpha}_{s,t} = \frac{\alpha_s}{\alpha_t}$$
, and $p(x_s = j \mid x_t = i) = \left(1 - \frac{\alpha_s}{\alpha_t}\right) \delta_{mj} + \frac{\alpha_s}{\alpha_t} \delta_{ij}$. (45)

From here, the rate matrix of the noising process is

$$A_t(i,j) = \frac{\partial p(x_s = j \mid x_t = i)}{\partial s} \bigg|_{s=t} = \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} (\delta_{ij} - \delta_{mj}). \tag{46}$$

A.4 Reverse-time Masked Diffusion

 For the inverse time $\tau=1-t$, we flip the marginals $q_{\tau}(i)\coloneqq p_{1-\tau}(i)$ and take the derivative w.r.t. τ

$$\frac{\partial q_{\tau}(i)}{\partial \tau} = \frac{\partial p_{1-\tau}(i)}{\partial \tau} = -\frac{\partial p_t(i)}{\partial t} \bigg|_{t=1-\tau}$$
(47)

$$= -\sum_{j \neq i} (A_{1-\tau}(j,i)p_{1-\tau}(j) - A_{1-\tau}(i,j)p_{1-\tau}(i))$$
(48)

$$= \sum_{i \neq i} \left(A_{1-\tau}(i,j) \frac{p_{1-\tau}(i)}{q_{\tau}(j)} q_{\tau}(j) - A_{1-\tau}(j,i) \frac{p_{1-\tau}(j)}{q_{\tau}(i)} q_{\tau}(i) \right)$$
(49)

$$= \sum_{j \neq i} (B_{\tau}(j, i) q_{\tau}(j) - B_{\tau}(i, j) q_{\tau}(i)), \quad B_{\tau}(i, j) \coloneqq A_{1-\tau}(j, i) \frac{p_{1-\tau}(j)}{p_{1-\tau}(i)}. \tag{50}$$

Note that here we define only the off-diagonal elements and the diagonal elements are

$$B_{\tau}(i,i) = -\sum_{j \neq i} B_{\tau}(i,j) = -\sum_{j \neq i} A_{1-\tau}(j,i) \frac{p_{1-\tau}(j)}{p_{1-\tau}(i)}.$$
 (51)

In particular, for the masked diffusion, we have

$$B_{\tau}(i,j) = \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} (\delta_{ij} - \delta_{mi}) \frac{p_t(j)}{p_t(i)}, \quad i \neq j$$
(52)

$$= -\frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)} \delta_{mi}, \qquad (53)$$

$$B_{\tau}(i,i) = -\sum_{i \neq i} B_{\tau}(i,j) = \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{1 - p_t(m)}{p_t(m)} \delta_{mi}.$$
 (54)

A.5 De-masking parameterization

Furthermore, analogously to the derivation from (Shi et al., 2024) (Appendix H.3), we have

$$\frac{p_t(j)}{p_t(m)} = \sum_{i} \frac{p_0(i)}{p_t(m)} p(x_t = j \mid x_0 = i)$$
(55)

$$= \sum_{i} \frac{p_0(i)p(x_t = m \mid x_0 = i)}{p_t(m)p(x_0 = i \mid x_t = m)} \frac{p(x_0 = i \mid x_t = m)}{p(x_t = m \mid x_0 = i)} p(x_t = j \mid x_0 = i)$$
 (56)

$$= \sum_{i} \frac{p(x_0 = i \mid x_t = m)}{p(x_t = m \mid x_0 = i)} p(x_t = j \mid x_0 = i)$$
(57)

$$= \sum_{i} \frac{p(x_0 = i \mid x_t = m)}{(1 - \alpha_t) + \alpha_t \delta_{im}} ((1 - \alpha_t) \delta_{mj} + \alpha_t \delta_{ij})$$

$$(58)$$

$$= \frac{1}{1 - \alpha_t} \sum_{i} ((1 - \alpha_t) \delta_{mj} + \alpha_t \delta_{ij}) p(x_0 = i \mid x_t = m)$$
 (59)

$$= \delta_{mj} + \frac{\alpha_t}{1 - \alpha_t} p(x_0 = j \mid x_t = m).$$
 (60)

where we used the fact that $p(x_0 = m) = 0$.

A.6 Multidimensional case

For the multi-dimensional case, we consider the masking process applied independently to each coordinate, i.e.

$$p(x_s = [j_1 \dots j_d] \mid x_t = [i_1 \dots i_d]) = \prod_{k=1}^d p(x_s[k] = j_k \mid x_t[k] = i_k)$$
(61)

$$= \prod_{k=1}^{d} \left(\left(1 - \frac{\alpha_s}{\alpha_t} \right) \delta_{mj_k} + \frac{\alpha_s}{\alpha_t} \delta_{i_k j_k} \right), \tag{62}$$

which defines the following rate matrix

$$A_t([i_1 \dots i_d], [j_1 \dots j_d]) = \frac{\partial p(x_s = [j_1 \dots j_d] \mid x_t = [i_1 \dots i_d])}{\partial s} \Big|_{s=t}$$
(63)

$$= \sum_{k=1}^{d} \prod_{l \neq k} p(x_t[l] = j_l \,|\, x_t[l] = i_l) \frac{\partial p(x_s[k] = j_k \,|\, x_t[k] = i_k)}{\partial s} \bigg|_{s=t}$$
(64)

$$= \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{k=1}^d \prod_{l \neq k} \delta_{j_l i_l} (\delta_{i_k j_k} - \delta_{m j_k}). \tag{65}$$

For the off-diagonal elements of the reverse-time matrix, we have

$$B_t([i_1 \dots i_d], [j_1 \dots j_d]) = A_t([j_1 \dots j_d], [i_1 \dots i_d]) \frac{p_t([j_1 \dots j_d])}{p_t([i_1 \dots i_d])}$$
(66)

$$= \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t([j_1 \dots j_d])}{p_t([i_1 \dots i_d])} \sum_{k=1}^d \prod_{l \neq k} \delta_{j_l i_l} (\delta_{i_k j_k} - \delta_{m i_k})$$
 (67)

$$= -\frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t([j_1 \dots j_d])}{p_t([i_1 \dots i_d])} \sum_{k=1}^d \prod_{l \neq k} \delta_{j_l i_l} \delta_{m i_k}.$$
 (68)

B DISCRETE FEYNMAN-KAC CORRECTORS Proofs

B.1 Annealing of FKE

Theorem 3.1. [Temperature Annealing] Consider the forward Kolmogorov equation from Eq. (2) describing the time-evolution of the marginals $p_t(i)$ with the rate matrix $A_t(i,j)$. For the temperature annealed marginals $q_t(i) \propto p_t(i)^{\beta}$, the following equation holds

$$\frac{\partial q_t(i)}{\partial t} = \sum_{i \neq i} \left(A_t^{\text{anneal}}(j, i) q_t(j) - A_t^{\text{anneal}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{q_t(j)} g_t(j) \right), \tag{12}$$

where
$$A_t^{\text{anneal}}(i,j) := \beta A_t(i,j) \frac{p_t^{1-\beta}(i)}{p_t^{1-\beta}(j)}, \ g_t(i) := \sum_{i \neq i} \left(A_t^{\text{anneal}}(i,j) - \beta A_t(i,j) \right).$$
 (13)

Proof. Consider the forward Kolmogorov equation for the given rate matrix $A_t(i,j)$

$$\frac{\partial p_t(i)}{\partial t} = \sum_{j \neq i} A_t(j, i) p_t(j) - \sum_{j \neq i} A_t(i, j) p_t(i)$$
(69)

$$\frac{\partial}{\partial t} \log p_t(i) = \sum_{j \neq i} A_t(j, i) \frac{p_t(j)}{p_t(i)} - \sum_{j \neq i} A_t(i, j) = \sum_{j \neq i} \left(A_t(j, i) \frac{p_t(j)}{p_t(i)} - A_t(i, j) \right). \tag{70}$$

Then the annealed target $q_t(i) := p_t^{\beta}(i)/Z_t$ follows

$$\frac{\partial}{\partial t} \log q_t(j) = \beta \frac{\partial}{\partial t} \log p_t(i) - \frac{\partial}{\partial t} \log Z_t \tag{71}$$

$$= \sum_{j \neq i} \left(\beta A_t(j, i) \frac{p_t(j)}{p_t(i)} - \beta A_t(i, j) \right) - \frac{\partial}{\partial t} \log Z_t$$
 (72)

$$= \sum_{j \neq i} \left(\underbrace{\beta A_t(j, i) \frac{p_t^{1-\beta}(j)}{p_t^{1-\beta}(i)}}_{q_t(i)} \frac{q_t(j)}{q_t(i)} - A_t^{\text{anneal}}(i, j) \right) + \tag{73}$$

$$+ \sum_{j \neq i} \left(A_t^{\text{anneal}}(i,j) - \beta A_t(i,j) \right) - \frac{\partial}{\partial t} \log Z_t \,. \tag{74}$$

Denoting the second term as $g_t(j)$, we have

$$\frac{\partial q_t(i)}{\partial t} = \sum_{i \neq i} \left(A_t^{\text{anneal}}(j, i) q_t(j) - A_t^{\text{anneal}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \frac{\partial}{\partial t} \log Z_t \right), \quad (75)$$

$$A_t^{\text{anneal}}(j,i) := \beta A_t(j,i) \frac{p_t^{1-\beta}(j)}{p_t^{1-\beta}(i)}, \quad g_t(i) := \sum_{j \neq i} \left(A_t^{\text{anneal}}(i,j) - \beta A_t(i,j) \right). \tag{76}$$

From the definition of $q_t(i)$ we have

$$\sum_{i} q_t(i) = 1, \quad \forall t, \tag{77}$$

hence,

$$\sum_{i} \frac{\partial q_t(i)}{\partial t} = 0 \implies \sum_{i} q_t(i) \left(g_t(i) - \frac{\partial}{\partial t} \log Z_t \right) = 0, \tag{78}$$

which immediately yields

$$g_t(i) - \frac{\partial}{\partial t} \log Z_t = g_t(i) - \mathbb{E}_{i \sim q_t(i)} g_t(i).$$
(79)

However, one can also verify this through the definition of the normalization constant

$$\frac{\partial}{\partial t} \log Z_t = \frac{1}{Z_t} \sum_i \frac{\partial p_t^{\beta}(i)}{\partial t} = \sum_i \frac{p_t^{\beta}(i)}{Z_t} \beta \frac{\partial}{\partial t} \log p_t(i)$$
 (80)

$$= \sum_{i} q_t(i) \sum_{j \neq i} \left(\beta A_t(j, i) \frac{p_t(j)}{p_t(i)} - \beta A_t(i, j) \right), \tag{81}$$

and, correspondingly

$$\sum_{i} q_t(i)g_t(i) - \frac{\partial}{\partial t} \log Z_t = \sum_{i} q_t(i) \sum_{j \neq i} \left(\beta A_t(i,j) \frac{p_t^{1-\beta}(i)}{p_t^{1-\beta}(j)} - \beta A_t(j,i) \frac{p_t(j)}{p_t(i)} \right)$$
(82)

$$= \frac{\beta}{Z_t} \sum_{i} \sum_{j \neq i} \left(A_t(i, j) \frac{p_t(i)}{p_t^{1-\beta}(j)} - A_t(j, i) \frac{p_t(j)}{p_t^{1-\beta}(i)} \right)$$
(83)

$$= \frac{\beta}{Z_t} \left(\sum_{i} \sum_{j \neq i} \hat{A}_t(i, j) - \sum_{i} \sum_{j \neq i} \hat{A}_t(j, i) \right)$$
(84)

$$= \frac{\beta}{Z_t} \left(\sum_{i,j} \hat{A}_t(i,j) - \sum_{i,j} \hat{A}_t(j,i) \right) = 0,$$
 (85)

where we denote $\hat{A}_t(i,j) := A_t(i,j) \frac{p_t(i)}{p_t^{1-\beta}(j)}$.

Thus, we have

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{anneal}}(j, i) q_t(j) - A_t^{\text{anneal}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{q_t(j)} g_t(j) \right), \tag{86}$$

$$A_t^{\text{anneal}}(j,i) := \beta A_t(j,i) \frac{p_t^{1-\beta}(j)}{p_t^{1-\beta}(i)}, \quad g_t(i) := \sum_{j \neq i} \left(A_t^{\text{anneal}}(i,j) - \beta A_t(i,j) \right). \tag{87}$$

Corollary B.1. [Annealed Masked Diffusion] For the rate matrix of the reverse-time masked diffusion from Eq. (10), Theorem 3.1 yields the following off-diagonal elements of the rate matrix and the corresponding weight function

$$B_{\tau}^{\text{anneal}}(i,j) = -\delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)}, \ g_{\tau}(i) = \delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{i} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)} \right). \tag{14}$$

Proof. The reverse-time rate matrix is

$$B_t(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)}, \ i \neq j.$$
 (88)

Then, according to Theorem 3.1, the rate matrix of the annealed process is

$$B_t^{\text{anneal}}(i,j) = \beta B_t(i,j) \frac{p_t^{1-\beta}(i)}{p_t^{1-\beta}(j)} = -\delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)} \frac{p_t^{1-\beta}(i)}{p_t^{1-\beta}(j)} = -\delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)}$$
(89)

And the weighting term is

$$g_t(i) = \sum_{j \neq i} \left(B_t^{\text{anneal}}(i, j) - \beta B_t(i, j) \right) = \delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j \neq i} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)} \right)$$
(90)

$$= \delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j \neq m} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)} \right) = \delta_{mi} \frac{\beta}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t^{\beta}(j)}{p_t^{\beta}(m)} \right)$$
(91)

B.2 Product of FKEs

Theorem 3.3. [Product of FKEs] Consider two forward Kolmogorov equations (from Eq. (2)) with different rate matrices $A_t^1(i,j)$ and $A_t^2(i,j)$ describing the evolution of marginals $p_t^1(i)$ and $p_t^2(i)$. For the product of marginals $q_t(i) \propto p_t^1(i)p_t^2(i)$, the following equation holds

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{prod}}(j, i) q_t(j) - A_t^{\text{prod}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{j \sim q_t(j)} g_t(j) \right), \tag{16}$$

$$A_t^{\text{prod}}(i,j) \coloneqq A_t^1(i,j) \frac{p_t^2(j)}{p_t^2(i)} + A_t^2(i,j) \frac{p_t^1(j)}{p_t^1(i)} \,, \ g_t(i) \coloneqq \sum_{j \neq i} \left(A_t^{\text{prod}}(i,j) - A_t^1(i,j) - A_t^2(i,j) \right).$$

Proof. Consider two forward Kolmogorov equations with different rate matrices $A_t^1(i,j)$ and $A_t^2(i,j)$. For both we have the equations of the form

$$\frac{\partial p_t^{1,2}(i)}{\partial t} = \sum_{j \neq i} A_t^{1,2}(j,i) p_t^{1,2}(j) - \sum_{j \neq i} A_t^{1,2}(i,j) p_t^{1,2}(i)$$
(92)

$$\frac{\partial}{\partial t} \log p_t^{1,2}(i) = \sum_{j \neq i} A_t^{1,2}(j,i) \frac{p_t^{1,2}(j)}{p_t^{1,2}(i)} - \sum_{j \neq i} A_t^{1,2}(i,j)$$
(93)

$$= \sum_{j \neq i} \left(A_t^{1,2}(j,i) \frac{p_t^{1,2}(j)}{p_t^{1,2}(i)} - A_t^{1,2}(i,j) \right). \tag{94}$$

Correspondingly, for the density $q_t(i) := p_t^1(i)p_t^2(i)/Z_t$, we have

$$\frac{\partial}{\partial t} \log q_t(i) = \frac{\partial}{\partial t} \log p_t^1(i) + \frac{\partial}{\partial t} \log p_t^2(i) - \frac{\partial}{\partial t} \log Z_t$$
(95)

$$= \sum_{j \neq i} \left(A_t^1(j, i) \frac{p_t^1(j)}{p_t^1(i)} - A_t^1(i, j) + A_t^2(j, i) \frac{p_t^2(j)}{p_t^2(i)} - A_t^2(i, j) \right) - \frac{\partial}{\partial t} \log Z_t$$
 (96)

$$= \sum_{j \neq i} \left(A_t^1(j,i) \frac{p_t^2(i)}{p_t^2(j)} \frac{q_t(j)}{q_t(i)} + A_t^2(j,i) \frac{p_t^1(i)}{p_t^1(j)} \frac{q_t(j)}{q_t(i)} - A_t^1(i,j) - A_t^2(i,j) \right) - \frac{\partial}{\partial t} \log Z_t$$
 (97)

$$= \sum_{j \neq i} \left(\underbrace{\left[A_t^1(j,i) \frac{p_t^2(i)}{p_t^2(j)} + A_t^2(j,i) \frac{p_t^1(i)}{p_t^1(j)} \right]}_{:= A_t^{\text{prod}}(j,i)} \frac{q_t(j)}{q_t(i)} - A_t^1(i,j) - A_t^2(i,j) \right) - \frac{\partial}{\partial t} \log Z_t$$
 (98)

$$= \sum_{j \neq i} \left(A_t^{\text{prod}}(j, i) \frac{q_t(j)}{q_t(i)} - A_t^{\text{prod}}(i, j) \right) + \tag{99}$$

$$+\underbrace{\sum_{j\neq i} \left(A_t^{\text{prod}}(i,j) - A_t^1(i,j) - A_t^2(i,j) \right)}_{:=q_t(i)} - \frac{\partial}{\partial t} \log Z_t.$$

$$(100)$$

Finally, we have to show that the weights are self-normalized, i.e.

$$g_t(i) - \frac{\partial}{\partial t} \log Z_t = g_t(i) - \mathbb{E}_{i \sim q_t(j)} g_t(j).$$
(101)

Expanding the derivative of the normalization constant, we have

$$\frac{\partial}{\partial t} \log Z_{t} = \frac{1}{Z_{t}} \sum_{i} \left(p_{t}^{1}(i) \frac{\partial p_{t}^{2}(i)}{\partial t} + p_{t}^{2}(i) \frac{\partial p_{t}^{1}(i)}{\partial t} \right) = \sum_{i} q_{t}(i) \left(\frac{\partial}{\partial t} \log p_{t}^{2}(i) + \frac{\partial}{\partial t} \log p_{t}^{1}(i) \right) \\
= \sum_{i} q_{t}(i) \sum_{j \neq i} \left(A_{t}^{1}(j, i) \frac{p_{t}^{1}(j)}{p_{t}^{1}(i)} - A_{t}^{1}(i, j) + A_{t}^{2}(j, i) \frac{p_{t}^{2}(j)}{p_{t}^{2}(i)} - A_{t}^{2}(i, j) \right).$$
(102)

Thus, we have

$$\sum_{i} q_{t}(i)g_{t}(i) - \frac{\partial}{\partial t} \log Z_{t} = \sum_{i} q_{t}(i) \sum_{j \neq i} \left(A_{t}^{\text{prod}}(i,j) - A_{t}^{1}(j,i) \frac{p_{t}^{1}(j)}{p_{t}^{1}(i)} - A_{t}^{2}(j,i) \frac{p_{t}^{2}(j)}{p_{t}^{2}(i)} \right) \\
= \sum_{i} q_{t}(i) \sum_{j \neq i} \left(A_{t}^{1}(i,j) \frac{p_{t}^{2}(j)}{p_{t}^{2}(i)} + A_{t}^{2}(i,j) \frac{p_{t}^{1}(j)}{p_{t}^{1}(i)} - A_{t}^{1}(j,i) \frac{p_{t}^{1}(j)}{p_{t}^{1}(i)} - A_{t}^{2}(j,i) \frac{p_{t}^{2}(j)}{p_{t}^{2}(i)} \right) \tag{103}$$

$$= \frac{1}{Z_{t}} \sum_{i} \sum_{j \neq i} \left(A_{t}^{1}(i,j) p_{t}^{1}(i) p_{t}^{2}(j) + A_{t}^{2}(i,j) p_{t}^{1}(j) p_{t}^{2}(i) - \right) \tag{104}$$

$$-A_t^1(j,i)p_t^1(j)p_t^2(i) - A_t^2(j,i)p_t^1(i)p_t^2(j)$$
(105)

Denoting

$$\hat{A}_t(i,j) := A_t^1(i,j)p_t^1(i)p_t^2(j) + A_t^2(i,j)p_t^1(j)p_t^2(i), \qquad (106)$$

we can show

$$\sum_{i} q_t(i)g_t(i) - \frac{\partial}{\partial t} \log Z_t = \frac{1}{Z_t} \sum_{i} \sum_{j \neq i} \left(\hat{A}_t(i,j) - \hat{A}_t(j,i) \right)$$
(107)

$$= \frac{1}{Z_t} \sum_{i,j} \left(\hat{A}_t(i,j) - \hat{A}_t(j,i) \right) = 0.$$
 (108)

Thus, we have the result of the theorem, i.e.

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{prod}}(j, i) \ q_t(j) - A_t^{\text{prod}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{j \sim q_t(j)} g_t(j) \right), \quad (109)$$

where
$$A_t^{\text{prod}}(i,j) := A_t^1(i,j) \frac{p_t^2(j)}{p_t^2(i)} + A_t^2(i,j) \frac{p_t^1(j)}{p_t^1(i)}$$
, (110)

$$g_t(i) := \sum_{j \neq i} \left(A_t^{\text{prod}}(i,j) - A_t^1(i,j) - A_t^2(i,j) \right). \tag{111}$$

Corollary B.2. [Product of Masked Diffusions] For the rate matrix of the reverse-time masked diffusion from Eq. (10), Theorem 3.3 yields

$$B_{\tau}^{\text{prod}}(i,j) = -2\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^1(j)}{p_t^1(m)} \frac{p_t^2(j)}{p_t^2(m)}, g_{\tau}(i) = \frac{\delta_{mi}}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_j \frac{p_t^1(j)}{p_t^1(m)} + \frac{p_t^2(j)}{p_t^2(m)} - 2\frac{p_t^1(j)}{p_t^1(m)} \frac{p_t^2(j)}{p_t^2(m)}$$

Proof. The reverse-time rate matrices are

$$B_t^1(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^1(j)}{p_t^1(m)}, \quad B_t^2(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^2(j)}{p_t^2(m)}. \tag{112}$$

Then, according to Theorem 3.3, the rate matrix for the product is

$$B_t^{\text{prod}}(i,j) = B_t^1(i,j) \frac{p_t^2(j)}{p_t^2(i)} + B_t^2(i,j) \frac{p_t^1(j)}{p_t^1(i)}$$
(113)

$$= -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^1(j)}{p_t^1(m)} \frac{p_t^2(j)}{p_t^2(i)} - \delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^2(j)}{p_t^2(m)} \frac{p_t^1(j)}{p_t^1(i)}$$
(114)

$$= -\delta_{mi} \frac{2}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^1(j)}{p_t^1(m)} \frac{p_t^2(j)}{p_t^2(m)}.$$
 (115)

And the weighting term is

$$g_t(i) = \sum_{j \neq i} \left(B_t^{\text{prod}}(i, j) - B_t^1(i, j) - B_t^2(i, j) \right)$$
 (116)

$$= \delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j} \left(\frac{p_t^1(j)}{p_t^1(m)} + \frac{p_t^2(j)}{p_t^2(m)} - 2 \frac{p_t^1(j)}{p_t^1(m)} \frac{p_t^2(j)}{p_t^2(m)} \right). \tag{117}$$

B.3 Geometric Average of FKEs

Theorem B.3. [Geometric Average of FKEs] Consider N forward Kolmogorov equations with marginals $p_t^n(i)$ and corresponding rate matrices $A_t^n(i,j)$. For the geometric average of

marginals $q_t(i) \propto \prod_{n=1}^N p_t^n(i)^{\beta_n}$, with $\sum_{i=1}^N \beta_n = 1$, the following equation holds

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{geom}}(j, i) \ q_t(j) - A_t^{\text{geom}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{q_t(j)} g_t(j) \right), \tag{118}$$

where
$$A_t^{\text{geo}}(i,j) := \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(i)}\right)^{\beta_n} \sum_{n=1}^N \beta_n A_t^n(j,i) \frac{p_t^n(j)}{p_t^n(i)},$$
 (119)

$$g_t(i) := \sum_{i \neq i} \left(A_t^{\text{geo}}(i, j) - \sum_{n=1}^N \beta_n A_t^n(i, j) \right).$$
 (120)

Proof. We define the target marginals as

$$q_t(i) := \frac{1}{Z_t} \prod_{n=1}^N p_t^n(i)^{\beta_n} , \quad Z_t = \sum_i \prod_{n=1}^N p_t^n(i)^{\beta_n} . \tag{121}$$

Hence, the time derivative of the marginals is

$$\frac{\partial}{\partial t} \log q_t(i) = \sum_{n=1}^{N} \beta_n \frac{\partial}{\partial t} \log p_t^n(i) - \frac{\partial}{\partial t} \log Z_t$$
 (122)

$$= \sum_{i \neq i} \sum_{n=1}^{N} \beta_n \left(A_t^n(j, i) \frac{p_t^n(j)}{p_t^n(i)} - A_t^n(i, j) \right) - \frac{\partial}{\partial t} \log Z_t$$
 (123)

$$= \sum_{j \neq i} \left(\underbrace{\sum_{n=1}^{N} \beta_n A_t^n(j,i) \frac{p_t^n(j)}{p_t^n(i)} \frac{q_t(i)}{q_t(j)}}_{:=A_t^{\text{geom}}(j,i)} \underbrace{q_t(j)}_{q_t(i)} - A_t^{\text{geom}}(i,j) \right) + \tag{124}$$

$$+\sum_{i\neq i} \left(A_t^{\text{geom}}(i,j) - \sum_{n=1}^N \beta_n A_t^n(i,j) \right) - \frac{\partial}{\partial t} \log Z_t \,. \tag{125}$$

Denoting

$$A_t^{\text{geom}}(i,j) := \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(i)} \right)^{\beta_n} \sum_{n=1}^N \beta_n A_t^n(i,j) \frac{p_t^n(i)}{p_t^n(j)}, \text{ and}$$
 (126)

$$g_t(i) := \sum_{j \neq i} \left(A_t^{\text{geom}}(i, j) - \sum_{n=1}^N \beta_n A_t^n(i, j) \right), \tag{127}$$

we can describe the evolution of the marginals $q_t(i)$ as

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} (A_t^{\text{geom}}(j, i) q_t(j) - A_t^{\text{geom}}(i, j) q_t(i)) + q_t(i) (g_t(i) - \mathbb{E}_{j \sim q_t(j)} g_t(j)). \tag{128}$$

Corollary B.4. [Geometric Average of Masked Diffusions] For the rate matrix of the reversetime masked diffusion from Eq. (10), Theorem B.3 yields

$$B_t^{\text{geom}}(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(m)} \right)^{\beta_n}, \ i \neq j$$
 (129)

$$g_t(i) = \delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j \neq i} \left(\sum_{n=1}^N \beta_n \frac{p_t^n(j)}{p_t^n(m)} - \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(m)} \right)^{\beta_n} \right). \tag{130}$$

Proof. For the reverse-time masked diffusion, we have

$$B_t^n(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t^n(j)}{p_t^n(m)}, \ i \neq j, \ n = 1, \dots, N.$$
(131)

Using the result of Theorem B.3, we have

$$B_t^{\text{geom}}(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(i)} \right)^{\beta_n} \sum_{n=1}^N \beta_n B_t^n(i,j) \frac{p_t^n(i)}{p_t^n(j)}$$
(132)

$$= -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(i)} \right)^{\beta_n} \sum_{n=1}^N \beta_n \frac{p_t^n(j)}{p_t^n(m)} \frac{p_t^n(i)}{p_t^n(j)}$$
(133)

$$= -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(m)} \right)^{\beta_n}, \tag{134}$$

where in the last transition we have used the fact that the expression is zero unless i=m and $\sum_{n=1}^{N} \beta_n = 1$. Correspondingly, the weights are

$$g_t(i) = \sum_{j \neq i} \left(B_t^{\text{geom}}(i,j) - \sum_{n=1}^N \beta_n B_t^n(i,j) \right)$$
(135)

$$= \delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j \neq i} \left(\sum_{n=1}^N \beta_n \frac{p_t^n(j)}{p_t^n(m)} - \prod_{n=1}^N \left(\frac{p_t^n(j)}{p_t^n(m)} \right)^{\beta_n} \right). \tag{136}$$

B.4 Reward-Tilted FKE

Theorem 3.5. [Reward-tilted FKE] Consider the forward Kolmogorov equation from Eq. (2) describing the time evolution of the marginals $p_t(i)$ with the rate matrix $A_t(i,j)$. For the reward-tilted marginals $q_t(i) \propto p_t(i) \exp(\beta_t r(i))$, the following equation holds

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{reward}}(j, i) q_t(j) - A_t^{\text{reward}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{q_t(j)} g_t(j) \right), \quad (17)$$

$$A_t^{\text{reward}}(i,j) := A_t(i,j) \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(i))}, \ g_t(i) := \sum_{i \neq j} \left(A_t^{\text{reward}}(i,j) - A_t(i,j) \right) + \frac{\partial \beta_t}{\partial t} r(i). \tag{18}$$

Proof. We define

$$q_t(i) := \frac{1}{Z_t} p_t(i) \exp(\beta_t r(i)), \quad Z_t = \sum_i p_t(i) \exp(\beta_t r(i))$$
(137)

The derivative of the log-probability is

$$\frac{\partial}{\partial t} \log q_t(i) = \sum_{j \neq i} \left(A_t(j, i) \frac{p_t(j)}{p_t(i)} - A_t(i, j) \right) + \frac{\partial \beta_t}{\partial t} r(i) - \frac{\partial}{\partial t} \log Z_t$$
 (138)

$$= \sum_{j \neq i} \left(\underbrace{A_t(j,i) \frac{\exp(\beta_t r(i))}{\exp(\beta_t r(j))}}_{:=A_t^{\text{reward}}(j,i)} \frac{q_t(j)}{q_t(i)} - A_t(i,j) \right) + \frac{\partial \beta_t}{\partial t} r(i) - \frac{\partial}{\partial t} \log Z_t$$
 (139)

$$= \sum_{j \neq i} \left(A_t^{\text{reward}}(j, i) \frac{q_t(j)}{q_t(i)} - A_t^{\text{reward}}(i, j) \right) + \tag{140}$$

$$+\underbrace{\sum_{j\neq i} (A_t^{\text{reward}}(i,j) - A_t(i,j)) + \frac{\partial \beta_t}{\partial t} r(i)}_{:=g_t(i)} - \frac{\partial}{\partial t} \log Z_t$$
(141)

To show the following equality

$$g_t(i) - \frac{\partial}{\partial t} \log Z_t = g_t(i) - \mathbb{E}_{i \sim q_t(j)} g_t(j), \qquad (142)$$

one can either use the definition of $q_t(i)$ and its normalization, or explicitly calculate the derivative of the normalizing constant, i.e.

$$\frac{\partial}{\partial t} \log Z_t = \frac{1}{Z_t} \sum_i \frac{\partial}{\partial t} \Big(p_t(i) \exp(\beta_t r(i)) \Big)$$
(143)

$$= \sum_{i} q_{t}(i) \left(\frac{\partial}{\partial t} \log p_{t}(i) + \frac{\partial \beta_{t}}{\partial t} r(i) \right)$$
(144)

$$= \sum_{i} q_t(i) \left(\sum_{j \neq i} \left(A_t(j, i) \frac{p_t(j)}{p_t(i)} - A_t(i, j) \right) + \frac{\partial \beta_t}{\partial t} r(i) \right)$$
(145)

Thus, we have

$$\sum_{i} q_{t}(i) g_{t}(i) - \frac{\partial}{\partial t} \log Z_{t} = \sum_{i} q_{t}(i) \left(\left(\sum_{i \neq i} A_{t}^{\text{reward}}(i, j) - A_{t}(i, j) \right) + \frac{\partial \beta_{t}}{\partial t} r(i) \right)$$
(146)

$$-\left(\sum_{j\neq i} A_t(j,i) \frac{p_t(j)}{p_t(i)} - A_t(i,j)\right) - \frac{\partial \beta_t}{\partial t} r(i)\right)$$
(147)

$$= \sum_{i} q_t(i) \sum_{j \neq i} \left(A_t^{\text{reward}}(i,j) - A_t(j,i) \frac{p_t(j)}{p_t(i)} \right)$$
(148)

$$= \sum_{i} q_{t}(i) \sum_{j \neq i} \left(A_{t}(i,j) \frac{\exp(\beta_{t} r(j))}{\exp(\beta_{t} r(i))} - A_{t}(j,i) \frac{p_{t}(j)}{p_{t}(i)} \right)$$
(149)

$$= \frac{1}{Z_t} \sum_{i} \sum_{j \neq i} \left(A_t(i,j) \exp(\beta_t r(j)) p_t(i) - A_t(j,i) \exp(\beta_t r(i)) p_t(j) \right)$$
(150)

$$= \frac{1}{Z_t} \sum_{i} \sum_{j \neq i} \left(\hat{A}_t(i,j) - \hat{A}_t(j,i) \right) = \frac{1}{Z_t} \sum_{i,j} \left(\hat{A}_t(i,j) - \hat{A}_t(j,i) \right) = 0, \quad (151)$$

where we denote

$$\hat{A}_t(i,j) := A_t(i,j) \exp(\beta_t r(j)) p_t(i). \tag{152}$$

Finally, we have

$$\frac{\partial q_t(i)}{\partial t} = \sum_{j \neq i} \left(A_t^{\text{reward}}(j, i) q_t(j) - A_t^{\text{reward}}(i, j) q_t(i) \right) + q_t(i) \left(g_t(i) - \mathbb{E}_{j \sim q_t(j)} g_t(j) \right), \quad (153)$$

$$A_t^{\text{reward}}(i,j) \coloneqq A_t(i,j) \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(i))}, \quad g_t(i) \coloneqq \sum_{j \neq i} \left(A_t^{\text{reward}}(i,j) - A_t(i,j) \right) + \frac{\partial \beta_t}{\partial t} r(i).$$

Corollary B.5. [Reward-tilted Masked Diffusion] For the rate matrix of the reverse-time masked diffusion from Eq. (10), Theorem 3.5 yields

$$B_{\tau}^{\text{reward}}(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)} \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(m))},$$
(19)

$$g_{\tau}(i) = \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \delta_{mi} \sum_{j} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t(j)}{p_t(m)} \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(m))} \right) + \frac{\partial \beta_t}{\partial t} r(i).$$
 (20)

Proof. The reverse-time rate matrix is

$$B_t(i,j) = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)}.$$
 (154)

Then the reward-weighted matrix is

$$B_t^{\text{reward}}(i,j) = B_t(i,j) \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(i))} = -\delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \frac{p_t(j)}{p_t(m)} \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(m))}, \quad (155)$$

and the weighting term is

$$g_t(i) = \sum_{j \neq i} \left(B_t^{\text{reward}}(i, j) - B_t(i, j) \right) + \frac{\partial \beta_t}{\partial t} r(i)$$
(156)

$$= \delta_{mi} \frac{1}{\alpha_t} \frac{\partial \alpha_t}{\partial t} \sum_{j} \left(\frac{p_t(j)}{p_t(m)} - \frac{p_t(j)}{p_t(m)} \frac{\exp(\beta_t r(j))}{\exp(\beta_t r(m))} \right) + \frac{\partial \beta_t}{\partial t} r(i)$$
 (157)

C Experimental Details

Code is available at https://anonymous.4open.science/r/discrete_fkc-40B8/.

C.1 Amortized Linear Regression

C.1.1 Theoretical Justification

The posterior over parameters factors as:

$$p(\theta|\mathcal{X}) \propto p(\theta)p(\mathcal{X}|\theta) = p(\theta) \prod_{k=1}^{K} p(\mathcal{X}_{k}|\theta) \propto p(\theta)^{1-K} \prod_{k=1}^{K} p(\theta|\mathcal{X}_{k})$$
(158)

For a uniform prior $p(\theta)$, this results in the product we applied $p(\theta|\mathcal{X}) \propto \prod_{k=1}^K p(\theta|\mathcal{X}_k)$.

C.1.2 Experimental Setup

All experiments were done on a single A100 GPU.

For each experiment, the dataset \mathcal{X} was generated using $(\theta_0, \theta_1) = (3.0, 4.0)$, with x spaced linearly between [-10, 10], and $y_i = \theta_1^* x_i + \theta_0^* + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 0.1^2)$.

For inference with LLaDA, a temperature of 1.0 was used, and the random remasking strategy was applied. All predictions were made in a single block, and the generation length was capped at 128 tokens.

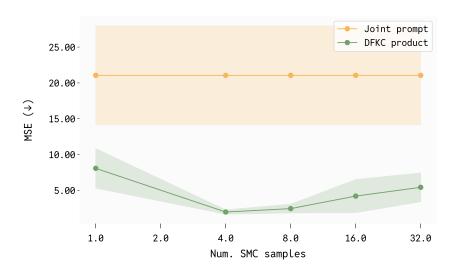


Figure A1: Increasing the number of SMC samples for DFKC improves over no SMC resampling; gain is largest with 4 or 8 samples. Taking the product has a lower (better) mean squared error (MSE) than joint prompting, and resampling with DFKC significantly improves this further.

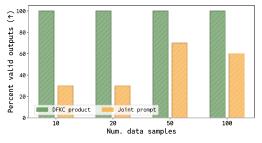
The prompt used to generate predictions is of the form: "Assume a model of the form y = a * x + b, where a and b are the parameters of the model. The observations are given as (x,y) points, where y has Gaussian noise with standard deviation 0.1 added. Predict the parameters of linear regression for (x,y) points: " + (x_1,y_1) , ..., (x_N,y_N) + " Output the final answer as: "The best estimate for parameters of the model are: $a = _$, and $b = _$ " where $_$ is replaced with the values of a and then b."

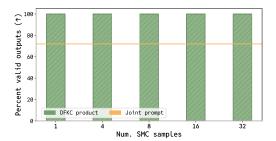
C.1.3 Additional Results for Amortized Learning

We include an ablation over the number of SMC samples, for a fixed number of products in Fig. A1. We can observe that more SMC samples improves performance, up to a threshold of 8 samples.

We additionally include a comparison of how well the outputs adhered to the specified prompt format in Fig. A2.

Some selected samples from the product and joint prompting strategies are included in Table A1. We can note that outputs using joint prompting often fail to adhere to the output format specified in the prompt, and sometimes cannot be parsed for values of (θ_0, θ_1) . This issue wasn't observed for the product prompt (using any number of particles).





(a) DFKC generates a higher percentage of valid, parseable outputs compared with joint prompting at all data sizes.

(b) DFKC generates consistently generates 100% valid, parseable outputs at all SMC sample sizes while joint prompting only generates 72% valid prompts on average.

Figure A2: Effect of data quantity on predicting linear regression parameters.

C.2 Multi-constraint Story Generation

All experiments were performed on a single L40 GPU.

For inference with LLaDA, the next token to unmask was chosen randomly (as opposed to picking highest confidence one) due to an issue where the model would sample an end of text token often, using the latter setting. All experiments used a temperature of T=1.0.

The prompt for the story generation is composed as follows: the base prompt is "Write a story.". The conditions are sampled at random from a set of 50 conditions, containing mutually compatible constraints such as:

```
    "It should include a curious child."
    "It should describe a small village."
    "It should feature a dense forest."
    ...
```

C.3 Protein Sequence Generation

All experiments were done on a single L40 GPU. For each length and particle type, we generate 50 samples.

The log-reward for a sequence x with length L is defined using the ESM2 model f_{θ} . We first compute a score S(x) by passing the entire sequence x to the model, and then averaging the log-likelihoods evaluated at the amino acid sequence:

$$S(x) = \frac{1}{L} \sum_{i=1}^{L} f_{\theta}(x)[x_i]$$
 (159)

This approach allows us to compute the toy reward in one single pass.

The score is scaled by a hyperparameter γ to obtain the log-reward r(x):

$$r(x) = \gamma S(x) \tag{160}$$

In our experiments for unconditional protein sequence, we set a hyperparameter $\gamma = 200$ across all lengths and particles.

The base discrete diffusion model used is DPLM1 650M (Wang et al., 2024b). For a sequence of length l, l generation steps are used, and once a token is unmasked, it is not remasked in future steps (to align more closely to the traditional masked diffusion generation process, and as opposed to the remasking strategies used in (Wang et al., 2024b)).

A linear annealing schedule $\beta_t = 1 - t$ is used for the reward (where generation starts at t = 1 and proceeds to t = 0).

C.4 Additional Experimental Results for Annealing the Ising Model

Dataset for experiment 1

```
source: synthetic Ising model configurations
size: 100,000 samples after burn-in
sampling method: Swendsen-Wang
   burn-in length: 10,000 steps
   thinning interval: 5
beta: 0.25
lattice size: 16
Dataset for experiment 2
```

```
1508    source: synthetic Ising model configurations
1509    size: 10,000 samples after burn-in
1510    sampling method: Glauber dynamics
1511    burn-in length: 10,000 steps
    thinning interval: 1
```

```
1512
      beta: 0.2 and 0.3
1513
      lattice size: 16
1514
1515
      Model
1516
      architecture: UNet
1517
      activation: SiLU
1518
      channels: [64, 128, 256]
1519
      resblocks per stage: 2
1520
      attention: applied at 4×4 resolution
1521
      initialization: Xavier uniform
      time embedding: sinusoidal embedding
1522
1523
      Training
1524
1525
      optimizer: Adam
1526
          learning rate: 1e-4
          betas: (0.9, 0.999)
1527
      batch size: 256
1528
      epochs: 600
1529
      learning rate schedule: constant with warmup
1530
      hardware: 1 × NVIDIA A100 GPU (40 GB memory)
1531
      loss: denoising score entropy
1532
1533
      Evaluation
1534
      metrics for global structure: 2-Wasserstein metric between
1535
      distributions of
1536
      energy and distributions of magnetization.
1537
      metrics for local structure: MSE for correlation function.
1538
      sample size: 10,000
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
```

Data size	Joint Prompt Output	Product Prompt Output
N=10	The best estimate for parameters of the model are: $a = 4.337$, and $b = -34.049$	The best estimate for parameters of the model are: a = 3.000, and b = 10.004
N=20	Based on the observed data points, we can see a trend that y is directly proportional to x. The best estimate for the parameters a and b is: a = 1.0, and b = 0.0.	The best estimate for parameters of the model are: a = 3.82, and b = 10.12.
N=50	To obtain the best estimates for the parameters (a and b), you need to follow the detailed steps of building a linear regression model using Ordinary Least Squares (also namedIM, and guide, filter). These steps involve typically a program such as R or a statistical tool among others. The objective is to predict parameters, but after an ensemble calculation, we are going to use, known as the sum of residuals, to estimate the model's parameters. The sum of residuals helps us evaluate the discrepancy of model with a given residuals. Once I've made these predictions, I'll be able to provide more precise feedback on parameter estimates.	The best estimate for parameters of the model are: a = 1.344, and b = -22.331
N=100	The best estimate for parameters of the model are: a = 0x583C622F 052D29A9 + 00EA6F242949D26F and b = 0x 41796E30 0027A200 - 76CF406498D45505. Note: These values of a and b are with 95% confidence taking into account the Gaussian balls added to Python and Python recovery points.	The best estimate for parameters of the model are: a = 0.8313, and b = 0.0564.

Table A1: Comparison between curated joint and product prompt outputs at varying data sizes.