

# EFFICIENT PERPLEXITY BOUND AND RATIO MATCHING IN DISCRETE DIFFUSION LANGUAGE MODELS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

While continuous diffusion models excel in modeling continuous distributions, their application to categorical data has been less effective. Recent work has shown that ratio-matching through *score-entropy* within a continuous-time discrete Markov chain (CTMC) framework serves as a competitive alternative to autoregressive models in language modeling. To enhance this framework, we first introduce three new theorems concerning the KL divergence between the data and learned distribution. Our results serve as the discrete counterpart to those established for continuous diffusion models and allow us to derive an improved upper bound of the perplexity. Second, we empirically show that ratio-matching performed by minimizing the *denoising cross-entropy* between the clean and corrupted data enables models to outperform those utilizing score-entropy with up to 10% lower perplexity/generative-perplexity, and 15% faster training steps. To further support our findings, we introduce and evaluate a novel CTMC transition-rate matrix that allows prediction refinement, and derive the analytic expression for its matrix exponential which facilitates the computation of conditional ratios thus enabling efficient training and generation.

## 1 INTRODUCTION

Modeling data distributions is a fundamental task in machine learning. In the case of continuous data distributions, recent advancements in continuous diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020b) have demonstrated impressive capabilities in generating data samples and performing density estimation (Song et al., 2021a;c). Despite these achievements, the application of such models to categorical data distributions, like language, remains limited, as continuous diffusion models generally underperform compared to autoregressive models in these scenarios (Chen et al., 2023; Gulrajani & Hashimoto, 2024; Li et al., 2022; Dieleman et al., 2022; Strudel et al., 2022). To address this, recent research has focused on the development of discrete diffusion models (Austin et al., 2021; Campbell et al., 2022; Meng et al., 2022; Lou et al., 2024; Sahoo et al., 2024; Ou et al., 2024) which offer distinct advantages compared to autoregressive models, such as the ability to infill various parts of a sequence non-sequentially and have the potential to reduce computing time and expenses in generating lengthy sequences (Deschenaux & Gulcehre, 2024; Christopher et al., 2024).

Evaluating discrete diffusion models, however, presents a practical challenge due to the difficulty in calculating the perplexity, unlike in autoregressive models where this computation is straightforward. Although a recent perplexity bound has been proposed by Lou et al. (2024), no tightness guarantees exist. In this paper, we present three theorems concerning the Kullback-Leibler (KL) divergence between the data distribution and the learned distribution in discrete diffusion models (Lou et al., 2024). These results serve as the discrete analogue of the continuous diffusion theorems provided in Song et al. (2021b). One of our key contributions is Theorem 4, which provides an upper bound ( $J_2$ ) on the cross-entropy between the data and learned distributions, offering a more direct way of bounding the perplexity. This bound is computationally more efficient than the existing bound in Lou et al. (2024), and empirical results suggest that it is also tighter.

In addition, inspired by model reparametrizations (Ho et al., 2020; Karras et al., 2022; Lou et al., 2024), this paper examines the ratio-matching training objective in Lou et al. (2024). We highlight that the sole unknowns implicitly learned by the model are the per-token marginal probabilities

054 across the vocabulary, conditioned on the current perturbed sequence. Consequently, rather than  
 055 modeling the ratios directly, we employ cross-entropy discrete diffusion (CEDD) to model these  
 056 probabilities, paralleling the methodology used in autoregressive models. This method coincides  
 057 with the training approach  $L_{ll}$  proposed in Campbell et al. (2022). Moreover, this strategy mirrors  
 058 the cross-entropy loss utilized in continuous diffusion models, as described in Dieleman et al. (2022),  
 059 thereby drawing another connection between discrete and continuous diffusion models. Similar re-  
 060 sults for the cross-entropy loss have been reported in studies of *absorb* discrete diffusion, as shown  
 061 in concurrent research by Sahoo et al. (2024); Ou et al. (2024); Shi et al. (2024). These particu-  
 062 lar results, up to a scaling factor, are specific cases within the broader CEDD framework, which is  
 063 applicable to any discrete diffusion process, independent of the transition-rate matrix. CEDD im-  
 064 proves upon SEDD (Lou et al., 2024) by circumventing the learning of conditional ratios—which  
 065 are analytically determinable—focusing instead on learning the mixing weights that constitute the  
 066 necessary marginal ratios for generation. This focus not only conserves modeling resources but is  
 067 also particularly advantageous when the distribution of conditional ratios is complex.

068 In order to illustrate our point, we design a new transition-rate matrix named *roulette* diffusion, and  
 069 derive its matrix exponential. The roulette diffusion is an interpolation between the absorb and the  
 070 uniform diffusion. In the forward process, a token can transition to any state until it hits the absorb  
 071 state, that is, until it is masked. In return, the reverse process begins with a sequence of masked to-  
 072 kens, which are gradually unmasked, and where the unmasked tokens can be refined. Intuitively, this  
 073 capability should be important for discrete diffusion models (Deschenaux & Gulcehre, 2024). More-  
 074 over, during the reverse process, the scores/ratios of unmasked tokens have much larger magnitudes  
 075 than those of masked tokens, posing a significant learning challenge for the network due to output  
 076 scale variability. However, employing cross-entropy (CEDD) training mitigates this challenge as  
 demonstrated experimentally.

077 In summary, the main **contributions** of this paper include:

- 079 • Providing 3 new theorems concerning the KL divergence and cross entropy between the data and  
 080 the learned distribution. Improving model evaluation through the bound provided in Theorem 4.
- 081 • Introducing a new transition-rate matrix (roulette diffusion) that allows token correction after  
 082 unmasking in the reverse process. Deriving its matrix exponential which enables efficient training  
 083 using SEDD and generation when CEDD is employed.
- 084 • Comparing the performance of SEDD and CEDD experimentally in the task of language mod-  
 085 elling on absorb, uniform and roulette diffusion models. Showing that with simple modifications  
 086 of the reparametrization, CEDD outperforms SEDD in all cases in terms of perplexity.

## 088 2 PRELIMINARIES AND NOTATION

### 089 2.1 MARKOV CHAINS OVER FINITE-STATE SPACES

092 A discrete-time Markov Chain in a finite-state space is a stochastic process  $X_1, X_2, \dots, X_{\bar{T}}$ , where  
 093 each state  $X_t$  depends solely on the preceding one. The states  $X_t$  can take values from  $\{1, 2, \dots, S\}$ ,  
 094 and  $\bar{T}$  represents the number of time steps. The probability of being in state  $x$  at time  $t$  is

$$095 p_t(X_t = x) = \sum_{y=1}^S p(X_t = x, X_{t-1} = y) = \sum_{y=1}^S p_{t|t-1}(X_t = x|X_{t-1} = y)p_{t-1}(X_{t-1} = y). \quad (1)$$

098 Placing all such probabilities  $p_t(X_t = x)$  in a vector  $\mathbf{s}_t$ , such that  $\mathbf{s}_t(x) = p_t(X_t = x)$ , gives

$$099 \mathbf{s}_t = \mathbf{P}\mathbf{s}_{t-1}, \text{ where } \mathbf{P}(x, y) = p_{t|t-1}(X_t = x|X_{t-1} = y). \quad (2)$$

101 One can generalize such processes into Continuous Time Markov Chains (CTMCs) where  $t \in [0, \bar{T}]$ ,  
 102 (Anderson, 2012). For simplicity, we make the choice  $\bar{T} = 1$ . To construct a CTMC, one first  
 103 chooses a transition-rate matrix  $\mathbf{Q}_t$ , which has the property that its non-diagonal elements are non-  
 104 negative, and the elements in each of its columns add to zero (Suhov & Kelbert, 2008). Given an  
 105 initial probability distribution  $\mathbf{s}_0$ , the equation below fully determines the evolution of the probabil-  
 106 ity with respect to time:

$$107 \frac{d\mathbf{s}_t}{dt} = \mathbf{Q}_t\mathbf{s}_t. \quad (3)$$

In addition, we choose  $\mathbf{Q}_t = \sigma'(t)\mathbf{Q}$ , where  $\mathbf{Q}$  is itself a constant transition-rate matrix and where function  $\sigma$  is monotonically increasing, and satisfies  $\sigma(0) = 0$  as well as  $\lim_{t \rightarrow 1} \sigma(t) = T$ . In this setting, the distribution over states at time  $t$  is the solution of the linear ODE in Equation (3), that is,  $\mathbf{s}_t = e^{\sigma(t)\mathbf{Q}}\mathbf{s}_0$ .

Matrices  $\mathbf{Q}_t$  are chosen such that: a) the matrix exponential  $e^{\sigma(t)\mathbf{Q}}$  is easy to calculate, which is essential as  $p_{t|0}(x|y) = e^{\sigma(t)\mathbf{Q}}(x, y)$ ; and b)  $\mathbf{s}_1$  is an easy reference distribution to sample from (Austin et al., 2021; Campbell et al., 2022).

Finally, similar to diffusion processes in continuous spaces, the continuous-time Markov chain in Equation (3) also admits a reverse process (Kelly, 1979; Sun et al., 2022):

$$\frac{d\mathbf{s}_{1-t}}{dt} = \bar{\mathbf{Q}}_{1-t}\mathbf{s}_{1-t}, \quad (4)$$

where  $\bar{\mathbf{Q}}_t(x, y) = \mathbf{Q}_t(y, x)\frac{p_t(x)}{p_t(y)}$  for  $x \neq y$ , and  $\bar{\mathbf{Q}}_t(x, x) = -\sum_{y \neq x} \bar{\mathbf{Q}}_t(y, x)$ . Since we can easily sample from the reference distribution, the only unknowns preventing us from being able to run backwards are the ratios  $\frac{p_t(x)}{p_t(y)}$  also known as concrete scores (Meng et al., 2022; Lou et al., 2024). Once such ratios are modeled using a neural network, we can generate samples from the learned data distribution  $p_0^\theta$  by discretizing Equation (4) as follows:

$$p(x_{t-\epsilon} = y \mid x_t = x) = \delta_x(y) + \bar{\mathbf{Q}}_t(y, x)\epsilon + O(\epsilon^2). \quad (5)$$

Additional details are provided in Appendix D.

## 2.2 SEDD: ESTIMATING THE RATIOS VIA SCORE ENTROPY

As pointed out in the previous subsection, we wish to model the ratios  $\frac{p_t(y)}{p_t(x_t)}$  via a neural network  $s_\theta(x_t, t)_y$ , for example by minimizing the score entropy loss (Lou et al., 2024):

$$\mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} w_{x_t, y} \ell \left( \frac{p_t(y)}{p_t(x_t)}, s_\theta(x_t, t)_y \right), \text{ for } \ell(a, b) = (b - a \log b + K(a)), \quad (6)$$

and  $K(a) = a(\log a - 1)$ . In Lou et al. (2024),  $w_{x_t, y} = \mathbf{Q}_t(x_t, y)$ , and furthermore they show that an equivalent loss is the following:

$$\mathbb{E}_{x_0 \sim p_0, x_t \sim p_{t|0}(\cdot|x_0)} \sum_{y \neq x_t} w_{x_t, y} \ell \left( \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)}, s_\theta(x_t, t)_y \right), \quad (7)$$

which side-steps the problem of not knowing the marginal ratios  $\frac{p_t(y)}{p_t(x_t)}$ , by employing  $p_{t|0}(i|j) = e^{\sigma(t)\mathbf{Q}}(i, j)$ . A more detailed derivation of Equation (7) can be found in Appendix A.2.

## 2.3 DISCRETE DIFFUSION FOR LANGUAGE MODELING - TOKEN LEVEL TRANSITIONS

In the case of Language Modeling, we write a sequence of length  $L$  from the data distribution as  $\mathbf{x}_0$ , where  $\mathbf{x}_0 = (x_0^1, x_0^2, \dots, x_0^L)$  and  $x_0^i \in \text{Vocab} = \{1, 2, \dots, V\}$ . The number of possible sequences, that is, the number of states  $S$  is  $V^L$ . Unfortunately, this implies that it is not computationally feasible to model the ratios of probabilities between the current state  $\mathbf{x}_t$  and all other states  $\mathbf{y}$ , since the output of our neural network would have to be  $V^L$  dimensional (Campbell et al., 2022).

We follow the usual approach (Campbell et al., 2022; Lou et al., 2024) to mitigate this issue, which is to select a sparse matrix  $\mathbf{Q}_t(S \times S)$ , such that each entry  $\mathbf{Q}_t(\mathbf{x}, \mathbf{y})$  for two sequences  $\mathbf{x}, \mathbf{y}$  that differ in more than one token will be zero. The forward process that such a  $\mathbf{Q}_t$  defines, can equivalently be described as follows: at each discretized step, only one uniformly randomly chosen token from the current sequence can be modified, according to a token level forward diffusion process  $\mathbf{Q}_t^{\text{tok}}(V \times V)$ . More formally, for  $\mathbf{x} = (x_0^1, \dots, x^i, \dots, x_0^L)$  and  $\mathbf{y} = (x_0^1, \dots, \hat{x}^i, \dots, x_0^L)$ , if  $x^i \neq \hat{x}^i$  we have  $\mathbf{Q}_t(\mathbf{x}, \mathbf{y}) = \mathbf{Q}_t^{\text{tok}}(x^i, \hat{x}^i)$ , otherwise  $\mathbf{Q}_t$  is zero in other non-diagonal entries. For such a sparse choice of  $\mathbf{Q}_t$ , and  $\mathbf{y}$  which only differs from  $\mathbf{x}$  at a single position  $i$ , Expression (5) becomes

$$p(x_{t-\epsilon} = \mathbf{y} \mid x_t = \mathbf{x}) = \begin{cases} 1 - \sum_{z^i \in \text{Vocab} \setminus x^i} \mathbf{Q}_t^{\text{tok}}(x^i, z^i) \frac{p_t(z^i)}{p_t(x^i)} \epsilon + O(\epsilon^2). & \text{if } \mathbf{y} = \mathbf{x} \\ \mathbf{Q}_t^{\text{tok}}(x^i, \hat{x}^i) \frac{p_t(\hat{x}^i)}{p_t(x^i)} \epsilon + O(\epsilon^2). & \text{if } \mathbf{y} \neq \mathbf{x} \end{cases} \quad (8)$$

where  $z$  denotes a sequence that is identical to  $x$  everywhere, but position  $i$ . Thus, the usual approach entails only modeling the ratios between  $x$  and ‘neighbours’  $y$  which only differ from  $x$  by one token. The number of such neighbours is  $L \times V$ , that is  $V$  per each of the  $L$  positions, hence the output of the network is  $L \times V$  coinciding with that of transformers in autoregressive language models. It should be pointed out that one can indeed use Expression (7) for training, due to the fact that tokens are perturbed independently from one another in the forward process  $p_{t|0}(x_t|x_0) = \prod_j p_{t|0}(x_t^j|x_0^j)$ , and thus

$$\frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)} = \frac{\prod p_{t|0}(y^j|x_0^j)}{\prod p_{t|0}(x_t^j|x_0^j)} = \prod \frac{p_{t|0}(y^j|x_0^j)}{p_{t|0}(x_t^j|x_0^j)} = \frac{p_{t|0}(y^i|x_0^i)}{p_{t|0}(x_t^i|x_0^i)} = \frac{e^{\mathcal{Q}_t^{tok}(y^i, x_0^i)}}{e^{\mathcal{Q}_t^{tok}(x_t^i, x_0^i)}}, \quad (9)$$

Finally, the noise schedule  $\sigma_t$  is typically loglinear  $-\log(1 - (1 - \epsilon)t)$  or geometric  $\sigma_{min}^{1-t} \cdot \sigma_{max}^t$ .

### 3 METHODOLOGY AND THEORETICAL RESULTS

In Subsection 3.1, we provide results related to the cross entropy and the KL divergence between the data and the learned distribution in the CTMC (discrete diffusion) framework. The first three theorems therein can be considered as the discrete diffusion analog of the ones given in (Song et al., 2021b). Importantly, Theorem 4 provides an upper bound ( $J_2$ ) on the cross entropy between the data and learned distribution which can be used to bound the perplexity, and which does not depend on the function  $K$  (Equation 6). We emphasize that the results hold for general CTMCs, and not only in the special case of token-level transitions. From the second subsection onwards, we operate in the token-level transition framework. More precisely, in Subsection 3.2, we introduce the roulette transition-rate matrix, and provide an expression for its exponential. In Subsection 3.3, we state Proposition 6, which enables a more efficient estimation of  $J_2$ . In Subsection 3.4, we highlight that similarly to the continuous case (Dieleman et al., 2022), the ratios can be modeled using  $L_{ll}$  from Campbell et al. (2022). We then present how this approach is adapted in our experimental setup.

#### 3.1 CROSS ENTROPY AND KL DIVERGENCE RESULTS

We begin by finding an upper bound for the KL divergence between the data and the learned distribution. The proofs are provided in Appendix A.1.

**Theorem 1.** *Define a CTMC with transition matrix  $\mathbf{Q}_t$  that runs from time 0 to 1. The true reverse process defines a probability evolution  $p_t$  from  $p_1$  to the data distribution  $p_0$ , while the learned reverse process induces the evolution  $p_t^\theta$  from the reference distribution  $p_1^\theta = p_r$  to the approximation of the data distribution  $p_0^\theta$ . In this setting, the following KL divergence bound holds*

$$D_{KL}(p_0||p_0^\theta) \leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \ell \left( \frac{p_t(y)}{p_t(x_t)}, s_\theta(x_t, t)_y \right) dt + D_{KL}(p_1||p_r). \quad (10)$$

where  $\ell(a, b) = (b - a \log b + K(a))$  and  $K(a) = a(\log a - 1)$ .

The following theorem provides an expression for the entropy of the data distribution. Furthermore, it provides sufficient conditions for the bound given above to become tight.

**Theorem 2.** *Denote the intermediate distributions at time  $t$  determined by the true reverse process, and by the learned reverse process with  $p_t$  and  $p_t^\theta$ , respectively. We can write the entropy of the data distribution  $H(p_0)$  as*

$$H(p_0) = H(p_1) - \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) K \left( \frac{p_t(y)}{p_t(x_t)} \right) dt. \quad (11)$$

*In addition, if the learned ratios  $s_\theta(x_t, t)_y$  equal  $\frac{p_t^\theta(y)}{p_t^\theta(x_t)}$  and  $p_1 = p_1^\theta := p_r$ , then the inequality in Theorem 1, becomes an equality.*

A particular case where the conditions of the theorem above hold is when  $s_\theta(x_t, t)_y = \frac{p_t(y)}{p_t(x_t)}$  as then  $\frac{p_t^\theta(y)}{p_t^\theta(x_t)} = \frac{p_t(y)}{p_t(x_t)} = s_\theta(x_t, t)_y$ . The third theorem gives an upper bound of the negative log-likelihood at a single point. This is a central result in Lou et al. (2024, Theorem 3.6), but we restate it here for completeness, and provide an alternative, more detailed proof in Appendix A.1.

**Theorem 3.** Let  $p_0^\theta$  denote the learned distribution from which the reverse process samples. The negative log-likelihood of a state  $x_0$  being sampled by the reverse process can be bounded from above as follows,

$$-\log p_0^\theta(x_0) \leq \int_0^1 \mathbb{E}_{x_t \sim p_{t|0}(\cdot|x_0)} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \ell \left( \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)}, s_\theta(x_t, t)_y \right) dt + D_{KL}(p_{1|0}(\cdot|x_0) || p_r). \quad (12)$$

Since the noise schedule is chosen such that  $p_1 \approx p_r$ , so that  $D_{KL}(p_{1|0} || p_r)$  (Lou et al., 2024; Ou et al., 2024), and  $H(p_1, p_r) - H(p_r)$  are negligible, one can use Inequality (12) to bound the cross entropy  $\frac{1}{L} H(p_0, p_0^\theta) = \mathbb{E}_{x_0} [-\frac{1}{L} \log(p_0^\theta(x_0))]$ . One approach is to compute the RHS in Expression (12) for each point and then average results (Appendix B.5), which can be computationally expensive. Instead, we can divide by  $L$  and take the expectation with regards to data distribution on both sides of Expression (12) as in Ou et al. (2024), and calculate

$$J_1 = \frac{1}{L} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{x_0 \sim p_0(x_0)} \mathbb{E}_{x_t \sim p_{t|0}(\cdot|x_0)} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \ell \left( \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)}, s_\theta(x_t, t)_y \right). \quad (13)$$

Using Theorem 1 and 2, we provide another, direct upper bound on the cross-entropy between the data and learned distributions, which evades the computation of  $K$ .

**Theorem 4.** Under the conditions stated in Theorem 1, the following inequality for the cross entropy between the data and the learned distribution holds:

$$H(p_0, p_0^\theta) \leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \bar{\ell} \left( \frac{p_t(y)}{p_t(x_t)}, s_\theta(x_t, t)_y \right) dt - \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(y, x_t) dt + H(p_1, p_r), \text{ where } \bar{\ell}(a, b) = (b - a \log b). \quad (14)$$

We highlight that  $H(p_1, p_r) \approx H(p_r)$  and  $-\int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(y, x_t) dt$  can be analytically computed as shown in Section 3.3, Proposition 6, while the first term can be rewritten as

$$\mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{x_0 \sim p_0(x_0)} \mathbb{E}_{x_t \sim p_{t|0}(\cdot|x_0)} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \bar{\ell} \left( \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)}, s_\theta(x_t, t)_y \right). \quad (15)$$

Therefore, due to Theorem 4, we instead propose to use

$$J_2 = \frac{1}{L} \left[ \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{x_0 \sim p_0(x_0)} \mathbb{E}_{x_t \sim p_{t|0}(\cdot|x_0)} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \bar{\ell} \left( \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)}, s_\theta(x_t, t)_y \right) + H(p_r) - \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(y, x_t) dt \right]. \quad (16)$$

In both cases  $Perplexity = \exp(\mathbb{E}_{x_0} [-\frac{1}{L} \log(p_0^\theta(x_0))]) = \exp(\frac{1}{L} H(p_0, p_0^\theta)) \leq e^{J_1}, e^{J_2}$ .

### 3.2 ROULETTE DISCRETE DIFFUSION

Typically, matrices  $\mathbf{Q}^{tok}$  are defined as  $\mathbf{Q}^{tok} = \mathbf{P}^{tok} - \mathbf{I}$ , where  $\mathbf{P}^{tok}$  is idempotent, since this implies that  $(\mathbf{Q}^{tok})^2 = -\mathbf{Q}^{tok}$ . This last property of  $\mathbf{Q}^{tok}$  greatly simplifies the calculation of  $e^{\sigma(t)\mathbf{Q}^{tok}}$ . Indeed, using the Taylor series expansion we see that

$$e^{\sigma(t)\mathbf{Q}^{tok}} = \mathbf{I} + \sum_{k=1}^{\infty} \frac{(\mathbf{Q}^{tok})^k \sigma(t)^k}{k!} = \mathbf{I} + \sum_{k=1}^{\infty} \frac{(-1)^{k+1} \sigma(t)^k \mathbf{Q}^{tok}}{k!} = \mathbf{I} + \mathbf{Q}^{tok} (1 - e^{-\sigma(t)}). \quad (17)$$

Usually, the following two matrices  $\mathbf{P}^{tok}$  are chosen: a)  $\mathbf{P}_{uniform}^{tok} (V \times V)$  where each entry is equal to  $\frac{1}{V}$ , and b)  $\mathbf{P}_{absorb}^{tok} (n \times n)$  in which the last row is full of ones while all other elements are 0, and

where  $n = V + 1$ . We propose another transition-rate matrix whose exponential can be analytically calculated. To our knowledge this is the only matrix with such a property presented so far apart from the absorb and uniform ones. We refer to this new discrete diffusion process as *roulette* diffusion. The corresponding  $\mathbf{P}_{roulette}^{tok}(n \times n)$  is a matrix, such that  $\mathbf{P}_{roulette}^{tok}(i \neq n, j \neq n) = \frac{1}{V}(1 - p_m)$ ,  $\mathbf{P}_{roulette}^{tok}(n, j \neq n) = p_m$ ,  $\mathbf{P}_{roulette}^{tok}(i \neq n, n) = 0$  and  $\mathbf{P}_{roulette}^{tok}(n, n) = 1$ . We notice that for  $p_m = 1$ , roulette diffusion coincides with absorb diffusion, while for  $p_m = 0$  it coincides with the uniform diffusion. Intuitively, a token can transit from a non-absorb state to a non-absorb state with probability  $\frac{1}{V}(1 - p_m)$ , until it transits to the absorb state (with probability  $p_m$ ) and then remains there. While this matrix is not idempotent, one can still calculate its exponential as stated in the following proposition which is proved in Appendix A.3.3:

**Proposition 5.** *If we denote with  $\mathbf{Y}_t$  the matrix exponential of  $\sigma_t \mathbf{Q}_{roulette}^{tok} = \sigma_t (\mathbf{P}_{roulette}^{tok} - \mathbf{I})$ , then  $\mathbf{Y}_t(i \notin \{j, n\}, j \neq n) = e^{-\sigma_t p_m} \frac{1}{n-1} (1 - e^{-(1-p_m)\sigma_t})$ ,  $\mathbf{Y}_t(i \neq n, i \neq n) = e^{-\sigma_t p_m} (1 - \frac{n-2}{n-1} (1 - e^{-(1-p_m)\sigma_t}))$ ,  $\mathbf{Y}_t(n, j \neq n) = 1 - e^{-\sigma_t p_m}$ ,  $\mathbf{Y}_t(i \neq n, n) = 0$ , and  $\mathbf{Y}_t(n, n) = 1$ .*

The noise schedule used is the roulette-loglinear noise  $-\frac{1}{p_m} \log(1 - (1 - \epsilon)t)$ . In the reverse process, when a token is unmasked it can still be corrected with probability directly related to  $p_m$  as shown in Appendix A.3.4. A generalization for time-evolving  $p_m$  is given in Appendix A.3.5, and the corresponding  $\mathbf{Q}_{eroulette}^{tok}(t)$  is named *eroulette*. Therein (Proposition 7), it is shown that the exponential matrix can be calculated as in the previous proposition, by substituting  $p_m$  with  $p_m(t)$ .

### 3.3 EFFICIENT ESTIMATION OF $J_2$

In this subsection we provide Proposition 6, which shows that the second and third term on the RHS of Expression (14) can be computed efficiently:

**Proposition 6.** *In the case of the roulette diffusion with roulette-loglinear noise,  $H(p_r) = 0$  and  $-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) dt = \left(1 - \frac{1-p_m}{n-1}\right) \frac{L}{p_m} (\epsilon - 1)$ . For the absorb diffusion, we have  $H(p_r) = 0$  and  $-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) dt = L(\epsilon - 1)$ . Finally, in the case of uniform diffusion,  $H(p_r) = L \log(V)$  and  $-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) dt = -\left(1 - \frac{1}{V}\right) L \int_0^1 \sigma'(t) dt$ .*

### 3.4 ESTIMATING RATIOS VIA THE CROSS ENTROPY DISCRETE DIFFUSION LOSS

For sequences  $\mathbf{x}, \mathbf{y}$  which only differ at some position  $i$ , we can write

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{y}^i | h)}{p_{t|0}(\mathbf{x}_t^i | h)} p_{0|t}^i(h | \mathbf{x}_t), \quad (18)$$

where  $p_{0|t}^i(\mathbf{x}_0^i | \mathbf{x}_t) = \int p_{0|t}(\mathbf{x}_0 | \mathbf{x}_t) d\mathbf{x}_0^1 \dots d\mathbf{x}_0^{i-1} d\mathbf{x}_0^{i+1} \dots d\mathbf{x}_0^L$ , as shown in Equation (86), Appendix A.2.2. Since conditional ratios  $\frac{p_{t|0}(\mathbf{y}^i | h)}{p_{t|0}(\mathbf{x}_t^i | h)}$  are known, we can choose to reparametrize the score as

$$s_\theta^i(\mathbf{x}_t, t)_{\mathbf{y}} = \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{y}^i | h)}{p_{t|0}(\mathbf{x}_t^i | h)} f_\theta^i(\mathbf{x}_t, t)[h], \quad (19)$$

where  $f_\theta(\mathbf{x}_t, t)$ , is a neural network, with a matrix output of shape  $L \times V$ , whose elements in each row  $i$  add to 1, that is  $\sum_h f_\theta^i(\mathbf{x}_t, t)[h] = 1$ . Intuitively, given the current perturbed sequence  $\mathbf{x}_t$ , predictor  $f_\theta^i(\mathbf{x}_t, t)[h]$  gives the probability that the pre-perturbation token at position  $i$  used to be  $h$ . In Appendix A.2.2, we explain how the loss in Expression (7) can be optimized by minimizing the following cross-entropy loss:

$$L_{ll} = -\mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\cdot | \mathbf{x}_0)} \sum_{i=1}^L w(t) \log f_\theta^i(\mathbf{x}_t, t)[\mathbf{x}_0^i]. \quad (20)$$

Thus, we can learn the ratios via  $L_{ll}$  (Campbell et al., 2022) for any type of diffusion model. This approach is analogous to the one in continuous diffusion models for language modeling used in Dieleman et al. (2022), as shown in Appendix A.2.3. In addition, in Appendix B.3.2, we provide our original motivation for the reparametrization given in (19).

When  $w(t) = 1$  we refer to the method as CEDD, while when  $w(t) = \log(e + \frac{0.3}{t})$  we use CEDD\*. Lastly, in the uniform case, models using the direct reparametrization in (19) underperform SEDD in terms of perplexity. Indeed, the model is too confident in its predictions when  $t \rightarrow 0$  as it does not benefit from the neural network regularization due to the incorporation of the true conditional ratios. Thus, inside the ratios  $\frac{p_{t|0}(\mathbf{y}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)}$  we rescale  $\sigma_t^\theta < 0.0015$ , by setting  $\sigma_t^\theta = 0.0015$ . We also perform rescaling in the case of roulette diffusion dynamics, as when time is close to 0 most tokens are unmasked. The same strategy is applied, as before, only to ratios corresponding to unmasked token, by rescaling  $\sigma_t^\theta$  when  $\sigma_t^\theta < 0.5$  as follows:  $\sigma_t^{scaled,\theta} = \log(1.1\sigma_t^\theta + 1.1)$ . Naturally, for the sake of rigor, these are also the models we employ to generate samples, whose quality is measured in terms of generative perplexity. Further, in order to evade metric hacking, we take care to not modify the  $\sigma_t$  in the metrics  $J_1$  and  $J_2$  used to evaluate the models. Additional details and motivations for such rescalings are provided in the last paragraph of Appendix B.3.3.

## 4 EXPERIMENTS

We now empirically validate the approaches and theoretical contributions presented in the previous section. In Subsection 4.1, we compare the generative perplexities of models trained on OpenWebText (Gokaslan & Cohen, 2019) with SEDD and SEDD scaled (SEDDs, see Appendix B.3.2 for details) versus those trained with CEDD and CEDD\*. We keep all other variables unchanged for a fair comparison, finding that CEDD outperforms SEDD in all cases. The tests are conducted for the absorb, uniform and roulette diffusion dynamics. In Subsection 4.2, we evaluate the perplexity of the models, by calculating the upper bound on 5 different datasets, namely: 1BW, LAMBADA, PTB, Wikitext2 and Wikitext103 (Chelba et al., 2013; Paperno et al., 2016; Marcus et al., 1993; Merity et al., 2016). For the sake of reproducing the results of previous work we use  $J_1$  as a metric. Finally, having computed the results using  $J_1$ , in Subsection 4.3 we also re-evaluate the models via  $J_2$ . Our findings suggest that the bound provided by  $J_2$  is slightly tighter. In the last subsection we compare SEDD and CEDD\* on a spell-checking task.

Our model employs the transformer architecture as described by Lou et al. (2024), with no modifications; more details can be found in Appendix B.2. The algorithms for training via SEDD, and sampling unconditionally and conditionally can be found in the Appendix of Lou et al. (2024). On the other hand, we provide the algorithm for training using CEDD in Algorithm 1, Appendix B.1. In all cases, samples are generated using tau-leaping (Gillespie, 2001; Campbell et al., 2022), which performs an update at each position simultaneously for each reverse time step. All models were trained for 400k parameter updates unless stated otherwise.

### 4.1 GENERATIVE PERPLEXITY COMPARISONS

We compare the generative perplexities (GenPerp) of identical networks trained using SEDD, SEDD scaled (SEDDs), CEDD and CEDD\*. To evaluate the generative perplexity of a model, we generate samples from that model, and use a GPT-2 large model in order to assess the likelihood of the generated samples. However, this metric can be unreliable, as models such as GPT-2 large are not perfect themselves, and they tend to assign high probability to some unlikely sequences, such as those that contain repetitive tokens. Such biased samples can be generated by increasing the step size, while maintaining the number of reverse steps. To ensure a fair comparison, we evade such approaches. Furthermore, recently Zheng et al. (2024) showed that the sampling procedure in Lou et al. (2024) suffers from numerical precision issues. To address this, they proposed fixing the categorical sampling to 64-bit floating-point precision—a strategy that we have also adopted. In Appendix B.6, we also provide perplexity results as evaluated by LLama 3.1 8B (Dubey et al., 2024).

As an initial experiment, we perform a grid search to find the best  $p_m$  for the roulette transition-rate matrix. Out of the 4 values, 0.95, 0.65, 0.35, 0.05, we found  $p_m = 0.95$  performs best, but we expect that the optimum is reached when  $p_m = 1$ . The results can be found in Appendix B.6, Table 5. Therefore, in what follows, when the roulette case is concerned, it is implied that  $p_m = 0.95$ . In Table 1, we provide the experimental results in the case that the sequence length is 128 and the number of reverse steps is 1024, while the sampler used was the analytic one. Generated samples can be found in Appendix C.1.

## 4.2 PERPLEXITY COMPARISONS

In this subsection, we compare the performance of identical networks trained using SEDD, SEDDs, CEDD and CEDD\*, in terms of  $J_1$ . We do not shuffle the test set. Results for sequence length  $L = 128$  are provided in Table 1. Plots illustrating cumulative performance trajectories across the testing sets are provided in Appendix B.6.

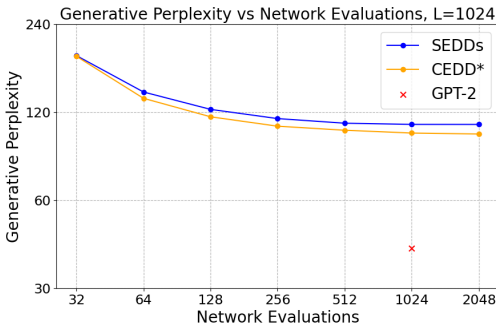


Figure 1: Scaling of Generative Perplexity vs sampling steps for SEDDs (loaded) and CEDD\* absorb.

SEDDs	<b>Seeing them is a treat</b> , because I'm coaching. It's a full week for me. "On Saturday night, Gonzalez spoke to media and said he thinks he could work out the fine fit to make the made
CEDD*	<b>Seeing them is a treat</b> , though it was only in her presence. Just being able to think about all the soldiers out here helped make the place even look like their countries and the rest of
GPT-2	<b>Seeing them is a treat</b> , and ultimately all the practical work should be well underway for everyone to enjoy. All Milan, Italy-based staff members will be given access to the dev portal

Figure 2: Filtered samples from SEDDs and CEDD\* absorb, L=1024. The conditional part is highlighted in bold.

Table 1: Results comparing SEDD, SEDD scaled (SEDDs), CEDD and CEDD\*. Lower is better.

Model (L=128)	GenPerp	LAMBADA	WikiText2	PTB	WikiText103	1BW
SEDD Absorb	172.35	70.07	75.20	240.43	74.79	88.99
SEDDs Absorb	166.35	67.05	69.37	208.69	69.17	83.87
CEDD Absorb	148.21	65.18	65.66	199.69	65.62	79.83
CEDD* Absorb	<b>143.86</b>	<b>64.60</b>	<b>65.04</b>	<b>192.99</b>	<b>64.69</b>	<b>79.81</b>
SEDD Roulette	178.94	72.07	80.13	230.74	79.68	93.45
SEDDs Roulette	172.93	69.10	74.38	<b>209.12</b>	74.16	88.02
CEDD Roulette	167.67	69.77	72.91	227.16	72.49	<b>86.55</b>
CEDD* Roulette	<b>158.56</b>	<b>67.84</b>	<b>70.54</b>	216.91	<b>70.18</b>	86.76
SEDD Uniform	169.66	80.74	91.79	252.81	91.40	102.75
SEDDs Uniform	163.88	81.13	89.21	<b>228.37</b>	88.56	100.80
CEDD Uniform	<b>161.84</b>	<b>80.27</b>	<b>87.91</b>	279.65	<b>87.46</b>	<b>99.34</b>
CEDD* Uniform	175.42	82.54	89.68	289.09	88.90	106.32
DFM $k_t = t$	<b>145.48</b>	<b>71.90</b>	<b>71.20</b>	221.15	<b>70.84</b>	<b>82.63</b>
DFM $k_t = t^2$	152.70	72.31	72.87	<b>215.30</b>	72.55	85.82

We also compare our approach against models trained utilizing Discrete Flow matching (DFM) (Gat et al., 2024). In Table 1, we present results when comparing against flows with convex interpolants, where we chose schedules  $k_t = t$ , as in Campbell et al. (2024), as well as  $k_t = t^2$ . The perplexity bound for these models is calculated using Expression (24) in Haxholli et al. (2024).

It can be seen that CEDD\* absorb performs best overall, thus we compare this model, against CEDDT, that is, CEDD with the scaling loss used in the SOTA discrete diffusion model (Sahoo et al., 2024). Interestingly, our scaling CEDD\* outperforms that of CEDDT, despite its theoretical support with regards to the score entropy loss. The results can be found in Table 2.

We also train 3 absorb models, namely SEDDs, CEDD\*, CEDDT, as well as GPT2, with a sequence length of 1024. Results are provided in Table 2, where it can be seen that overall GPT-2 performs best. The gaps between SEDDs, CEDD\*, CEDDT are reduced, likely since by seeing more tokens they all approach their optimal performance. However, models trained with CEDD/CEDD\* converge faster to the optimum in terms of number of parameter updates. In Appendix B.6, Table 7 and Figure 8, we show the difference in performance between CEDD\* and SEDDs absorb during



and after training for 20k parameter updates. In addition, training with CEDD (and its variants) is roughly 15% faster per iteration, due to the simplified loss function. Furthermore, by incorporating the  $f'_\theta(\mathbf{x}_t, t)$  to  $s'_\theta(\mathbf{x}_t, t)$  scaling in the timestep, absorb models trained with CEDD/CEDD\* generate sequences 2% faster than those trained with SEDDs.

Table 2: Results comparing SEDDs (retrained), CEDD\*, CEDDT and GPT-2 (retrained) in terms of generative perplexity, and perplexity on 5 test sets. Number of generation steps is 1024.

Model (Absorb)	GenPerp	LAMBADA	WikiText2	PTB	WikiText103	1BW
SEDDs L=128	166.35	67.05	69.37	208.69	69.17	83.87
CEDD* L=128	<b>143.86</b>	<b>64.60</b>	<b>65.04</b>	<b>192.99</b>	<b>64.69</b>	<b>79.81</b>
CEDDT L=128	154.04	68.24	68.61	204.76	68.10	81.81
SEDDs L=1024	105.27	<b>52.18</b>	42.02	117.00	41.83	80.79
CEDD* L=1024	<b>101.83</b>	52.70	<b>41.57</b>	<b>115.99*</b>	<b>41.31</b>	<b>77.96</b>
CEDDT L=1024	108.88	53.20	42.24	121.05	42.07	78.10
D3PM L=1024	-	93.47	77.28	200.82	75.16	138.92
PLAID L=1024	-	57.28	51.80	142.60	50.86	91.12
GPT-2 L=1024	41.02	<b>49.02*</b>	<b>37.68*</b>	134.13	<b>37.55*</b>	<b>58.92*</b>

#### 4.3 COMPARING THE TWO UPPER BOUNDS

Lastly, we compare the two upper bounds  $J_1$  and  $J_2$ . The bound  $J_2$  is shown empirically to be tighter than  $J_1$ , supporting the importance of Theorem 4 and Proposition 6. We estimate the integral with respect to time in both  $J_1$  and  $J_2$  by randomly sampling time points from a uniform distribution in the interval  $(e^{-4}, 1 - e^{-4})$ , and the averaging the loss. The procedure of comparing these bounds is explained in detail in Appendix B.4. We present the results for several models in Table 3, while we show the testing plots for CEDD\* absorb (L=1024) in Figure 3. The proposed bound  $J_2$  gives a lower bound in every single case, as it can also be seen in Table 6, Appendix B.6, where we provide the equivalent of Table 1 when  $J_2$  is utilized.

Table 3: Results comparing  $J_1$  and  $J_2$  for the best performing models of each category.

Model/L/Perplexity-Bound	LAMBADA	WikiText2	PTB	WikiText103	1BW
SEDDs absorb/1024/ $\exp(J_1)$	52.18	42.02	117.00	41.83	80.79
SEDDs absorb/1024/ $\exp(J_2)$	<b>51.78</b>	<b>41.76</b>	<b>115.97</b>	<b>41.51</b>	<b>80.53</b>
CEDD* absorb/1024/ $\exp(J_1)$	52.70	41.57	115.99	41.31	77.96
CEDD* absorb/1024/ $\exp(J_2)$	<b>52.10</b>	<b>41.18</b>	<b>115.03</b>	<b>40.98</b>	<b>77.28</b>
CEDD* absorb/128/ $\exp(J_1)$	64.60	65.04	192.99	64.69	79.81
CEDD* absorb/128/ $\exp(J_2)$	<b>64.11</b>	<b>64.54</b>	<b>191.38</b>	<b>64.30</b>	<b>79.17</b>
CEDD* roulette/128/ $\exp(J_1)$	67.84	70.54	216.91	70.18	86.76
CEDD* roulette/128/ $\exp(J_2)$	<b>67.27</b>	<b>69.61</b>	<b>213.90</b>	<b>69.45</b>	<b>85.64</b>
CEDD uniform/128/ $\exp(J_1)$	80.27	87.91	279.65	87.46	99.34
CEDD uniform/128/ $\exp(J_2)$	<b>79.46</b>	<b>86.82</b>	<b>276.61</b>	<b>86.52</b>	<b>98.44</b>

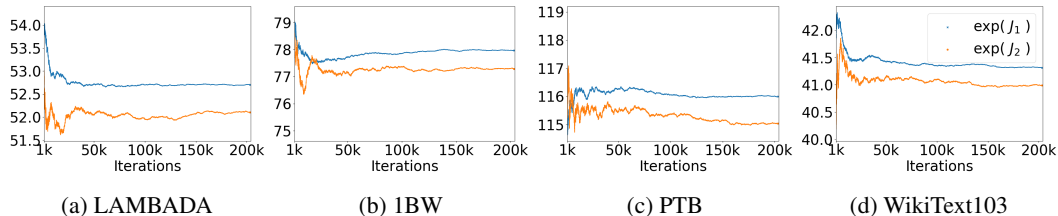


Figure 3: Upper bounds  $J_1$  and  $J_2$  of CEDD\* absorb L=1024 for different testing sets.

#### 4.4 SPELLING CORRECTION

We evaluate our uniform and roulette diffusion models on a character-level unsupervised spell-checking task (Hoogeboom et al., 2021). We train CEDD\* and SEDD models on 'War and Peace' (3.3M tokens), and test on 'Crime and Punishment' (CAP) and 'Pride and Prejudice' (PAP). The test set is contaminated with mistakes (5% of characters), and the model chooses the most likely correction. Results when training for 25k/50k iterations are provided in Table 4 and Appendix B.6.

Table 4: Correction accuracy percentages. 25k training iterations, batch size of 32 and  $L = 128$ .

Model (L=128)	CEDD* Uniform	SEDD Uniform	CEDD* Roulette	SEDD Roulette
PAP	89.5	86.9	<b>89.7</b>	85.1
CAP	89.5	87.5	<b>90.3</b>	85.8

## 5 RELATED WORK AND FUTURE OUTLOOK

**The continuous diffusion** approach has demonstrated excellent performance in modeling continuous data distributions (Song et al., 2020a; 2021c; Kingma et al., 2021; Nichol & Dhariwal, 2021; Saharia et al., 2022; Ramesh et al., 2022), leading to numerous efforts to adapt it for language modeling tasks (Chen et al., 2023; Gulrajani & Hashimoto, 2024; Li et al., 2022; Dieleman et al., 2022; Strudel et al., 2022; Gong et al., 2022; Mahabadi et al., 2023). Although initial results were not competitive, recent advancements have reduced the performance gap with autoregressive models.

**Discrete diffusion models** offer an alternative approach for modeling categorical distributions like language data. Originally proposed by Hoogeboom et al. (2021); Austin et al. (2021), the framework was extended to continuous time by Campbell et al. (2022). Both strands of work employ a combination of the variational lower bound and cross entropy loss, the latter being central to our training approach and corresponding to the strategy employed by Dieleman et al. (2022) in the continuous case. The cross-entropy loss is also derived in Sahoo et al. (2024); Ou et al. (2024); Shi et al. (2024), but only for the absorb transition-rate matrix. In addition, the cross entropy loss is similar to the loss employed in Sun et al. (2022), however in their case, one conditions on the current sequence  $\mathbf{x}_t^{-i}$  without the token at position  $i$ , and attempts to maximize the likelihood of  $\mathbf{x}_t^i$ . In contrast, Lou et al. (2024) proposed modeling ratios of probabilities (Meng et al., 2022) directly using score entropy.

**The evaluation** of such models can be performed by using Theorem 3, as originally derived and used in Lou et al. (2024). Inspired by this result, we formulate and prove the discrete version of the rest of the theorems in Song et al. (2021b), which provide important information regarding the KL divergence between the data and learned distributions in CTMCs, and which justify the usage of  $J_2$ .

**Transition-rate matrices** are an essential component in CTMCs, as demonstrated by the performance discrepancy between the absorb and uniform matrices (Lou et al., 2024; Campbell et al., 2022). We introduced roulette matrix, an interpolation between the two and derived its exponential.

**Regarding future work**, numerous avenues remain open for optimization within the diffusion model framework such as the exploration of the Eroulette matrix. Additionally, the exploration of noise schedules in discrete diffusion models remains relatively nascent, as similar to Lou et al. (2024), we did not systemically explore noise schedules. In our experiments, cross entropy loss is modulated by heuristically chosen time-dependent weighting. Investigating and establishing a general optimal weighting schedule could further refine performance metrics. Finally, studying scaling laws for discrete diffusion models and establishing their practical utility in downstream tasks is crucial.

## 6 CONCLUSION

In this work, we provided three new theorems concerning the KL divergence between the data and the learned distribution, improving model evaluation through the bound presented in Theorem 4. We also introduced a new transition-rate matrix that allows for token correction after unmasking in the reverse process, and derived its exponential matrix to enable efficient training/sampling. Finally, we proposed favouring denoising cross entropy loss over score entropy for training discrete diffusion models due to the findings in the experiments we conducted.

## 540 ETHICS STATEMENT

541  
542 This paper introduces research that progresses the field of natural language generation. Beyond  
543 the pre-existing ethical concerns associated with this domain, such as e.g. bias, toxicity, and the  
544 generation of deceptive content, our methodology poses no unique risks. This is because our work  
545 is primarily theoretical and not of a magnitude that could cause any distinct issues.

## 547 REFERENCES

- 548  
549 William J Anderson. *Continuous-time Markov chains: An applications-oriented approach*. Springer,  
550 2012. doi: <https://doi.org/10.1007/978-1-4612-3038-0>.
- 551  
552 Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured  
553 denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing*  
554 *Systems*, 34:17981–17993, 2021.
- 555  
556 Andrew Campbell, Joe Benton, Valentin De Bortoli, Thomas Rainforth, George Deligiannidis, and  
557 Arnaud Doucet. A continuous time framework for discrete denoising models. *Advances in Neural*  
558 *Information Processing Systems*, 35:28266–28279, 2022.
- 559  
560 Andrew Campbell, Jason Yim, Regina Barzilay, Tom Rainforth, and Tommi Jaakkola. Generative  
561 flows on discrete state-spaces: Enabling multimodal flows with applications to protein co-design.  
562 *arXiv preprint arXiv:2402.04997*, 2024.
- 563  
564 Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony  
565 Robinson. One billion word benchmark for measuring progress in statistical language modeling.  
566 *arXiv preprint arXiv:1312.3005*, 2013.
- 567  
568 Ting Chen, Ruixiang ZHANG, and Geoffrey Hinton. Analog bits: Generating discrete data using  
569 diffusion models with self-conditioning. In *The Eleventh International Conference on Learning*  
570 *Representations*, 2023.
- 571  
572 Jacob K Christopher, Brian R Bartoldson, Bhavya Kailkhura, and Ferdinando Fioretto. Specu-  
573 lative diffusion decoding: Accelerating language generation through diffusion. *arXiv preprint*  
574 *arXiv:2408.05636*, 2024.
- 575  
576 Justin Deschenaux and Caglar Gulcehre. Promises, outlooks and challenges of diffusion language  
577 modeling. *arXiv preprint arXiv:2406.11473*, 2024.
- 578  
579 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
580 deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and  
581 Tamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the*  
582 *Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and*  
583 *Short Papers)*, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- 584  
585 Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin, Pierre H  
586 Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, et al. Continuous diffu-  
587 sion for categorical data. *arXiv preprint arXiv:2211.15089*, 2022.
- 588  
589 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
590 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
591 *arXiv preprint arXiv:2407.21783*, 2024.
- 592  
593 Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and  
594 Yaron Lipman. Discrete flow matching. *arXiv preprint arXiv:2407.15595*, 2024.
- 595  
596 Daniel T. Gillespie. Approximate accelerated stochastic simulation of chemically reacting systems.  
597 *The Journal of Chemical Physics*, 115(4):1716–1733, 07 2001. ISSN 0021-9606. doi: 10.1063/  
598 1.1378322.
- 599  
600 Aaron Gokaslan and Vanya Cohen. Openwebtext corpus, 2019.

- 594 Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and LingPeng Kong. Diffuseq: Sequence to  
595 sequence text generation with diffusion models. *arXiv preprint arXiv:2210.08933*, 2022.
- 596
- 597 Ishaan Gulrajani and Tatsunori B Hashimoto. Likelihood-based diffusion language models. *Ad-*  
598 *vances in Neural Information Processing Systems*, 36, 2024.
- 599 Etrit Haxholli, Yeti Z Gürbüz, Oğul Can, and Eli Waxman. Minibatch optimal transport and per-  
600 plexity bound estimation in discrete flow matching. *arXiv preprint arXiv:2411.00759*, 2024.
- 601
- 602 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In  
603 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Infor-*  
604 *mation Processing Systems*, volume 33, pp. 6840–6851. Curran Associates, Inc., 2020.
- 605 Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows  
606 and multinomial diffusion: Learning categorical distributions. In A. Beygelzimer, Y. Dauphin,  
607 P. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*,  
608 2021.
- 609 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-  
610 based generative models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun  
611 Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- 612
- 613 F.P. Kelly. *Reversibility and Stochastic Networks*. Probability and Statistics Series. Wiley, 1979.
- 614 Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. On density estimation with diffu-  
615 sion models. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), *Advances*  
616 *in Neural Information Processing Systems*, 2021.
- 617
- 618 Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-  
619 lm improves controllable text generation. *Advances in Neural Information Processing Systems*,  
620 35:4328–4343, 2022.
- 621 Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios  
622 of the data distribution. In *Forty-first International Conference on Machine Learning*, 2024.
- 623
- 624 Rabeeh Karimi Mahabadi, Hamish Ivison, Jaesung Tae, James Henderson, Iz Beltagy, Matthew E  
625 Peters, and Arman Cohan. Tess: Text-to-text self-conditioned simplex diffusion. *arXiv preprint*  
626 *arXiv:2305.08379*, 2023.
- 627 Mitchell Marcus, Beatrice Santorini, and Mary Marcinkiewicz. Building a large annotated corpus  
628 of English : the Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993.
- 629
- 630 Chenlin Meng, Kristy Choi, Jiaming Song, and Stefano Ermon. Concrete score matching: General-  
631 ized score matching for discrete data. *Advances in Neural Information Processing Systems*, 35:  
632 34532–34545, 2022.
- 633
- 634 Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture  
635 models. *arXiv preprint arXiv:1609.07843*, 2016.
- 636
- 637 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.  
638 In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on*  
*Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8162–8171.  
PMLR, 2021.
- 639
- 640 Manfred Opper and Guido Sanguinetti. Variational inference for markov jump processes. In J. Platt,  
641 D. Koller, Y. Singer, and S. Roweis (eds.), *Advances in Neural Information Processing Systems*,  
642 volume 20. Curran Associates, Inc., 2007.
- 643
- 644 Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan  
645 Li. Your absorbing discrete diffusion secretly models the conditional distributions of clean data.  
*arXiv preprint arXiv:2406.03736*, 2024.
- 646
- 647 Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi,  
Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset:  
Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*, 2016.

- 648 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*  
649 *the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 4195–4205, October  
650 2023.
- 651 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
652 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 653 Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad  
654 Norouzi. Image super-resolution via iterative refinement. *IEEE transactions on pattern anal-*  
655 *ysis and machine intelligence*, 45(4):4713–4726, 2022.
- 656 Subham Sekhar Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin T  
657 Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language  
658 models. *arXiv preprint arXiv:2406.07524*, 2024.
- 659 Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis K Titsias. Simplified and gener-  
660 alized masked diffusion for discrete data. *arXiv preprint arXiv:2406.04329*, 2024.
- 661 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised  
662 learning using nonequilibrium thermodynamics. In Francis Bach and David Blei (eds.), *Proceed-*  
663 *ings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of*  
664 *Machine Learning Research*, pp. 2256–2265, Lille, France, 2015. PMLR.
- 665 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
666 *preprint arXiv:2010.02502*, 2020a.
- 667 Yang Song, Sahaj Garg, Jiaxin Shi, and Stefano Ermon. Sliced score matching: A scalable approach  
668 to density and score estimation. In Ryan P. Adams and Vibhav Gogate (eds.), *Proceedings of The*  
669 *35th Uncertainty in Artificial Intelligence Conference*, volume 115 of *Proceedings of Machine*  
670 *Learning Research*, pp. 574–584. PMLR, 2020b.
- 671 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-  
672 based diffusion models. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman  
673 Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 1415–1428.  
674 Curran Associates, Inc., 2021a.
- 675 Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of  
676 score-based diffusion models. *Advances in neural information processing systems*, 34:1415–  
677 1428, 2021b.
- 678 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
679 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-*  
680 *ional Conference on Learning Representations*, 2021c.
- 681 Robin Strudel, Corentin Tallec, Florent Alché, Yilun Du, Yaroslav Ganin, Arthur Mensch, Will  
682 Grathwohl, Nikolay Savinov, Sander Dieleman, Laurent Sifre, et al. Self-conditioned embedding  
683 diffusion for text generation. *arXiv preprint arXiv:2211.04236*, 2022.
- 684 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: En-  
685 hanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 686 Y.M. Suhov and M. Kelbert. *Probability and Statistics by Example: Volume 2, Markov Chains:*  
687 *A Primer in Random Processes and Their Applications*. Probability and Statistics by Example.  
688 Cambridge University Press, 2008. ISBN 9780521847674.
- 689 Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time  
690 discrete diffusion models. *arXiv preprint arXiv:2211.16750*, 2022.
- 691 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
692 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*  
693 *tion processing systems*, 30, 2017.
- 694 Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked  
695 diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical  
696 sampling. *arXiv preprint arXiv:2409.02908*, 2024.

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

---

**APPENDIX**


---

A.1	KL Divergence Theorems . . . . .	15
A.2	Cross Entropy Discrete Diffusion (CEDD) . . . . .	20
A.2.1	Rederiving the Denosing Score Entropy Loss . . . . .	20
A.2.2	Deriving Cross-Entropy from Ratio Matching . . . . .	20
A.2.3	Relation to Cross Entropy in Continuous Diffusion . . . . .	22
A.3	Roulette Discrete Diffusion . . . . .	23
A.3.1	Deriving the exponential matrix of the absorb diffusion . . . . .	23
A.3.2	Deriving the exponential of the uniform diffusion . . . . .	23
A.3.3	Deriving the exponential of the roulette diffusion . . . . .	24
A.3.4	Controlling the number of the corrected tokens in the reverse process . . . . .	25
A.3.5	Time Evolving Roulette . . . . .	26
B.1	Algorithms . . . . .	27
B.2	Network Architecture and Hyper-parameters . . . . .	27
B.3	Efficient Implementation in Practice . . . . .	27
B.3.1	Efficient Estimation of the Score in Practice . . . . .	27
B.3.2	SEDD Scaling . . . . .	29
B.3.3	Score Estimation Through CEDD in Practice . . . . .	31
B.3.4	Analytic Sampling in Practice . . . . .	32
B.4	Procedure for generating the plots in Figure 4 . . . . .	33
B.5	Per sequence estimation of $J_1$ . . . . .	34
B.6	Additional results . . . . .	34
C.1	Generated text . . . . .	37
D.1	Discrete-Time Markov Chains Over Finite-State Spaces . . . . .	49
D.2	Continuous-Time Markov Chains Over Finite-State Spaces . . . . .	49

---

## A THEORETICAL RESULTS

### A.1 KL DIVERGENCE THEOREMS

**Lemma 1.** Let  $q_1^*$  and  $\hat{q}_1$  denote two initial distributions and define two Continuous-Time Markov Chains (CTMCs) running from time 1 to 0, with transition-rate matrices  $\mathbf{Q}_t^*(y, x) = \mathbf{Q}_t(x, y)r(x, y)$  and  $\hat{\mathbf{Q}}_t(y, x) = \mathbf{Q}_t(x, y)\hat{r}(x, y)$ , where  $r(x, y)$  and  $\hat{r}(x, y)$  are chosen such that they ensure  $\mathbf{Q}_t^*$  and  $\hat{\mathbf{Q}}_t$  are indeed transition-rate matrices. When the first process (with matrix  $\mathbf{Q}_t^*(y, x)$ ) is applied to initial distribution  $q_1^*$  and the second (with matrix  $\hat{\mathbf{Q}}_t(y, x)$ ) to initial distribution  $\hat{q}_1$ , they define distributions  $q_t^*$  and  $\hat{q}_t$  at time  $t$ , and  $q_0^*$ ,  $\hat{q}_0$  at the end of the diffusion, for which

$$D_{KL}(q_0^*||\hat{q}_0) \leq \int_0^1 \mathbb{E}_{x_t \sim q_t^*} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \ell(r(x_t, y), \hat{r}(x_t, y)) dt + D_{KL}(q_1^*||\hat{q}_1), \quad (21)$$

where  $\ell(a, b) = (b - a \log b + K(a))$  and  $K(a) = a(\log a - 1)$ .

*Proof.* From the data processing inequality applied to the reverse of the processes defined in the Lemma, we have  $D_{KL}(q_0^*||\hat{q}_0) \leq D_{KL}(q^*||\hat{q})$  where  $q^*$  is the measure over the space of paths generated by the first process, while  $\hat{q}$  is the measure over the space of paths generated by second one. Using the expression for the KL divergence between path measures of two CTMCs found in (Oppèr & Sanguinetti, 2007, Section 2.1), and by substituting  $q_t(x) := q^*(x_t)$ ,  $g(x'|x) = \mathbf{Q}_t^*(y, x_t)$  and  $f(x'|x) = \hat{\mathbf{Q}}_t(y, x_t)$ , we have:

$$D_{KL}(q^*||\hat{q}) = \int_0^T \mathbb{E}_{x_t \sim q_t^*} \sum_{y \neq x_t} \left( \mathbf{Q}_t^*(y, x_t) \log \frac{\mathbf{Q}_t^*(y, x_t)}{\hat{\mathbf{Q}}_t(y, x_t)} + \hat{\mathbf{Q}}_t(y, x_t) - \mathbf{Q}_t^*(y, x_t) \right) dt + D_{KL}(q_1^*||\hat{q}_1). \quad (22)$$

Now we focus on the expression inside the sum:

$$\mathbf{Q}_t^*(y, x_t) \log \frac{\mathbf{Q}_t^*(y, x_t)}{\hat{\mathbf{Q}}_t(y, x_t)} + \hat{\mathbf{Q}}_t(y, x_t) - \mathbf{Q}_t^*(y, x_t) = \quad (23)$$

$$\mathbf{Q}_t(x_t, y)r(x_t, y) \log \frac{\mathbf{Q}_t(x_t, y)r(x_t, y)}{\mathbf{Q}_t(x_t, y)\hat{r}(x_t, y)} + \mathbf{Q}_t(x_t, y)\hat{r}(x_t, y) - \mathbf{Q}_t(x_t, y)r(x_t, y) = \quad (24)$$

$$\mathbf{Q}_t(x_t, y)r(x_t, y) (\log r(x_t, y) - \log \hat{r}(x_t, y)) + \mathbf{Q}_t(x_t, y)\hat{r}(x_t, y) - \mathbf{Q}_t(x_t, y)r(x_t, y) = \quad (25)$$

$$\mathbf{Q}_t(x_t, y)r(x_t, y) (\log r(x_t, y) - 1) - \mathbf{Q}_t(x_t, y)r(x_t, y) \log \hat{r}(x_t, y) + \mathbf{Q}_t(x_t, y)\hat{r}(x_t, y) = \quad (26)$$

$$\mathbf{Q}_t(x_t, y) [K(r(x_t, y)) + (\hat{r}(x_t, y) - r(x_t, y) \log \hat{r}(x_t, y))] = \mathbf{Q}_t(x_t, y) \ell(r(x_t, y), \hat{r}(x_t, y)), \quad (27)$$

which concludes the proof. We highlight that the CTMCs defined in the Lemma are completely arbitrary and not necessarily related to one-another, as the choices of  $r(x, y)$  and  $\hat{r}(x, y)$  can completely overwrite the matrix  $\mathbf{Q}_t(x, y)$ .  $\square$

**Theorem 1.** Define a CTMC with transition matrix  $\mathbf{Q}_t$  that runs from time 0 to 1. The true reverse process defines a probability evolution  $p_t$  from  $p_1$  to the data distribution  $p_0$ , while the learned reverse process induces the evolution  $p_t^\theta$  from the reference distribution  $p_1^\theta = p_r$  to the approximation of the data distribution  $p_0^\theta$ . In this setting, the following KL divergence bound holds

$$D_{KL}(p_0||p_0^\theta) \leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \ell \left( \frac{p_t(y)}{p_t(x_t)}, s_\theta(x_t, t)_y \right) dt + D_{KL}(p_1||p_r). \quad (28)$$

*Proof.* In Lemma 1, we set  $r(x, y) = \frac{p_t(y)}{p_t(x)}$  and  $\hat{r}(x_t, y) = s_\theta(x_t, t)_y$ . This implies that  $q_t^* = p_t$  and  $\hat{q}_t = p_t^\theta$ . In particular,  $q_0^* = p_0$ ,  $\hat{q}_0 = p_0^\theta$  and  $q_1^* = p_1$ ,  $\hat{q}_1 = p_1^\theta := p_r$ . Plugging everything in Expression (21) gives the desired result.  $\square$

**Theorem 3.** Let  $p_0^\theta$  denote the learned distribution from which the reverse process samples. The negative log-likelihood of a state  $x_0$  being sampled by the reverse process can be bounded from above as follows,

$$-\log p_0^\theta(x_0) \leq \int_0^1 \mathbb{E}_{x_t \sim p_{t|0}(\cdot|x_0)} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \ell \left( \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)}, s_\theta(x_t, t)_y \right) dt + D_{KL}(p_1 || p_r). \quad (29)$$

*Proof.* In Lemma 1 we set  $r(x, y) = \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)}$  and  $\hat{r}(x_t, y) = s_\theta(x_t, t)_y$ , where  $p_{t|0}$  is the probability over states at time  $t$  determined by a CTMC with transition-rate matrix  $\mathbf{Q}_t$  applied to initial distribution  $\delta(x = x_0)$ . As such, we have  $q_t^* = p_{t|0}$  and  $\hat{q}_t = p_t^\theta$ . In particular,  $q_0^* = \delta(x = x_0)$ ,  $\hat{q}_0 = p_0^\theta$  and  $q_1^* = p_{1|0}$ ,  $\hat{q}_1 = p_1^\theta := p_r$ . Plugging everything in Expression (21):

$$D_{KL}(\delta(x = x_0) || p_0^\theta) \leq \int_0^1 \mathbb{E}_{x_t \sim p_{t|0}(\cdot|x_0)} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \left[ K \left( \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)} \right) + \left( s_\theta(x_t, t)_y - \frac{p_{t|0}(y|x_0)}{p_{t|0}(x_t|x_0)} \log s_\theta(x_t, t)_y \right) \right] dt + D_{KL}(p_{1|0}(x_0) || p_r).$$

This concludes the proof as  $D_{KL}(\delta(x = x_0) || p_0^\theta) = -\log(p_0^\theta(x_0))$ .  $\square$

**Theorem 2.** Denote the intermediate distributions at time  $t$  determined by the true reverse process, and by the learned reverse process with  $p_t$  and  $p_t^\theta$ , respectively. We can write the entropy of the data distribution  $H(p_0)$  as

$$H(p_0) = H(p_1) - \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) K \left( \frac{p_t(y)}{p_t(x_t)} \right) dt. \quad (30)$$

In addition, if the learned ratios  $s_\theta(x_t, t)_y$  equal  $\frac{p_t^\theta(y)}{p_t^\theta(x_t)}$  and  $p_1 = p_1^\theta := p_r$ , then the inequality in Theorem 1, becomes an equality.

*Proof.* The cross entropy between the true data distribution, and the modeled data distribution satisfies the following:

$$H(p_0, p_0^\theta) - H(p_1, p_1^\theta) = \int_1^0 \frac{\partial}{\partial t} H(p_t, p_t^\theta) dt \quad (31)$$

To keep things clear we focus on  $\frac{\partial}{\partial t} H(p_t, p_t^\theta)$ .

$$\frac{\partial}{\partial t} H(p_t, p_t^\theta) = \frac{\partial}{\partial t} \int -p_t(x_t) \log(p_t^\theta(x_t)) dx_t = \int -\frac{\partial}{\partial t} [p_t(x_t) \log(p_t^\theta(x_t))] dx_t \quad (32)$$

$$= - \left( \int \frac{\partial p_t(x_t)}{\partial t} \log(p_t^\theta(x_t)) dx_t + \int \frac{p_t(x_t)}{p_t^\theta(x_t)} \frac{\partial p_t^\theta(x_t)}{\partial t} dx_t \right). \quad (33)$$

We can use the Kolmogorov forward equations (Equation 3), that is,  $\sum_y \mathbf{Q}_t(x_t, y) p_t(y) = \frac{\partial p_t(x_t)}{\partial t}$  and  $\frac{\partial p_t^\theta(x_t)}{\partial t} = \sum_y \mathbf{Q}_t(x_t, y) p_t^\theta(y)$  to get

$$= - \left( \int \sum_y \mathbf{Q}_t(x_t, y) p_t(y) \log(p_t^\theta(x_t)) dx_t + \int \frac{p_t(x_t)}{p_t^\theta(x_t)} \sum_y \mathbf{Q}_t(x_t, y) p_t^\theta(y) dx_t \right) \quad (34)$$

$$= - \left( \int \sum_y \mathbf{Q}_t(x_t, y) \frac{p_t(y)}{p_t(x_t)} \log(p_t^\theta(x_t)) p(x_t) dx_t + \int p_t(x_t) \sum_y \mathbf{Q}_t(x_t, y) \frac{p_t^\theta(y)}{p_t^\theta(x_t)} dx_t \right) \quad (35)$$

$$= - \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( \frac{p_t(y)}{p_t(x_t)} \log(p_t^\theta(x_t)) + \frac{p_t^\theta(y)}{p_t^\theta(x_t)} \right) \quad (36)$$



$$= -\mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( \frac{p_t^\theta(y)}{p_t^\theta(x_t)} - \frac{p_t(y)}{p_t(x_t)} \log \frac{p_t^\theta(y)}{p_t^\theta(x_t)} \right) \quad (37)$$

$$+ \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left[ \frac{p_t(y)}{p_t(x_t)} \log p_t^\theta(y) \right]. \quad (38)$$

The second term above  $\mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left[ \frac{p_t(y)}{p_t(x_t)} \log p_t^\theta(y) \right]$  is zero. Indeed,

$$\mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left[ \frac{p_t(y)}{p_t(x_t)} \log p_t^\theta(y) \right] = \sum_{x_t} \sum_y \mathbf{Q}_t(x_t, y) p_t(x_t) \left[ \frac{p_t(y)}{p_t(x_t)} \log p_t^\theta(y) \right] \quad (39)$$

$$= \sum_{x_t} \sum_y \mathbf{Q}_t(x_t, y) p_t(y) \log(p_t^\theta(y)) = \sum_y \left( \sum_{x_t} \mathbf{Q}_t(x_t, y) \right) p_t(y) \log(p_t^\theta(y)) = 0, \quad (40)$$

since for any column of a transition-rate matrix the sum of that columns elements is 0, therefore  $\sum_x \mathbf{Q}_t(x_t, y) = 0$ . In total we have:

$$H(p_0, p_0^\theta) - H(p_1, p_1^\theta) = \int_1^0 -\mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( \frac{p_t^\theta(y)}{p_t^\theta(x_t)} - \frac{p_t(y)}{p_t(x_t)} \log \frac{p_t^\theta(y)}{p_t^\theta(x_t)} \right) dt \quad (41)$$

$$= \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( \frac{p_t^\theta(y)}{p_t^\theta(x_t)} - \frac{p_t(y)}{p_t(x_t)} \log \frac{p_t^\theta(y)}{p_t^\theta(x_t)} \right) dt. \quad (42)$$

This allows the derivation of the expression for entropy of the data distribution, by taking  $p_t^\theta = p_t$

$$H(p_0) = H(p_1) + \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( \frac{p_t(y)}{p_t(x_t)} - \frac{p_t(y)}{p_t(x_t)} \log \frac{p_t(y)}{p_t(x_t)} \right) dt, \quad (43)$$

thus

$$H(p_0) = H(p_1) - \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) K \left( \frac{p_t(y)}{p_t(x_t)} \right) dt. \quad (44)$$

If we assume  $\frac{p_t^\theta(y)}{p_t^\theta(x_t)} = s_\theta(x_t, t)_y$  and  $p_1 = p_1^\theta := p_b$ , from Equation (42) we get

$$H(p_0, p_0^\theta) - H(p_1) = \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log(s_\theta(x_t, t)_y) \right) dt, \quad (45)$$

and furthermore, from Equation (44) it is trivial that

$$-H(p_0) + H(p_1) = \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) K \left( \frac{p_t(y)}{p_t(x_t)} \right) dt, \quad (46)$$

thus adding this last equation and Equation (45) we get

$$D_{KL}(p_0 || p_0^\theta) \quad (47)$$

$$= \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left[ K \left( \frac{p_t(y)}{p_t(x_t)} \right) + \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log s_\theta(x_t, t)_y \right) \right] dt. \quad (48)$$

Since we know that  $\frac{p_t(x_t)}{p_t(x_t)} = 1$ , and  $s_\theta(x_t, t)_{x_t} = \frac{p_t(x_t)}{p_t(x_t)} = 1$ , we have

$$\left[ K \left( \frac{p_t(x_t)}{p_t(x_t)} \right) + \left( s_\theta(x_t, t)_{x_t} - \frac{p_t(x_t)}{p_t(x_t)} \log s_\theta(x_t, t)_{x_t} \right) \right] = -1 + 1 = 0 \quad (49)$$

so we can change the sums  $\sum_y$  to simply  $\sum_{y \neq x_t}$ , as follows

$$D_{KL}(p_0 || p_0^\theta) \quad (50)$$

$$= \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \left[ K \left( \frac{p_t(y)}{p_t(x_t)} \right) + \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log s_\theta(x_t, t)_y \right) \right] dt. \quad (51)$$

Finally, since  $p_1 = p_1^\theta := p_b$ , then  $D_{KL}(p_1 || p_r) = 0$ , therefore we can add it to the right side, finishing the proof.  $\square$

We remark that as  $s_\theta(x_t, t)_y \rightarrow \frac{p_t(y)}{p_t(x_t)}$ , then  $p_t(i)^\theta \rightarrow p_t(i)$ , and thus  $\frac{p_t^\theta(y)}{p_t^\theta(x_t)} \rightarrow \frac{p_t(y)}{p_t(x_t)}$ , as such  $s_\theta(x_t, t)_y \rightarrow \frac{p_t^\theta(y)}{p_t^\theta(x_t)}$ . Therefore, the bound in Theorem 1 becomes tighter as the model improves.

**Theorem 4.** *Under the conditions stated in Theorem 1, the following inequality for the cross entropy between the data and the learned distribution holds:*

$$\begin{aligned} H(p_0, p_0^\theta) &\leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \bar{\ell} \left( \frac{p_t(y)}{p_t(x_t)}, s_\theta(x_t, t)_y \right) dt \\ &\quad - \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(y, x_t) dt + H(p_1, p_r), \end{aligned} \quad (52)$$

where  $\bar{\ell}(a, b) = (b - a \log b)$ .

*Proof.* From Theorem 1, we have

$$\begin{aligned} &D_{KL}(p_0 \| p_0^\theta) - D_{KL}(p_1 \| p_r) \\ &\leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \left[ K \left( \frac{p_t(y)}{p_t(x_t)} \right) + \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log s_\theta(x_t, t)_y \right) \right] dt. \end{aligned} \quad (53)$$

Since we know that  $\frac{p_t(x_t)}{p_t(x_t)} = 1$ , we can manually set  $s_\theta(x_t, t)_{x_t} = 1$ , and we get

$$\left[ K \left( \frac{p_t(x_t)}{p_t(x_t)} \right) + \left( s_\theta(x_t, t)_{x_t} - \frac{p_t(x_t)}{p_t(x_t)} \log s_\theta(x_t, t)_{x_t} \right) \right] = -1 + 1 = 0, \quad (54)$$

therefore

$$\begin{aligned} &H(p_0, p_0^\theta) - H(p_0) - D_{KL}(p_1 \| p_r) \\ &\leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left[ K \left( \frac{p_t(y)}{p_t(x_t)} \right) + \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log s_\theta(x_t, t)_y \right) \right] dt. \end{aligned} \quad (55)$$

from which we deduce,

$$\begin{aligned} &H(p_0, p_0^\theta) \\ &\leq H(p_0) + \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) K \left( \frac{p_t(y)}{p_t(x_t)} \right) dt \\ &\quad + \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log s_\theta(x_t, t)_y \right) dt + D_{KL}(p_1 \| p_r). \end{aligned} \quad (56)$$

Since from Theorem 2 we have that

$$H(p_0) + \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x, y) K \left( \frac{p_t(y)}{p_t(x_t)} \right) dt = H(p_1), \quad (57)$$

hence we can write

$$\begin{aligned} &H(p_0, p_0^\theta) \\ &\leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_y \mathbf{Q}_t(x_t, y) \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log s_\theta(x_t, t)_y \right) dt + H(p_1) + D_{KL}(p_1 \| p_r), \end{aligned}$$

and therefore

$$\begin{aligned} H(p_0, p_0^\theta) &\leq \int_0^1 \mathbb{E}_{x_t \sim p_t} \sum_{y \neq x_t} \mathbf{Q}_t(x_t, y) \left( s_\theta(x_t, t)_y - \frac{p_t(y)}{p_t(x_t)} \log s_\theta(x_t, t)_y \right) dt \\ &\quad + \int_0^1 \mathbb{E}_{x_t \sim p_t} \mathbf{Q}_t(x_t, x_t) dt + H(p_1, p_r). \end{aligned} \quad (58)$$

The fact that  $\mathbf{Q}_t(x_t, x_t) = -\sum_{y \neq x_t} \mathbf{Q}_t(y, x_t)$  concludes the proof.  $\square$

**Proposition 6.** *In the case of the roulette diffusion with roulette-loglinear noise,  $H(p_r) = 0$  and  $-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) dt = \left(1 - \frac{1-p_m}{n-1}\right) \frac{L}{p_m} (\epsilon - 1)$ . For the absorb diffusion, that is when  $p_m \rightarrow 1$ , we have  $-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) dt = L(\epsilon - 1)$ . Finally, for the uniform diffusion,  $H(p_r) = L \log(V)$  and  $-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) dt = -\left(1 - \frac{1}{V}\right) L \int_0^1 \sigma'(t) dt$ .*

*Proof.* In all cases,

$$-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) dt = -\mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \sum_{i=1}^L \sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{y}^i, \mathbf{x}_t^i) \quad (59)$$

$$= -\mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{y}^i, \mathbf{x}_t^i) \quad (60)$$

$$= -\mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_t^i} p_t(\mathbf{x}_t^i) \sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{y}^i, \mathbf{x}_t^i) \quad (61)$$

$$= -\mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_0^i} p_0(\mathbf{x}_0^i) \sum_{\mathbf{x}_t^i} p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i) \sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{y}^i, \mathbf{x}_t^i) \quad (62)$$

For the roulette case, if  $\mathbf{x}_t^i$  is the masked token, that is,  $\mathbf{x}_t^i = n$ , then  $\sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{y}^i, \mathbf{x}_t^i) = 0$ .

Otherwise,  $\sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{y}^i, \mathbf{x}_t^i) = \sigma'_t \left(1 - \frac{1-p_m}{n-1}\right)$ . Thus,

$$-\int_0^1 \mathbb{E}_{\mathbf{x}_t \sim p_t} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{y}, \mathbf{x}_t) = \quad (63)$$

$$= -\left(1 - \frac{1-p_m}{n-1}\right) \mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_0^i} p_0(\mathbf{x}_0^i) \sum_{\mathbf{x}_t^i \neq n} p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i) \sigma'_t = \quad (64)$$

$$= -\left(1 - \frac{1-p_m}{n-1}\right) \mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_0^i} p_0(\mathbf{x}_0^i) \sigma'_t \sum_{\mathbf{x}_t^i \neq n} p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i) \quad (65)$$

$$= -\left(1 - \frac{1-p_m}{n-1}\right) \mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_0^i} p_0(\mathbf{x}_0^i) \sigma'_t (1 - p(\mathbf{x}_t^i = n | \mathbf{x}_0^i)) \quad (66)$$

$$= -\left(1 - \frac{1-p_m}{n-1}\right) \mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_0^i} p_0(\mathbf{x}_0^i) \sigma'_t e^{-\sigma_t p_m} = -\left(1 - \frac{1-p_m}{n-1}\right) L \int_0^1 \sigma'_t e^{-\sigma_t p_m} dt \quad (67)$$

$$= \left(1 - \frac{1-p_m}{n-1}\right) \frac{L}{p_m} (e^{-\sigma_1 p_m} - e^{-\sigma_0 p_m}) \quad (68)$$

Since  $\sigma_t$  is roulette-loglinear, that is  $\sigma_t = -\frac{1}{p_m} \log(1 - (1 - \epsilon)t)$ , we get the result stated in the Proposition.

In the Uniform case, we notice that Equation (62) can be rewritten as

$$\mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_0^i} \sum_{\mathbf{x}_t^i} p(\mathbf{x}_t^i, \mathbf{x}_0^i) \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{x}_0^i) \quad (69)$$

$$= \mathbb{E}_{t \sim U(0,1)} \sum_{i=1}^L \sum_{\mathbf{x}_0^i} \sum_{\mathbf{x}_t^i} p(\mathbf{x}_t^i, \mathbf{x}_0^i) \sigma'(t) \left(\frac{1}{V} - 1\right) = -L \left(1 - \frac{1}{V}\right) \int_0^1 \sigma'(t) dt. \quad (70)$$

□

## A.2 CROSS ENTROPY DISCRETE DIFFUSION (CEDD)

### A.2.1 REDERIVING THE DENOSING SCORE ENTROPY LOSS

In what follows, we rederive the Denosing Score Entropy Loss for completeness (Lou et al., 2024). Training a model using ratio matching is performed by minimizing the error of a network that receives as input  $\mathbf{x}_t$  and tries to predict the ratios of  $\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}$ , where  $\mathbf{y}$  is a neighbour of  $\mathbf{x}_t$  with Hamming distance of 1. In other words, one tries to minimize

$$\bar{\ell} \left( \frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right), \quad (71)$$

for all  $\mathbf{x}_t$  following distribution  $p_t(\mathbf{x}_t)$ , that is:

$$\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \bar{\ell} \left( \frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right). \quad (72)$$

The main bottleneck is that the ratios  $\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}$  are unknown, as  $p_t(\mathbf{x}_t) = \int p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)p_0(\mathbf{x}_0)d\mathbf{x}_0$  and  $p_0(\mathbf{x}_0)$  is what we are trying to model in the first place. Luckily, the denoising trick can be employed:

$$\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \bar{\ell} \left( \frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right) = \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \bar{\ell} \left( \sum_{\mathbf{x}_0} p_{t|0}(\mathbf{y}|\mathbf{x}_0) \frac{p_0(\mathbf{x}_0)}{p_t(\mathbf{x}_t)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right) = \quad (73)$$

$$= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \bar{\ell} \left( \sum_{\mathbf{x}_0} \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \frac{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)p_0(\mathbf{x}_0)}{p_t(\mathbf{x}_t)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right) \quad (74)$$

$$= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \bar{\ell} \left( \sum_{\mathbf{x}_0} \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} p_{0|i}(\mathbf{x}_0|\mathbf{x}_t), s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right) \quad (75)$$

$$= \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \bar{\ell} \left( \mathbb{E}_{\mathbf{x}_0 \sim p_{0|i}(\mathbf{x}_0|\mathbf{x}_t)} \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right). \quad (76)$$

When  $\bar{\ell}$  is linear with respect to ratios, like in the case of Score Entropy (Lou et al., 2024), then we can pull the inner expectation (sum) outside and have

$$\mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \bar{\ell} \left( \frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right) = \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p_{0|i}(\mathbf{x}_0|\mathbf{x}_t)} \bar{\ell} \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right), \quad (77)$$

which of course is equal to

$$\mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \bar{\ell} \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}, s_\theta(\mathbf{x}_t, t, \mathbf{y}) \right). \quad (78)$$

### A.2.2 DERIVING CROSS-ENTROPY FROM RATIO MATCHING

In order to go from Equation (73) to (75) we used the fact that

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0} \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} p_{0|i}(\mathbf{x}_0|\mathbf{x}_t). \quad (79)$$

Now, we select position  $i$  on the sequence with length  $L$ . We want to find an expression of the ratios  $\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}$  for all sequences  $\mathbf{y}$ , which differ with  $\mathbf{x}_t$  only on position  $i$ . From above we have

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0} \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} p_{0|i}(\mathbf{x}_0|\mathbf{x}_t), \quad (80)$$

where

$$\frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} = \frac{p_{t|0} \left( (\mathbf{y}^{(0)}, \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(i)}, \dots, \mathbf{y}^{(V-1)}) | (\mathbf{x}_0^{(0)}, \mathbf{x}_0^{(1)}, \dots, \mathbf{x}_0^{(i)}, \dots, \mathbf{x}_0^{(V-1)}) \right)}{p_{t|0} \left( (\mathbf{x}_t^{(0)}, \mathbf{x}_t^{(1)}, \dots, \mathbf{x}_t^{(i)}, \dots, \mathbf{x}_t^{(V-1)}) | (\mathbf{x}_0^{(0)}, \mathbf{x}_0^{(1)}, \dots, \mathbf{x}_0^{(i)}, \dots, \mathbf{x}_0^{(V-1)}) \right)} \quad (81)$$

and due to independence between entries in the forward process we have

$$\frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} = \frac{\prod p_{t|0}(\mathbf{y}^j|\mathbf{x}_0^j)}{\prod p_{t|0}(\mathbf{x}_t^j|\mathbf{x}_0^j)} = \prod \frac{p_{t|0}(\mathbf{y}^j|\mathbf{x}_0^j)}{p_{t|0}(\mathbf{x}_t^j|\mathbf{x}_0^j)} = \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)}, \quad (82)$$

as the rest of these ratios are 1, since  $\mathbf{y}$  differs with  $\mathbf{x}_t$  only on position  $i$ . Therefore we have

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0} \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} p_{0|t}(\mathbf{x}_0|\mathbf{x}_t), \quad (83)$$

where  $p_{0|t}(\mathbf{x}_0|\mathbf{x}_t) = p_{0|t}(\mathbf{x}_0^{(0)}, \mathbf{x}_0^{(1)}, \dots, \mathbf{x}_0^{(i)}, \dots, \mathbf{x}_0^{(V-1)}|\mathbf{x}_t)$  and where the expression inside the sum  $\frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)}$ , depends only on  $\mathbf{x}_0^i$  and not on the rest of  $\mathbf{x}_0^j$  for  $j \neq i$ . Thus,

$$\sum_{\mathbf{x}_0} \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} p_{0|t}(\mathbf{x}_0|\mathbf{x}_t) = \sum_{\mathbf{x}_0^i} \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} \sum_{\mathbf{x}_0^{(0)}, \dots, \mathbf{x}_0^{(i-1)}, \mathbf{x}_0^{(i+1)}, \dots, \mathbf{x}_0^{(V-1)}} p_{0|t}(\mathbf{x}_0|\mathbf{x}_t), \quad (84)$$

which implies that

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0^i} \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} p_{0|t}^i(\mathbf{x}_0|\mathbf{x}_t), \quad (85)$$

or more clearly

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{y}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)} p_{0|t}^i(h|\mathbf{x}_t). \quad (86)$$

If we learn the  $V$  probabilities  $[p_{0|t}^i(0|\mathbf{x}_t), p_{0|t}^i(1|\mathbf{x}_t), \dots, p_{0|t}^i((V-1)|\mathbf{x}_t)]$ , since we analytically have  $\frac{p_{t|0}(\mathbf{y}^i|h^i)}{p_{t|0}(\mathbf{x}_t^i|h^i)}$  (from the matrix exponential), then we will have the ratio  $\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)}$ . We can choose another sequence  $\mathbf{z}$  which also differs only at position  $i$  from  $\mathbf{x}_t$ . Then we will have

$$\frac{p_t(\mathbf{z})}{p_t(\mathbf{x}_t)} = \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{z}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)} p_{0|t}^i(h|\mathbf{x}_t), \quad (87)$$

which highlights the fact that even though the ratios might change, the same  $[p_{0|t}^i(0|\mathbf{x}_t), p_{0|t}^i(1|\mathbf{x}_t), \dots, p_{0|t}^i((V-1)|\mathbf{x}_t)]$  appear as long as the selected position  $i$  in the sequence remains unchanged. Therefore if we learn these probabilities for a position, we will have the ratios of all neighbours that differ only in that position, that is, we have  $V$  ratios, by modeling these  $V$  probabilities.

There is nothing special about position  $i$  and we can choose to model  $V$  probabilities for each  $L$  positions. In order to do so, we define a neural network  $f_\theta(\mathbf{x}_t, t)$  whose output is  $L \times V$ , where entry  $(i, h)$  of the output (that is  $f_\theta^i(\mathbf{x}_t, t)[h]$ ), predicts  $p_{0|t}^i(h|\mathbf{x}_t)$  the probability of the entry  $i$  of  $\mathbf{x}_t$  (position  $i$  of the current sequence), having being perturbed from token with id  $h$  in the original sequence  $\mathbf{x}_0$ , given that we are at sequence  $\mathbf{x}_t$  right now. That is, the model is directly trying to predict from where each of the tokens in the current sequence came from. Therefore we will reparametrize our ratio model as follows

$$s_\theta^i(\mathbf{x}_t, t)[h] = s_\theta^i(\mathbf{x}_t, t)_h = \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{z}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)} f_\theta^i(\mathbf{x}_t, t)[h], \quad (88)$$

where  $f_\theta^i(\mathbf{x}_t, t)$  are the outputs of the softmax at the end, and where for each  $i$ , distribution  $f_\theta^i(\mathbf{x}_t, t)$  should match  $p_{0|t}^i(\cdot|\mathbf{x}_t)$ . This happens when  $D_{KL}(p_{0|t}^i(\cdot|\mathbf{x}_t) || f_\theta^i(\mathbf{x}_t, t)) = 0$ , thus we propose to minimize,

$$\sum_{i=1}^L \mathbb{E}_{t \sim U(0,1)} w(t) \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} D_{KL}(p_{0|t}^i(\cdot|\mathbf{x}_t) || f_\theta^i(\mathbf{x}_t, t)). \quad (89)$$

Since  $D_{KL}(p_{0|t}^i(\cdot|\mathbf{x}_t) || f_\theta^i(\mathbf{x}_t, t)) = H(p_{0|t}^i(\cdot|\mathbf{x}_t), f_\theta^i(\mathbf{x}_t, t)) - C = \mathbb{E}_{p_{0|t}^i(h|\mathbf{x}_t)} \log f_\theta^i(\mathbf{x}_t, t)[h] - C$  the loss function above has the same gradients with regards to network parameters as

$$-\sum_{i=1}^L \mathbb{E}_{t \sim U(0,1)} w(t) \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \mathbb{E}_{h \sim p_{0|t}^i(h|\mathbf{x}_t)} \log f_\theta^i(\mathbf{x}_t, t)[h], \quad (90)$$

1134 which is equal to

$$1135 \quad - \sum_{i=1}^L \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_t \sim p_t(\mathbf{x}_t)} \mathbb{E}_{\mathbf{x}_0 \sim p_{0|t}(\mathbf{x}_0|\mathbf{x}_t)} w(t) \log f_{\theta}^i(\mathbf{x}_t, t)[\mathbf{x}_0^i]. \quad (91)$$

1138 Thus, as in (Campbell et al., 2022), the loss we utilize in our training is  $L_{ll}$ ,

$$1140 \quad L_{ll} = -\mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \sum_{i=1}^L w(t) \log f_{\theta}^i(\mathbf{x}_t, t)[\mathbf{x}_0^i]. \quad (92)$$

1142 We can also motivate this loss from the score entropy in Equation (78). Indeed, since  $\ell$  and  $\bar{\ell}$  differ by a constant, its gradients are the same as those of

$$1145 \quad \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \ell \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}, \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{y}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)} f_{\theta}^i(\mathbf{x}_t, t)[h] \right), \quad (93)$$

1147 and as before since  $\mathbf{y}$  and  $\mathbf{x}_t$  differ at only one token we get

$$1149 \quad = \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \ell \left( \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)}, \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{y}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)} f_{\theta}^i(\mathbf{x}_t, t)[h] \right). \quad (94)$$

1151 The function  $\ell(a, b)$  is clearly minimized when  $a = b$ , in which case  $\ell(a, b) = 0$ . In our case,

$$1152 \quad \ell \left( \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)}, \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{y}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)} f_{\theta}^i(\mathbf{x}_t, t)[h] \right) = 0, \quad (95)$$

1154 for

$$1155 \quad \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} = \sum_{h \in [V]} \frac{p_{t|0}(\mathbf{y}^i|h)}{p_{t|0}(\mathbf{x}_t^i|h)} f_{\theta}^i(\mathbf{x}_t, t)[h], \quad (96)$$

1156 which happens when  $f_{\theta}^i(\mathbf{x}_t, \theta)[h \neq \mathbf{x}_0^i] = 0$  and  $f_{\theta}^i(\mathbf{x}_t, \theta)[\mathbf{x}_0^i] = 1$ , that is when our probability prediction matches the one hot encoding used in cross entropy. Thus we train our model with cross entropy, which simplifies the job of the network of learning complex ratios, as the conditional ratios now do not participate in the loss. We only add them during sampling to the weighted sum in order to predict the marginal ratios, since we have analytic expressions for the conditional ratios. Therefore, we train the model by minimizing

$$1164 \quad -\mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)} \sum_{i=1}^L w(t) \log f_{\theta}^i(\mathbf{x}_t, t)[\mathbf{x}_0^i]. \quad (97)$$

### 1166 A.2.3 RELATION TO CROSS ENTROPY IN CONTINUOUS DIFFUSION

1168 Cross entropy has also been used in continuous diffusion models applied to Language Modelling (Dieleman et al., 2022). Here, we draw parallels between the two approaches. We denote a sequence of  $L$  tokens at time  $t$  as  $\mathbf{x}_t = (\mathbf{x}_t^1, \dots, \mathbf{x}_t^L)$ , where each token is embedded in a  $D$ -dimensional space  $\mathbf{x}_t^i \in \mathbb{R}^D$ . That is, each sequence can be seen as an  $L \times D$  vector, created from concatenating the embeddings of each token. The score that generates the reverse process, is therefore a  $L \times D$  vector  $s_{\theta}(\mathbf{x}_t, t)$ , which approximates:

$$1174 \quad \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = \int \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{x}_0) p_t(\mathbf{x}_0|\mathbf{x}_t) d\mathbf{x}_0 = \quad (98)$$

$$1176 \quad \int \nabla_{\mathbf{x}_t} \log e^{-\frac{\|\mathbf{x}_t - \mathbf{x}_0\|^2}{2\sigma_t^2}} p_t(\mathbf{x}_0|\mathbf{x}_t) d\mathbf{x}_0 = \int \frac{\mathbf{x}_0 - \mathbf{x}_t}{\sigma_t^2} p_t(\mathbf{x}_0|\mathbf{x}_t) d\mathbf{x}_0 = \int \frac{\mathbf{x}_0}{\sigma_t^2} p_t(\mathbf{x}_0|\mathbf{x}_t) d\mathbf{x}_0 - \frac{\mathbf{x}_t}{\sigma_t^2}. \quad (99)$$

1179 From above, it is clear that

$$1180 \quad \nabla_{\mathbf{x}_t^i} \log p_t(\mathbf{x}_t) = \int \frac{\mathbf{x}_0^i}{\sigma_t^2} p_t(\mathbf{x}_0|\mathbf{x}_t) d\mathbf{x}_0 - \frac{\mathbf{x}_t^i}{\sigma_t^2} = \int \frac{\mathbf{x}_0^i}{\sigma_t^2} \left( \int p_t(\mathbf{x}_0|\mathbf{x}_t) d\mathbf{x}_0^{-i} \right) d\mathbf{x}_0^i - \frac{\mathbf{x}_t^i}{\sigma_t^2}, \quad (100)$$

1182 where  $\mathbf{x}_0^{-i}$  denotes all entries of  $\mathbf{x}_0$  without the ones of  $\mathbf{x}_0^i$ . Therefore

$$1184 \quad s_{\theta}(\mathbf{x}_t, t)_{(i1, i2, \dots, iD)} \approx \nabla_{\mathbf{x}_t^i} \log p_t(\mathbf{x}_t) = \int \frac{\mathbf{x}_0^i}{\sigma_t^2} p_t(\mathbf{x}_0^i|\mathbf{x}_t) d\mathbf{x}_0^i - \frac{\mathbf{x}_t^i}{\sigma_t^2}. \quad (101)$$

1186 Clearly, all that need to be learned in order to model the score for the dimensions of the token at position  $i$  are the  $V$  probabilities  $p_t(\mathbf{x}_0^i|\mathbf{x}_t)$ . For  $L$  such tokens (to model the score of the entire sequence) one simply needs to model  $L \times V$  probabilities using cross entropy as in the discrete case.

### A.3 ROULETTE DISCRETE DIFFUSION

The roulette transition-rate matrix is

$$Q_{roulette}^{tok} = \begin{pmatrix} \frac{1}{n-1}(1-p_m) - 1 & \frac{1}{n-1}(1-p_m) & \dots & \frac{1}{n-1}(1-p_m) & 0 \\ \frac{1}{n-1}(1-p_m) & \frac{1}{n-1}(1-p_m) - 1 & \dots & \frac{1}{n-1}(1-p_m) & 0 \\ \dots & \dots & \dots & \dots & \dots \\ \frac{1}{n-1}(1-p_m) & \frac{1}{n-1}(1-p_m) & \dots & \frac{1}{n-1}(1-p_m) - 1 & 0 \\ p_m & p_m & \dots & p_m & 1-1 \end{pmatrix}_{n \times n}, \quad (102)$$

where  $n-1 = V$  is the number of tokens in our vocabulary. We add a special token for the absorb (mask) state, therefore increasing the number of total states to  $n$ , with the token id  $n$  corresponding to the absorbed state. Unfortunately, we cannot derive the matrix exponential as in the case of absorb and uniform diffusion, since the corresponding matrix  $P$  of  $Q$  is not idempotent. Therefore, we will first derive the exponential matrix of the uniform and absorb diffusion manually, in order to motivate the manual derivation of the exponential matrix in the roulette case.

#### A.3.1 DERIVING THE EXPONENTIAL MATRIX OF THE ABSORB DIFFUSION

The transition matrix in the absorb case is

$$Q_{abs}^{tok} = \begin{pmatrix} 0-1 & 0 & \dots & 0 & 0 \\ 0 & 0-1 & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0-1 & 0 \\ 1 & 1 & \dots & 1 & 1-1 \end{pmatrix}_{n \times n}. \quad (103)$$

Clearly this is  $Q_{roulette}^{tok}$  when  $p_m = 1$ .

If we are at state  $i$ , the probability of moving at state  $n$  (the absorb/mask state) for a time-step of size  $\epsilon$  is  $\epsilon$ . We discretize the time interval into discrete steps that are multiples of  $\epsilon$ , that is  $[0, \tau] \rightarrow \{\epsilon, \dots, j\epsilon, \dots, \lfloor \frac{\tau}{\epsilon} \rfloor \epsilon\}$  and denote the event of jumping from  $i \neq n$  to  $n$  at time  $j\epsilon$  as  $a_j$ . The intersection of any of these two events is empty and their union is the space of all possibilities, thus the probability of being at  $n$  at time  $\tau$  is

$$p(x_\tau = n) = p(\{x_\tau = n\} \cap \Omega) = p(\{x_\tau = n\} \cap \{\cup a_j\}) = p(\cup_{j\epsilon < \tau} \{a_j\}) \quad (104)$$

$$p(x_\tau = n) = \sum_{j < \frac{\tau}{\epsilon}} p(a_j) = \epsilon + \epsilon(1-\epsilon) + \epsilon(1-\epsilon)^2 + \dots + \epsilon(1-\epsilon)^{\lfloor \frac{\tau}{\epsilon} \rfloor - 1} \quad (105)$$

$$p(x_\tau = n) = \epsilon \frac{1 - (1-\epsilon)^{\lfloor \frac{\tau}{\epsilon} \rfloor}}{1 - (1-\epsilon)} = 1 - (1-\epsilon)^{\lfloor \frac{\tau}{\epsilon} \rfloor} \quad (106)$$

Taking the limit  $\epsilon \rightarrow 0$ , we get  $p(x_\tau = n) = 1 - e^{-\tau}$ . Thus, assuming that we start at state  $i$ , we have  $p_{\tau|0}(j \notin \{n, i\} | i \neq n) = 0$ ,  $p_{\tau|0}(j = i | i \neq n) = e^{-\tau}$  and  $p_{\tau|0}(j = n | i \neq n) = 1 - e^{-\tau}$ . This defines the column  $i \neq n$  of the exponential  $e^{\tau Q^{tok}}$ . Iterating through different  $i$  we construct the entire exponential matrix, except its last column which is a simple one hot encoding at  $n$ . Setting  $\tau = \sigma_t$  finishes the derivation.

#### A.3.2 DERIVING THE EXPONENTIAL OF THE UNIFORM DIFFUSION

$$Q_{unif}^{tok} = \begin{pmatrix} \frac{1}{n-1} - 1 & \frac{1}{n-1} & \dots & \frac{1}{n-1} \\ \frac{1}{n-1} & \frac{1}{n-1} - 1 & \dots & \frac{1}{n-1} \\ \dots & \dots & \dots & \dots \\ \frac{1}{n-1} & \frac{1}{n-1} & \dots & \frac{1}{n-1} - 1 \end{pmatrix}_{V \times V} \quad (107)$$

Clearly this is  $Q_{roulette}^{tok}$  when  $p_m = 0$ , without the absorption row and column, as our vocabulary size is  $V = n-1$  since we do not have the special token. As before we discretize the time interval into discrete steps that are multiples of  $\epsilon$ , that is  $[0, \tau] \rightarrow \{\epsilon, \dots, j\epsilon, \dots, \lfloor \frac{\tau}{\epsilon} \rfloor \epsilon\}$ . We assume that initially (at time 0) we start at position  $i$ . We wish to find the probability of being at state  $k \neq i$  at time  $\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon$ . For the sake of simplicity, we abuse notation by denoting  $p_{(\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon - m\epsilon)|0}(j|i)$  as  $\bar{p}_{\tau-m\epsilon}(j|i)$ . By symmetry we know that for every  $j \neq i$  the probability  $\bar{p}_{\tau-\epsilon}(j|i)$  is the same. Thus the probability of  $\bar{p}_\tau(j|i)$  is the sum of the following components:

- probability of being at  $j \notin \{i, k\}$  at time  $\tau - \epsilon$ , i.e.  $(\bar{p}_{\tau-\epsilon}(j|i))$ , times  $n - 3 = V - 2$  such states, times probability of moving ( $\epsilon$ ), times probability of hitting  $k$  on that move  $(\frac{1}{n-1})$ .
- probability of being at  $j = i$  at time  $\tau - \epsilon$ , i.e.  $(1 - (n - 2) \cdot \bar{p}_{\tau-\epsilon}(j|i))$ , times probability of moving ( $\epsilon$ ), times probability of hitting  $k$  on that move  $(\frac{1}{n-1})$ .
- probability of being at  $j = k$  at time  $\tau - \epsilon$ , i.e.  $(\bar{p}_{\tau-\epsilon}(j|i))$ , times of staying there  $(1 - \epsilon)$ .
- probability of being at  $j = k$  at time  $\tau - \epsilon$ , i.e.  $(\bar{p}_{\tau-\epsilon}(j|i))$ , times of moving ( $\epsilon$ ), times of hitting itself on this move  $(\frac{1}{n-1})$ .

We write them down mathematically and get:

$$\bar{p}_\tau(k|i) = \bar{p}_{\tau-\epsilon}(j|i)\epsilon\frac{1}{n-1}(n-3) + (1 - (n-2)\bar{p}_{\tau-\epsilon}(j|i))\epsilon\frac{1}{n-1} + \quad (108)$$

$$+ \bar{p}_{\tau-\epsilon}(j|i)(1 - \epsilon) + \bar{p}_{\tau-\epsilon}(j|i)\frac{1}{n-1}\epsilon, \quad (109)$$

thus

$$\bar{p}_\tau(k|i) = \frac{1}{n-1}\epsilon + \bar{p}_{\tau-\epsilon}(j|i)(1 - \epsilon). \quad (110)$$

Since by definition  $k \neq i$  and  $j \neq i$ , then by symmetry we have  $\bar{p}_\tau(k|i) = \bar{p}_\tau(j|i)$ . Therefore we get the recursion:

$$p_{\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon | 0}(k|i) = \frac{1}{n-1}\epsilon + p_{\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon - \epsilon | 0}(k|i), \quad (111)$$

which if we fully develop becomes:

$$p_{\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon | 0}(k|i) = \frac{1}{n-1}\epsilon(1 + (1 - \epsilon) + (1 - \epsilon)^2 + \dots + (1 - \epsilon)^{\lfloor \frac{\tau}{\epsilon} \rfloor - 1}), \quad (112)$$

which is equal to:

$$p_{\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon | 0}(k|i) = \frac{1}{n-1}\epsilon \frac{1 - (1 - \epsilon)^{\lfloor \frac{\tau}{\epsilon} \rfloor}}{1 - (1 - \epsilon)} = \frac{1}{n-1}(1 - (1 - \epsilon)^{\lfloor \frac{\tau}{\epsilon} \rfloor}). \quad (113)$$

Taking the limit  $\epsilon \rightarrow 0$  gives  $p_{\tau|0}(k \neq i|i) = \frac{1}{n-1}(1 - e^{-\tau})$  and  $p_{\tau|0}(i|i) = 1 - \frac{n-2}{n-1}(1 - e^{-\tau})$ , which gives the  $i$ th column of the exponential matrix. Setting  $\tau = \sigma_t$  finishes the derivation.

### A.3.3 DERIVING THE EXPONENTIAL OF THE ROULETTE DIFFUSION

First, we notice that we must start at a state that is different from the absorb state (there are no masked tokens in the training data). We can split the states into two groups, the non-absorbing states, and the absorbing state. We consider the non-absorbing states as a single super-state and the absorbing state as the other super-state. From our derivations in Appendix A.3.1, we know how to derive the probability of being at the absorb super-state at time  $\tau$ . The only difference is that previously, if we moved, it would be certain we would move to the absorption super-state, while in our case, we can move to a different state in our non-absorb super-state. The possibility of moving to the absorb super-state in a  $\epsilon$  time step, changes from  $\epsilon$  to  $\epsilon p_m$ , and the formula becomes

$$p(x_\tau = n) = \sum_{j < \frac{\tau}{\epsilon}} p(a_j) = \epsilon p_m + \epsilon p_m(1 - \epsilon p_m) + \epsilon p_m(1 - \epsilon p_m)^2 + \dots + \epsilon p_m(1 - \epsilon p_m)^{\lfloor \frac{\tau}{\epsilon} \rfloor - 1} \quad (114)$$

$$p(x_\tau = n) = \epsilon p_m \frac{1 - (1 - \epsilon p_m)^{\frac{1}{\epsilon p_m} \lfloor \frac{\tau}{\epsilon} \rfloor \epsilon p_m}}{1 - (1 - \epsilon p_m)} = 1 - (1 - \epsilon p_m)^{\frac{1}{\epsilon p_m} \lfloor \frac{\tau}{\epsilon} \rfloor \epsilon p_m} \quad (115)$$

which implies that  $p(x_\tau = n) = 1 - e^{-\tau p_m}$ . This means that at time  $\tau$  we are in one of the non-absorption states with a probability of  $e^{-\tau p_m}$ . Now, as before (Appendix A.3.2), we want to find the probability of being at each non-absorption state at time  $\tau$  given that we started a state  $i$ . We can construct these probabilities using a similar approach as before, but with two key differences. First, because we are conditioning on being in the non-absorption super-state, all probabilities must be multiplied by  $e^{-\tau p_m}$ . Second, this conditioning implicitly provides some information about the probability of moving when going from time  $\tau - \epsilon$  to time  $\tau$ . In the uniform case, for a time step  $\epsilon$  this probability was  $\epsilon$ , but now conditioning on the fact that we are not at the absorption state at time  $\tau$  provides extra information on whether we moved when going from  $\tau - \epsilon$  to  $\tau$ . Indeed, given



this extra information, we expect the probability of having moved when going from  $\tau - \epsilon$  to  $\tau$  to be reduced. This becomes obvious when  $p_m = 1$ , as then we are in the case of the absorb diffusion, and saying that at time  $\tau$  we are not at the absorption state, immediately implies that we did not move, as there are no other possibilities. As  $p_m$  decreases, the magnitude of this information decreases. So instead of writing the probability of moving (during an  $\epsilon$  time step) with  $\epsilon$  as we did previously, we write this probability with  $\delta(\epsilon)$ , and as before we can derive that

$$\bar{p}_\tau(k|i) = \bar{p}_{\tau-\epsilon}(j|i)\delta \frac{1}{n-1}(n-3) + (1 - (n-2)\bar{p}_{\tau-\epsilon}(j|i))\delta \frac{1}{n-1} + \quad (116)$$

$$+ \bar{p}_{\tau-\epsilon}(j|i)(1 - \delta) + \bar{p}_{\tau-\epsilon}(j|i)\frac{1}{n-1}\delta \quad (117)$$

thus

$$\bar{p}_\tau(k|i) = \frac{1}{n-1}\delta + \bar{p}_{\tau-\epsilon}(k|i)(1 - \delta) \quad (118)$$

and by the same recursion trick as before,

$$p_{\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon | 0}(k|i) = \frac{1}{n-1}\delta \frac{1 - (1 - \delta)^{\lfloor \frac{\tau}{\epsilon} \rfloor}}{1 - (1 - \delta)} = \frac{1}{n-1}(1 - (1 - \delta)^{\lfloor \frac{\tau}{\epsilon} \rfloor}). \quad (119)$$

Now we derive the expression of  $\delta$ . We write the move event when going from  $\tau - \epsilon$  to  $\tau$  with  $A$  and not going at the absorbed state at time  $\tau$  with  $B$

$$\delta = p(A|B) = \frac{p(A)}{p(B)}p(B|A) = \frac{p(B|A)p(A)}{p(B|A)p(A) + p(B|A^C)p(A^C)} = \quad (120)$$

$$= \frac{(1 - p_m)\epsilon}{(1 - p_m)\epsilon + 1(1 - \epsilon)} = \alpha\epsilon \quad (121)$$

where  $\alpha = \frac{(1-p_m)}{(1-p_m)\epsilon + 1(1-\epsilon)}$  goes to  $1 - p_m$  when  $\epsilon$  goes to 0, therefore finally

$$p_{\lfloor \frac{\tau}{\epsilon} \rfloor \epsilon | 0}(k|i) = \frac{1}{n-1}(1 - (1 - \delta)^{\lfloor \frac{\tau}{\epsilon} \rfloor}) = \frac{1}{n-1}(1 - (1 - \delta)^{\frac{1}{\delta}\delta\lfloor \frac{\tau}{\epsilon} \rfloor}) = \frac{1}{n-1}(1 - (1 - \delta)^{\frac{1}{\delta}\alpha\epsilon\lfloor \frac{\tau}{\epsilon} \rfloor}). \quad (122)$$

Therefore  $p_{\tau|0}(j \notin \{i, n\}|i) = e^{-\tau p_m} \frac{1}{n-1}(1 - e^{(1-p_m)\tau})$ , and  $p_{\tau|0}(n|i) = 1 - e^{-\tau p_m}$ , and  $p_{\tau|0}(i|i) = e^{-\tau p_m}(1 - \frac{n-2}{n-1}(1 - e^{(1-p_m)\tau}))$ , which gives the  $i$ th column of the exponential matrix, when  $i \neq n$ . In case  $i = n$ , the column of the exponential matrix is simply the one hot encoding at  $n$ . Setting  $\tau = \sigma_t$  finishes the derivation.

**Proposition 5** *If we denote with  $\mathbf{Y}_t$  the matrix exponential of  $\sigma_t \mathbf{Q}_{roulette}^{tok} = \sigma_t (\mathbf{I} - \mathbf{P}_{roulette}^{tok})$ , then  $\mathbf{Y}_t(i \notin \{j, n\}, j \neq n) = e^{-\sigma_t p_m} \frac{1}{n-1}(1 - e^{-(1-p_m)\sigma_t})$ ,  $\mathbf{Y}_t(i \neq n, i \neq n) = e^{-\sigma_t p_m}(1 - \frac{n-2}{n-1}(1 - e^{-(1-p_m)\sigma_t}))$ ,  $\mathbf{Y}_t(n, j \neq n) = 1 - e^{-\sigma_t p_m}$ ,  $\mathbf{Y}_t(i \neq n, n) = 0$ , and  $\mathbf{Y}_t(n, n) = 1$ .*

#### A.3.4 CONTROLLING THE NUMBER OF THE CORRECTED TOKENS IN THE REVERSE PROCESS

We wish to see how the choice of  $p_m$  corresponds to the probability of a token moving uniformly at least once before being masked, given some diffusion interval from  $[0, T]$ . This should correspond to the probability of a token being corrected after it is unmasked in the reverse process, therefore enabling us to control the expected number of tokens to be corrected.

As before, we discretize the time interval into subintervals of length  $\epsilon$ . The event of a token moving before getting masked will be denoted by  $X$ , and the event of a token being masked at time  $j\epsilon$  is denoted with  $A_j$ . Therefore

$$p(X) = p(X \cap \Omega) = p(X \cap \{\cup A_j\}) = p(\cup_{j\epsilon < T} X \cap A_j) \quad (123)$$

$$p(X) = \sum_{j < \lfloor \frac{T}{\epsilon} \rfloor} p(X|A_j)p(A_j). \quad (124)$$

$p(A_j)$  is the probability of being at the absorption state at  $j\epsilon$  and not being there at  $(j-1)\epsilon$ . This probability is  $p(A_j) = (1 - e^{-j\epsilon p_m}) - (1 - e^{-(j-1)\epsilon p_m}) = e^{-(j-1)\epsilon p_m} - e^{-j\epsilon p_m}$  therefore  $p(A_j) = e^{-(j-1)\epsilon p_m}(1 - e^{-\epsilon p_m}) \approx e^{-(j-1)\epsilon p_m} \epsilon p_m$ . We write  $\tau_j = j\epsilon$ , and thus we have

$$p(X) = \sum_{j < \lfloor \frac{T}{\epsilon} \rfloor} p(X|A_j)p(A_j) = \sum_{j < \lfloor \frac{T}{\epsilon} \rfloor} \frac{n-2}{n-1}(1 - e^{-(1-p_m)\tau_{j-1}}) \cdot e^{-\tau_{j-1} p_m} \epsilon p_m \quad (125)$$

1350 which is a Riemann sum. Taking the limit we get the following integral  
1351

$$1352 \quad p(X) = p_m \frac{n-2}{n-1} \int_0^T (1 - e^{-(1-p_m)x}) \cdot e^{-xp_m} dx, \quad (126)$$

1354 whose solution is

$$1355 \quad p(X) = p_m \frac{n-2}{n-1} (e^{-T} - \frac{1}{p_m} e^{-Tp_m} + \frac{1}{p_m} - 1). \quad (127)$$

1358 Setting  $t = \sigma_\tau$ , we get

$$1359 \quad p(X) = \frac{n-2}{n-1} (e^{-\sigma_1 p_m} - e^{-\sigma_1 p_m} + 1 - p_m). \quad (128)$$

1362 We can see that when  $\sigma_1$  and  $\sigma_1 p_m$  are relatively large then  $p(X) \approx 1 - p_m$ . That is, the ratio of  
1363 tokens moving before absorption is  $p(X) \approx 1 - p_m$ , or in other words the probability of tokens  
1364 being masked without ever moving is  $p_m$  which is precisely the probability of a token being masked  
1365 in our transition-rate matrix  $\mathbf{Q}^{tok}$ . Thus the last row of our matrix directly controls the percentage  
1366 of corrected tokens in the reverse process for large enough diffusion times.

### 1368 A.3.5 TIME EVOLVING ROULETTE

1369 Here, we generalize the results of the previous subsections for the case that  $p_m$  varies with respect  
1370 to time. We start by defining

$$1371 \quad \mathbf{Q}_{eroulette}^{tok}(t) = [(1 - p_m(t))\sigma(t)]'_t \mathbf{Q}_{uniform}^{tok} + [p_m(t)\sigma(t)]'_t \mathbf{Q}_{absorb}^{tok}. \quad (129)$$

1372 First, we point out that since the size of  $\mathbf{Q}_{eroulette}$  is  $(n \times n)$  and that of  $\mathbf{Q}_{uniform}$  is  $(V \times V)$ , where  
1373  $V = n - 1$ , in order for them to have the same size, we add a row of zeros and a column of zeros to the  
1374 bottom and right of  $\mathbf{Q}_{uniform}$  respectively. Clearly, by definition, these are transition rate matrices,  
1375 and therefore so is  $\mathbf{Q}_{eroulette}^{tok}(t)$ . Indeed, the elements in each of its columns add to 0, and the only  
1376 negative elements are in the diagonals. If we define,  $p_m(t)$  such that  $p_m(0) = 0$  and  $p_m(1) = 1$ ,  
1377 then the limiting distribution will be the one-hot encoding at the absorb state. Furthermore, we can  
1378 compute the exponential matrix of  $\mathbf{Q}_{eroulette}^{tok}(t)$ . It is easy to prove that  $(1 - p_m(t))\sigma(t)\mathbf{Q}_{uniform}^{tok}$   
1379 and  $p_m(t)\sigma(t)\mathbf{Q}_{absorb}^{tok}(t)$  commute with each other, thus

$$1380 \quad e^{\mathbf{Q}_{eroulette}^{tok}(t)} = e^{(1-p_m(t))\sigma(t)\mathbf{Q}_{uniform}^{tok} + p_m(t)\sigma(t)\mathbf{Q}_{absorb}^{tok}}. \quad (130)$$

$$1381 \quad = e^{p_m(t)\sigma_t \mathbf{Q}_{absorb}^{tok}} e^{(1-p_m(t))\sigma_t \mathbf{Q}_{uniform}^{tok}}. \quad (131)$$

1382 Writing  $\alpha_t = p_m(t)\sigma_t$  and  $\beta_t = (1 - p_m(t))\sigma_t$ , shows that we can calculate each of these exponen-  
1383 tial matrices using the same strategy as in Equation (17). Indeed,

$$1384 \quad e^{\alpha_t \mathbf{Q}_{absorb}^{tok}} = \mathbf{I} + \mathbf{Q}_{absorb}^{tok}(1 - e^{-\alpha_t}) \text{ and } e^{\beta_t \mathbf{Q}^{tok}} = \mathbf{I} + \mathbf{Q}_{uniform}^{tok}(1 - e^{-\beta_t}). \quad (132)$$

1385 Multiplying them together, we get the following proposition:

1386 **Proposition 7.** *If we denote with  $\mathbf{Y}_t$  the matrix exponential of  $\sigma_t \mathbf{Q}_{eroulette}^{tok}$ , then  $\mathbf{Y}_t(i \notin \{j, n\}, j \neq$   
1387  $n) = e^{-\sigma_t p_m(t)} \frac{1}{n-1} (1 - e^{-(1-p_m(t))\sigma_t})$ ,  $\mathbf{Y}_t(i \neq n, i \neq n) = e^{-\sigma_t p_m(t)} (1 - \frac{n-2}{n-1} (1 -$   
1388  $e^{-(1-p_m(t))\sigma_t})$ ,  $\mathbf{Y}_t(n, j \neq n) = 1 - e^{-\sigma_t p_m(t)}$ ,  $\mathbf{Y}_t(i \neq n, n) = 0$ , and  $\mathbf{Y}_t(n, n) = 1$ .*

1389 One can see that this is almost identical to Proposition 5, with the only difference being that  $p_m$   
1390 varies with  $t$ . Indeed, if we fix  $p_m$  so that it is constant with respect to time, Equation (129) becomes

$$1391 \quad \mathbf{Q}_{eroulette}^{tok}(t) = p_m \mathbf{Q}_{absorb}^{tok}(t) + (1 - p_m) \mathbf{Q}_{uniform}^{tok}(t) = \mathbf{Q}_{roulette}^{tok}(t), \quad (133)$$

1392 meaning that  $\mathbf{Q}_{eroulette}^{tok}(t)$ , coincides with  $\mathbf{Q}_{roulette}^{tok}(t)$  roulette in this case. This highlights that the  
1393 roulette diffusion is an interpolation between the roulette and uniform diffusion.

1394 One possible choice of  $p_m(t)$  is  $t^{\frac{1}{a}}$  for a positive constant  $a$ .

## B EXPERIMENTAL DETAILS

### B.1 ALGORITHMS

---

#### Algorithm 1 Cross Entropy Training Algorithm

---

**Require:** Network  $f_\theta$ , (total) noise schedule  $\sigma_t$ , data distribution  $p_{\text{data}}$ , token transition matrix  $\mathbf{Q}^{tok}$ , and time  $t \in [0, 1]$ .  
**Sample**  $\mathbf{x}_0 \sim p_0$ ,  $t \sim \mathcal{U}([0, 1])$ .  
**Construct**  $\mathbf{x}_t$  from  $\mathbf{x}_0$ . In particular,  $\mathbf{x}_t^i \sim p_{t|0}(\cdot|\mathbf{x}_0^i) = \exp(\sigma_t \mathbf{Q}^{tok})_{[:,\mathbf{x}_0^i]}$ .  
**Compute**  $L_{ll} = -\sum_{i=1}^L \log f_\theta^i(\mathbf{x}_t, t)[\mathbf{x}_0^i]$ .  
**Backpropagate**  $\nabla_{\theta} L_{ll}$ .  
**Run** optimizer.

---

### B.2 NETWORK ARCHITECTURE AND HYPER-PARAMETERS

This core model is grounded in the diffusion transformer architecture introduced by Peebles & Xie (2023), which integrates time conditioning into the conventional encoder-only transformer framework as established by Vaswani et al. (2017); Devlin et al. (2019). However, it includes minor adjustments, such as the use of rotary positional encoding (Su et al., 2024). The model contains approximately 5% more parameters than a standard transformer (utilized in the case of GPT-2), attributed to the incorporation of time conditioning. Additionally, the same tokenizers and data splits as in prior work are employed to avoid introducing artifacts.

The network is configured with 12 transformer blocks, each featuring 12 attention heads and a hidden size of 768, aligning with the "small" variant of GPT-2. It includes conditioning dimensions set at 128 to facilitate the diffusion process by encoding time-dependent features. Notably, the architecture excludes masking, typical of generative models that generate all tokens simultaneously rather than sequentially. It uses standard scaled dot-product attention mechanisms and incorporates a dropout rate of 0.1 to mitigate overfitting.

In terms of hyperparameters, the model was trained on a single H100 when the sequence length is set at 128, while in the case of sequence lengths of 1024 the model is trained using 8×H100 with a vocabulary size of 50,257 tokens. Training involves a batch size of 512. The training regime is designed for a total of 400,000 iterations.

Training utilizes the OpenWebText dataset, while evaluation is conducted on WikiText-103, with data managed locally to speed up access times. The noise schedule for the diffusion process is log-linear (uniform, absorb), and roulette log-linear (roulette) controlling the variance of noise added incrementally. In both cases we set  $\epsilon = 0.001$  as in (Lou et al., 2024). Sampling for evaluation during training employs an Euler predictor over 128 (and 1024 when  $L = 1024$ ) steps, with noise removal enabled.

For optimization, the model uses the AdamW optimizer with a learning rate of 0.0003, beta parameters of 0.9 and 0.999, and epsilon set to 1e-8. It features no weight decay, focusing on adapting learning without additional regularization. The optimizer includes a warm-up phase of 2,500 steps to stabilize learning dynamics, and employs gradient clipping at a threshold of 1 to prevent gradients from exploding during training. The log-linear noise schedule was used in the absorb and uniform case, while the roulette log-linear one was used in the roulette case.

### B.3 EFFICIENT IMPLEMENTATION IN PRACTICE

#### B.3.1 EFFICIENT ESTIMATION OF THE SCORE IN PRACTICE

For our choice of sparse matrices  $\mathbf{Q}_t$  which can only modify one position at each step, the loss function

$$\mathbb{E}_{t \sim \mathcal{U}(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\cdot|\mathbf{x}_0)} \sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{x}_t, \mathbf{y}) \ell \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}, s_\theta(\mathbf{x}_t, t) \mathbf{y} \right), \quad (134)$$

becomes

$$\mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p_0(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p_{t|0}(\cdot|\mathbf{x}_0)} \sum_{i=1}^L \sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i) \ell \left( \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)}, s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] \right). \quad (135)$$

Now we focus on calculating the loss at position  $i$ , as the total loss across the sequence, will be simply the sum of such individual losses at each position. We notice that the term  $\sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i) s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i]$  in the expression above is trivial to calculate. We simply add the exponentiated outputs of the neural network across the last dimension weighted by  $\mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i)$ , and then subtract  $\sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} s_{\theta}(\mathbf{x}_t, t)_{\mathbf{x}_t^i}^i \cdot \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{x}_t^i)$ . The expression

$$\sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i) \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} \log s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] \quad (136)$$

is more challenging to be calculated efficiently. We follow the approach of (Lou et al., 2024), when using SEDD for training. To reiterate, for each position  $i$ , if we use SEDD for training, we need to calculate the sum of the product  $\frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} \log s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i]$  over  $n$  ratios, where  $n = V$  if  $\mathbf{Q}^{tok} = \mathbf{Q}_{uniform}^{tok}$  and  $n = V + 1$  otherwise. Luckily, for the choices of  $\mathbf{Q}^{tok}$  in this paper (uniform, absorb and roulette), such ratios have relatively simple form and they are almost all the same. For example, in the case of the **absorb diffusion**, there are two cases, either  $\mathbf{x}_t^i$  has not moved ( $\mathbf{x}_t^i = \mathbf{x}_0^i$ ) or it has been masked ( $\mathbf{x}_t^i = n \neq \mathbf{x}_0^i$ ), which can be seen in Figure 4.

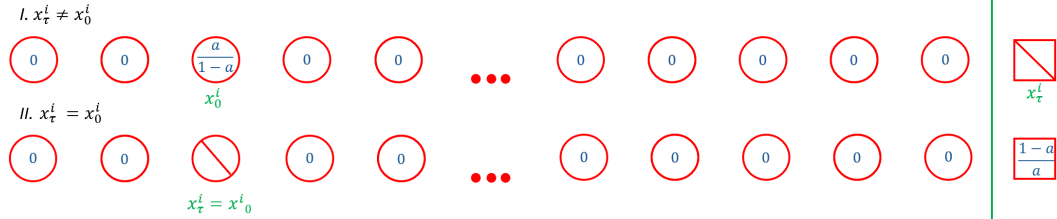


Figure 4: The conditional ratios at position  $i$  over the vocabulary in the two cases. The square represents the absorb state. The value of  $a$  is  $1 - e^{-\sigma\tau}$ .

To simplify notation, when iterating through different tokens in the vocabulary we write  $\mathbf{x}_0^i = k$ ,  $\mathbf{y}^i = j$  and  $\log s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = s_j$ . In the first case, when  $\mathbf{x}_t^i \neq \mathbf{x}_0^i = n$  we get

$$\sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i) \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = \frac{a}{1-a} s_k. \quad (137)$$

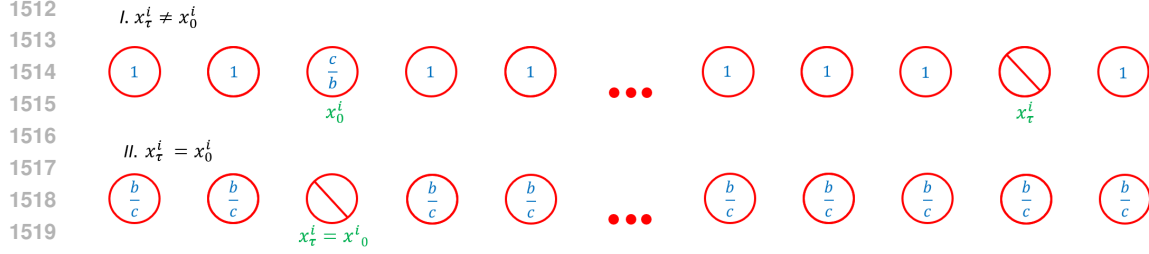
On the other hand, in the second case ( $\mathbf{x}_t^i = \mathbf{x}_0^i$ ) one has

$$\sum_{\mathbf{y}^i \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i) \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = \quad (138)$$

$$\sum_{j \notin \{\mathbf{x}_t^i, n\}} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, j) 0 s_j + \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, n) \frac{1-a}{a} s_n = 0, \quad (139)$$

as  $\mathbf{Q}_t^{tok}(\mathbf{x}_t^i, n) = 0$ . Thus when  $\mathbf{x}_t^i = \mathbf{x}_0^i \neq n$ , we can simply not calculate the loss.

For the **uniform diffusion**, there are also two cases, namely  $\mathbf{x}_t^i = \mathbf{x}_0^i$  and  $\mathbf{x}_t^i \neq \mathbf{x}_0^i$  as illustrated in Figure 5.



1521 Figure 5: The conditional ratios at position  $i$  over the vocabulary in the two cases. We define  
 1522  $b = \frac{1}{n=V}(1 - e^{-\sigma\tau})$  and  $c = 1 - (n - 1)b$ .

1524 If  $x_t^i \neq x_0^i$ , then

1525

1526 
$$\sum_{\mathbf{y}^i \neq \mathbf{x}_0^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i) \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_0^i | \mathbf{x}_0^i)} s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = \sum_{j \notin \{\mathbf{x}_t^i, k\}} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, j) s_j + \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, k) \frac{c}{b} s_k \quad (140)$$

1527

1528 
$$= \sum_{j \neq \mathbf{x}_t^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, j) s_j + \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, k) \left(\frac{c}{b} - 1\right) s_k = \sum_{j \neq \mathbf{x}_t^i} \frac{1}{n} s_j + \frac{1}{n} \left(\frac{c}{b} - 1\right) s_k. \quad (141)$$

1529

1530

1531 Otherwise for  $x_t^i = x_0^i$ , we have

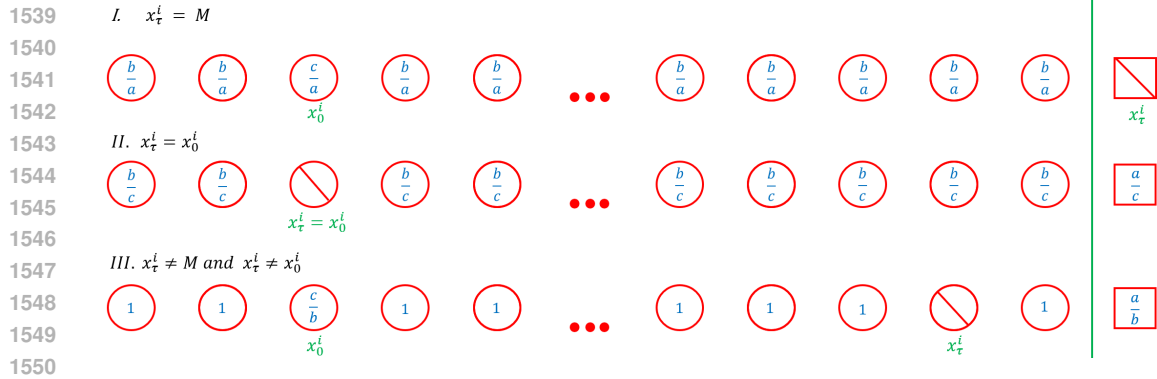
1532 
$$\sum_{\mathbf{y}^i \neq \mathbf{x}_0^i} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, \mathbf{y}^i) \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_0^i | \mathbf{x}_0^i)} s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = \frac{b}{c} \sum_{j \notin \{\mathbf{x}_t^i\}} \mathbf{Q}_t^{tok}(\mathbf{x}_t^i, j) s_j = \frac{1}{n} \sum_{j \neq \mathbf{x}_t^i} s_j. \quad (142)$$

1533

1534

1535 Finally writing  $S^i = \sum_j s_j$ , we get  $\frac{1}{n} S^i + \frac{1}{n} \left(\frac{c}{b} - 1\right) s_k - \frac{1}{n} s_{\mathbf{x}_t^i}$  in the first case, and  $\frac{1}{n} \frac{b}{c} S^i - \frac{1}{n} \frac{b}{c} s_{\mathbf{x}_t^i}$   
 1536 in the second one.

1537 For the **roulette diffusion** we proceed similarly. Figure 6 shows the ratios in all three possible cases.



1551 Figure 6: The conditional ratios at position  $i$  over the vocabulary in the three cases. Since  $x_0^i$  is  
 1552 a token from the data it cannot be masked. The square represents the absorb state. We define  
 1553  $a = 1 - e^{-p_m\sigma\tau}$ ,  $b = e^{-p_m\sigma\tau} \cdot \frac{1 - e^{-\sigma\tau(1-p_m)}}{n-1}$  and  $c = e^{-p_m\sigma\tau} \left[ \frac{1}{n-1} + \left(1 - \frac{1}{n-1}\right) e^{-\sigma\tau(1-p_m)} \right]$ .

1554

1555

1556 In the first case, the sum at position  $i$  is  $p_m \left[ \frac{b}{a} \left( S^i - s_k - s_{\mathbf{x}_t^i} \right) + \frac{c}{a} s_k \right]$ . In the second it is  
 1557  $\frac{1-p_m}{V} \frac{b}{c} \left( S^i - s_k - s_n \right)$ , while in the third case we have  $\frac{1-p_m}{V} \left( S^i - s_k - s_n - s_{\mathbf{x}_t^i} \right) + \frac{1-p_m}{V} s_k \frac{c}{b}$ .  
 1558 If one chooses to incorporate the constant term  $K(a)$ , in the loss, then a similar approach is fol-  
 1559 lowed. However, considering Theorem 4 for evaluation, and the fact that a constant term does not  
 1560 affect the gradient during training, implementing this term is not necessary.

### 1562 B.3.2 SEDD SCALING

1563

1564 During training, some scores  $\frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_0^i | \mathbf{x}_0^i)}$ , on expectation with respect to  $x_0^i$ , can vary significantly for  
 1565 different  $\mathbf{y}^i$ . This makes the task of the network more complex and reduces performance. In order to

alleviate this issue, in the case of the **absorb diffusion**, (Lou et al., 2024) add  $\log \sum_{\mathbf{x}_0^i} \frac{1}{n} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)}$  to the output of the network  $\log s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i]$ . In this case, since we only need to calculate scores when  $\mathbf{x}_t^i = n$ , this sum equals

$$\sum_{\mathbf{x}_0^i} \frac{1}{n-1} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)} = \sum_{\mathbf{x}_0^i} \frac{1}{n-1} \frac{p_{t|0}(\mathbf{y}^i \neq n | \mathbf{x}_0^i)}{p_{t|0}(n | \mathbf{x}_0^i)} = 0 + \frac{1}{n-1} \frac{p_{t|0}(\mathbf{y}^i \neq n | \mathbf{y}^i \neq n)}{p_{t|0}(n | \mathbf{y}^i \neq n)} \quad (143)$$

$$= \frac{1}{n-1} \frac{e^{-\sigma t}}{1 - e^{-\sigma t}} = \frac{1}{n-1} \frac{1}{e^{\sigma t} - 1}. \quad (144)$$

We notice that in effect, this changes the prediction of the score, from  $\log s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i]$  to  $\log s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] - \log(n-1) - \log(e^{\sigma t} - 1)$ .

We extend this strategy for the uniform and the roulette diffusion. In the case of the **uniform diffusion**, we get

$$\sum_{\mathbf{x}_0^i} \frac{1}{V} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)} = \sum_{\mathbf{x}_0^i \notin \{\mathbf{y}^i, \mathbf{x}_t^i\}} \frac{1}{V} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)} + \frac{1}{V} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{y}^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{y}^i)} + \frac{1}{V} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_t^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_t^i)}. \quad (145)$$

Now,

$$\frac{p_{t|0}(\mathbf{y}^i | \mathbf{y}^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{y}^i)} = \frac{1 - \frac{V-1}{V}(1 - e^{-\sigma t})}{\frac{1}{V}(1 - e^{-\sigma t})} = 1 + \frac{V}{e^{\sigma t} - 1}, \quad (146)$$

and

$$\frac{1}{V} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_t^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_t^i)} = \frac{\frac{1}{V}(1 - e^{-\sigma t})}{1 - \frac{V-1}{V}(1 - e^{-\sigma t})} = 1 - \frac{V}{e^{\sigma t} - 1 + V}. \quad (147)$$

Since  $\frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i \neq \mathbf{y}^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i \neq \mathbf{x}_t^i)} = 1$ , we conclude

$$\sum_{\mathbf{x}_0^i} \frac{1}{V} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)} = \frac{V-2}{V} + \frac{1}{V} + \frac{1}{e^{\sigma t} - 1} + \frac{1}{V} - \frac{1}{e^{\sigma t} - 1 + V} = 1 + \frac{1}{e^{\sigma t} - 1} - \frac{1}{e^{\sigma t} - 1 + V}.$$

We then modify the output of the network  $\log s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i]$  accordingly

$$\log s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] + \log \left( \sum_{\mathbf{x}_0^i} \frac{1}{V} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)} \right), \text{ where } n := V. \quad (148)$$

Similarly, we can calculate this scaling factor for the **roulette diffusion**. There are two cases, either  $\mathbf{x}_t^i = n$  or  $\mathbf{x}_t^i \neq n$ . In the first case, with little modifications from before we can show that

$$\sum_{\mathbf{x}_0^i} \frac{1}{n-1} \frac{p_{t|0}(\mathbf{y}^i \neq n | \mathbf{x}_0^i)}{p_{t|0}(n | \mathbf{x}_0^i)} = \frac{1}{n-1} \frac{1 - (1 - e^{-\sigma t p_m})}{1 - e^{-\sigma t p_m}} = \frac{1}{n-1} \frac{1}{e^{\sigma t p_m} - 1}. \quad (149)$$

In the second case,

$$\sum_{\mathbf{x}_0^i} \frac{1}{n-1} \frac{p_{t|0}(\mathbf{y}^i \neq n | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i \neq n | \mathbf{x}_0^i)} = 1 + \frac{1}{e^{(1-p_m)\sigma t} - 1} - \frac{1}{e^{(1-p_m)\sigma t} - 1 + n - 1}. \quad (150)$$

We then modify the output of the network  $\log s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i]$  accordingly

$$\log s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] + \log \left( \sum_{\mathbf{x}_0^i} \frac{1}{n} \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)} \right). \quad (151)$$

We take a moment to recall that if  $\mathbf{y}$  and  $\mathbf{x}_t$  differ only at one position  $i$  then

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0^i} p^i(\mathbf{x}_0^i | \mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}^i | \mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i | \mathbf{x}_0^i)}. \quad (152)$$

Thus, above we are scaling the output by a naive estimation of this expectation where  $p^i(\mathbf{x}_0^i | \mathbf{x}_t)$  is assumed to be  $\frac{1}{n-1}$ . In fact, the study of the scaling approach and the realization that the entity  $\frac{1}{n-1}$  can be more properly estimated, is what originally motivated the derivation and usage of CEDD in this paper.

### B.3.3 SCORE ESTIMATION THROUGH CEDD IN PRACTICE

In this subsection, similarly to Section B.3.1, we show how to convert the learned probabilities into learned ratios in practice. Suppose that the perturbed token at position  $i$  at time  $t$  is  $\mathbf{x}_t^i$ . If  $\mathbf{y}$  is a sequence that differs from  $\mathbf{x}_t$  only at position  $i$  then from before one can write:

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0^i} p^i(\mathbf{x}_0^i|\mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)}. \quad (153)$$

We point out that since  $\mathbf{x}_t$  and  $\mathbf{y}$  only differ at position  $i$ , then  $\mathbf{x}_t^i \neq \mathbf{y}^i$ . Thus,

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0^i \notin \{\mathbf{x}_t^i, \mathbf{y}^i\}} p^i(\mathbf{x}_0^i|\mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_0^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_0^i)} + p^i(\mathbf{y}^i|\mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}^i|\mathbf{y}^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{y}^i)} + p^i(\mathbf{x}_t^i|\mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}^i|\mathbf{x}_t^i)}{p_{t|0}(\mathbf{x}_t^i|\mathbf{x}_t^i)}. \quad (154)$$

In the case of the **absorb diffusion**, from Section B.3.1, we can see that we do not need the ratio when  $\mathbf{y}^i = n$ , and furthermore, we only need the ratios when  $\mathbf{x}_t^i = n$ . Since also the data is not masked, then  $\mathbf{x}_0^i \neq n$ . Therefore, the first term (the sum) in Equation (154) is 0, since  $\frac{p_{t|0}(\mathbf{y}^i \neq n | \mathbf{x}_0^i \notin \{n, \mathbf{y}^i\})}{p_{t|0}(n | \mathbf{x}_0^i \notin \{n, \mathbf{y}^i\})} = 0$ . Hence

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = p^i(\mathbf{y}^i|\mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}^i \neq n | \mathbf{y}^i \neq n)}{p_{t|0}(n | \mathbf{y}^i \neq n)} + p^i(\mathbf{x}_t^i|\mathbf{x}_t) \frac{p_{t|0}(\mathbf{y}^i \neq n | n)}{p_{t|0}(n | n)}, \quad (155)$$

where  $p_{t|0}(\mathbf{y}^i \neq n | n) = 0$  since a masked token cannot be unmasked in the forward process. Thus our approximation of the score is simply

$$s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] = f_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] \frac{e^{-\sigma_t}}{1 - e^{-\sigma_t}}. \quad (156)$$

To conclude  $s_\theta(\mathbf{x}_t, t) = f_\theta(\mathbf{x}_t, t) \frac{1}{e^{\sigma_t} - 1}$ .

In the case of the **uniform diffusion**, we write  $b = \frac{1}{n=V}(1 - e^{-\sigma_t})$  and  $c = 1 - (n - 1)b$  as in Section B.3.1, and Equation (154) becomes

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0^i \notin \{\mathbf{x}_t^i, \mathbf{y}^i\}} p^i(\mathbf{x}_0^i|\mathbf{x}_t) + p^i(\mathbf{y}^i|\mathbf{x}_t) \frac{c}{b} + p^i(\mathbf{x}_t^i|\mathbf{x}_t) \frac{b}{c}. \quad (157)$$

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0^i} p^i(\mathbf{x}_0^i|\mathbf{x}_t) + p^i(\mathbf{y}^i|\mathbf{x}_t) \left(\frac{c}{b} - 1\right) + p^i(\mathbf{x}_t^i|\mathbf{x}_t) \left(\frac{b}{c} - 1\right). \quad (158)$$

Therefore, our approximation of the score is simply

$$s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] = 1 + f_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] \left(\frac{c}{b} - 1\right) + f_\theta^i(\mathbf{x}_t, t)[\mathbf{x}_t^i] \left(\frac{b}{c} - 1\right), \quad (159)$$

or

$$s_\theta^i(\mathbf{x}_t, t) = \mathbf{1} + f_\theta^i(\mathbf{x}_t, t) \left(\frac{c}{b} - 1\right) + f_\theta^i(\mathbf{x}_t, t)[\mathbf{x}_t^i] \left(\frac{b}{c} - 1\right) \mathbf{1}, \quad (160)$$

in vector form.

Finally, in the case of **roulette diffusion**, we can discern two possibilities, namely  $\mathbf{x}_t^i = n$  and  $\mathbf{x}_t^i \neq n$ . In the first one, writing as in Section B.3.1,  $a = 1 - e^{-p_m \sigma_t}$ ,  $b = e^{-p_m \sigma_t} \cdot \frac{1 - e^{-\sigma_t(1-p_m)}}{n-1}$  and  $c = e^{-p_m \sigma_t} \left[ \frac{1}{n-1} + \left(1 - \frac{1}{n-1}\right) e^{-\sigma_t(1-p_m)} \right]$ , Equation 154 becomes

$$\frac{p_t(\mathbf{y})}{p_t(\mathbf{x}_t)} = \sum_{\mathbf{x}_0^i \notin \{\mathbf{x}_t^i, \mathbf{y}^i\}} p^i(\mathbf{x}_0^i|\mathbf{x}_t) \frac{b}{a} + p^i(\mathbf{y}^i|\mathbf{x}_t) \frac{c}{b} = \frac{b}{a} + p^i(\mathbf{y}^i|\mathbf{x}_t) \left(\frac{c}{b} - \frac{b}{a}\right). \quad (161)$$

Hence, if the perturbed token  $\mathbf{x}_t^i$  at position  $i$  is masked, the predicted ratios over the vocabulary are:

$$s_\theta^i(\mathbf{x}_t, t) = \mathbf{1} \frac{b}{a} + f_\theta^i(\mathbf{x}_t, t) \left(\frac{c}{b} - \frac{b}{a}\right). \quad (162)$$

Very similarly, if the perturbed token  $\mathbf{x}_t^i$  at position  $i$  is not masked, the predicted ratios over the vocabulary are:

$$s_\theta^i(\mathbf{x}_t, t) = \mathbf{1} + f_\theta^i(\mathbf{x}_t, t) \left( \frac{c}{b} - 1 \right) + f_\theta^i(\mathbf{x}_t, t)[\mathbf{x}_t^i] \left( \frac{b}{c} - 1 \right) \mathbf{1}, \quad (163)$$

This shows that we can calculate the ratios by scaling and translating the learned probabilities through CEDD.

#### Modification of the reparametrization in (19): Re-scaling conditional ratios when $t \rightarrow 0$ .

In the case of the uniform diffusion dynamics, when  $t \rightarrow 0$  then  $\frac{b}{c} \rightarrow 0$  and  $\frac{c}{b} \rightarrow \infty$ . This does not cause any issues when sampling, however, it does negatively impact the perplexity bound. We explain informally the issue below. We fix position  $i$  and assume we are at  $\mathbf{x}_t^i$  and that originally we were at  $\mathbf{x}_0^i$ . When the time  $t$  is close to 0, the model  $f_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i]$  is typically very confident and can predict that some token  $\mathbf{y}^i$  is  $\mathbf{x}_0^i$ , that is,  $f_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] \approx 1 - \epsilon$  and  $f^i$  is  $\approx \frac{\epsilon}{V-1}$  everywhere else. Furthermore, most of the time  $\mathbf{x}_t^i = \mathbf{x}_0^i$ . This means the score  $s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] \approx \frac{c}{b}$  and is  $\approx \frac{b}{c}$  everywhere else. If the prediction is correct then the loss becomes  $\approx \frac{1}{V-1}(V-2)[- \frac{b}{c} \log(\frac{b}{c}) + \frac{b}{c} + \frac{b}{c}(\log(\frac{b}{c}) - 1)] + [- \frac{c}{b} \log(\frac{c}{b}) + \frac{c}{b} + \frac{c}{b}(\log(\frac{c}{b}) - 1)] \approx 0$ . If it misses however, then we have  $\approx \frac{1}{V-1}(V-3)[- \frac{b}{c} \log(\frac{b}{c}) + \frac{b}{c} + \frac{b}{c}(\log(\frac{b}{c}) - 1)] + [- \frac{b}{c} \log(\frac{c}{b}) + \frac{c}{b} + \frac{b}{c}(\log(\frac{b}{c}) - 1)] + [- \frac{c}{b} \log(\frac{b}{c}) + \frac{b}{c} + \frac{c}{b}(\log(\frac{c}{b}) - 1)] \approx 2 \frac{c}{b} \log \frac{c}{b}$ . Considering the large magnitude of  $\frac{c}{b}$  when  $t \approx 0$ , this loss is extremely punitive. SEDD does not suffer from this issue, as the conditional ratios  $\frac{c}{b}$  are implicitly learned by the network, which due to its limited flexibility likely acts as a regularizer, not allowing the values of the ratios to rise steeply as  $t \rightarrow 0$ . Motivated by this, in the uniform case, we rescale  $\sigma_t < 0.0015$ , by setting  $\sigma_t = 0.0015$ . This significantly reduces the magnitude of  $\frac{c}{b}$  as  $t \rightarrow 0$ . We highlight that only the  $\sigma_t$  that is used to calculate  $\frac{c}{b}$  and  $\frac{b}{c}$  are scaled while the  $\sigma_t$  that is fed to the model is not touched. Naturally, for the sake of rigor, this is also the model we employ to generate samples. Finally, we note that this problem also appears in the case of the roulette diffusion dynamics, as when time is close to 0 most tokens are unmasked. The same strategy is applied, as before, only to ratios  $\frac{c}{b}$  and  $\frac{b}{c}$ , by rescaling  $\sigma_t$  when  $\sigma_t < 0.5$  as follows:  $\sigma_t^{scaled} = \log(1.1\sigma_t + 1.1)$ .

#### B.3.4 ANALYTIC SAMPLING IN PRACTICE

In (Lou et al., 2024) an alternative sampling scheme (Equation (18)) is provided. This method is called the analytic method and it performs better than Euler sampling, in particular when the number of sampling steps is small:

$$p_{t-\epsilon|t}(\mathbf{x}_{t-\epsilon}^i | \mathbf{x}_t^i) = \left( e^{\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(\mathbf{x}_t^i, \mathbf{x}_{t-\epsilon}^i) \right) \sum_{\mathbf{y}^i=1}^n \left( e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(\mathbf{x}_{t-\epsilon}^i, \mathbf{y}^i) \right) s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i], \quad (164)$$

where  $\sigma_t^{\Delta t} = \sigma_t - \sigma_{t-\epsilon}$ .

Since we have an analytic expression of the matrix exponential  $e^{\sigma_t^{\Delta t} \mathbf{Q}^{tok}}$ , we can easily derive the expression of  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}$ , by simply substituting  $\sigma_t^{\Delta t}$  with  $-\sigma_t^{\Delta t}$  in each entry. One can efficiently calculate the sum above by following the strategy of (Lou et al., 2024) as below.

In the case of the **absorb diffusion**,  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(j \neq n, j \neq n) = 1 - \bar{a}$ ,  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(n, n) = 1$ ,  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(n, j \neq n) = \bar{a}$ , with all other entries being 0, where  $\bar{a} = 1 - e^{\sigma_t^{\Delta t}}$ . As such if  $\mathbf{x}_{t-\epsilon} \neq n$

$$\sum_{\mathbf{y}^i=1}^n \left( e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(\mathbf{x}_{t-\epsilon}^i, \mathbf{y}^i) \right) s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] = (1 - \bar{a}) s_\theta^i(\mathbf{x}_t, t)[\mathbf{x}_{t-\epsilon}^i], \quad (165)$$

and if  $\mathbf{x}_{t-\epsilon} = n$

$$\sum_{\mathbf{y}^i=1}^n \left( e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(\mathbf{x}_{t-\epsilon}^i, \mathbf{y}^i) \right) s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] = \bar{a} \sum_{\mathbf{y}^i=1}^{n-1} s_\theta^i(\mathbf{x}_t, t)[\mathbf{y}^i] + s_\theta^i(\mathbf{x}_t, t)[n] = \quad (166)$$



$$= \bar{a} \sum_{\mathbf{y}^i=1}^n s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] + (1 - \bar{a}) s_{\theta}^i(\mathbf{x}_t, t)[n] = \bar{a} S_{\theta}^i(\mathbf{x}_t, t) + (1 - \bar{a}) s_{\theta}^i(\mathbf{x}_t, t)[n], \quad (167)$$

where  $S_{\theta}^i(\mathbf{x}_t, t) = \sum_{\mathbf{y}^i=1}^n s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i]$ .

In the case of the **uniform diffusion**  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(j, j) = \bar{c}$  while the rest of the entries are  $\bar{b}$ , where  $\bar{b} = \frac{1}{n-V}(1 - e^{\sigma_t^{\Delta t}})$  and  $\bar{c} = 1 - (n-1)\bar{b}$ . Therefore

$$\sum_{\mathbf{y}^i=1}^n \left( e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(\mathbf{x}_{t-\epsilon}^i, \mathbf{y}^i) \right) s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = (\bar{c} - \bar{b}) s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{x}_{t-\epsilon}^i] + \bar{b} S_{\theta}^i(\mathbf{x}_t, t) = \quad (168)$$

$$= \frac{s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{x}_{t-\epsilon}^i]}{e^{-\sigma_t^{\Delta t}}} + \frac{e^{-\sigma_t^{\Delta t}} - 1}{ne^{-\sigma_t^{\Delta t}}} S_{\theta}^i(\mathbf{x}_t, t). \quad (169)$$

Finally, in the case of the **roulette diffusion**, we have  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(j \neq n, j \neq n) = \bar{c}$ ,  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(j = n, j \neq n) = \bar{a}$ ,  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(j = n, j = n) = 1$ ,  $e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(j \neq n, j = n) = 0$  and the rest of the entries are  $\bar{b}$ , where  $\bar{a} = 1 - e^{p_m \sigma_t^{\Delta t}}$ ,  $\bar{b} = e^{p_m \sigma_t^{\Delta t}} \cdot \frac{1 - e^{\sigma_t^{\Delta t}}}{n-1}$  and  $\bar{c} = e^{p_m \sigma_t^{\Delta t}} \left[ \frac{1}{n-1} + \left(1 - \frac{1}{n-1}\right) e^{\sigma_t^{\Delta t} (1-p_m)} \right]$ . As in the absorb case, there are two cases, the first one being  $\mathbf{x}_{t-\epsilon} \neq n$ :

$$\sum_{\mathbf{y}^i=1}^n \left( e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(\mathbf{x}_{t-\epsilon}^i, \mathbf{y}^i) \right) s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = (\bar{c} - \bar{b}) s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{x}_{t-\epsilon}^i] + \bar{b} S_{\theta}^i(\mathbf{x}_t, t) - \bar{b} s_{\theta}^i(\mathbf{x}_t, t)[n], \quad (170)$$

while for  $\mathbf{x}_{t-\epsilon} = n$ , one derives:

$$\sum_{\mathbf{y}^i=1}^n \left( e^{-\sigma_t^{\Delta t} \mathbf{Q}^{tok}}(\mathbf{x}_{t-\epsilon}^i, \mathbf{y}^i) \right) s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{y}^i] = \bar{a} S_{\theta}^i(\mathbf{x}_t, t) + (1 - \bar{a}) s_{\theta}^i(\mathbf{x}_t, t)[n] = \quad (171)$$

$$(\bar{c} - \bar{b}) s_{\theta}^i(\mathbf{x}_t, t)[n] + \bar{b} S_{\theta}^i(\mathbf{x}_t, t) - \bar{b} s_{\theta}^i(\mathbf{x}_t, t)[n] + (\bar{a} - \bar{b}) S_{\theta}^i(\mathbf{x}_t, t) + s_{\theta}^i(\mathbf{x}_t, t)[n](-\bar{a} + 2\bar{b} - \bar{c} + 1). \quad (172)$$

It is important to remark that in the case of the absorb and roulette diffusion,  $s_{\theta}^i(\mathbf{x}_t, t)[n]$  plays a role in the quantities above only when  $\mathbf{x}_t^i = n$ . In the previous section, it is mentioned that  $s_{\theta}^i(\mathbf{x}_t, t)[n]$  is not learned, but in this case when  $\mathbf{x}_t^i = n$  is needed. However, this particular case is not problematic as  $s_{\theta}^i(\mathbf{x}_t, t)[n] = s_{\theta}^i(\mathbf{x}_t, t)[\mathbf{x}_t^i] = 1$ , thus we set  $s_{\theta}^i(\mathbf{x}_t, t)[n]$  to 1 manually.

#### B.4 PROCEDURE FOR GENERATING THE PLOTS IN FIGURE 4

We sample a batch of 16 test points, and we perturb each datapoint with respect to a different  $t$ . Then, for example in the absorb case, we calculate the expressions

$$\sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{x}_t, \mathbf{y}) \ell \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}, s_{\theta}(\mathbf{x}_t, t)_{\mathbf{y}} \right), \quad (173)$$

and

$$\sum_{\mathbf{y} \neq \mathbf{x}_t} \mathbf{Q}_t(\mathbf{x}_t, \mathbf{y}) \bar{\ell} \left( \frac{p_{t|0}(\mathbf{y}|\mathbf{x}_0)}{p_{t|0}(\mathbf{x}_t|\mathbf{x}_0)}, s_{\theta}(\mathbf{x}_t, t)_{\mathbf{y}} \right) - 1 + \epsilon, \quad (174)$$

respectively for  $J_1$  and  $J_2$ , where  $\epsilon = 10^{-4}$ . In each case, this returns a single point, therefore we repeat this procedure  $64N$  times, where  $N$  is the number of testing points. This ensures that in all cases ( $J_1$  and  $J_2$ ), we are estimating the loss in the same number of perturbed points. However, when  $N$  is small, we keep testing  $J_1$  and  $J_2$  (extending their plots) for more points in order to provide to the reader more information about their (limiting) behaviour.

Table 5: Results comparing the roulette transition-matrices for different  $p_m$  in terms of generative perplexity.

	$p_m = 0.95$	$p_m = 0.65$	$p_m = 0.35$	$p_m = 0.05$
Roulette	72.31	85.72	124.19	284.55

## B.5 PER SEQUENCE ESTIMATION OF $J_1$

Here we explain a different method of calculating  $J_1$ . When applying this method one first samples a point (sequence) from the test set, and perturbs it for (say) 1024 time values  $t$ . Then the 1024 perturbed points are separated into (e.g.) 64 batches of size 16, and Expression 173 is computed for each batch which returns 64 values. This process is then repeated for the next test point (sequence). If the test set has  $N$  points then the process produces  $64N$  values. These values can be averaged, divided by  $L$  and exponentiated.

## B.6 ADDITIONAL RESULTS

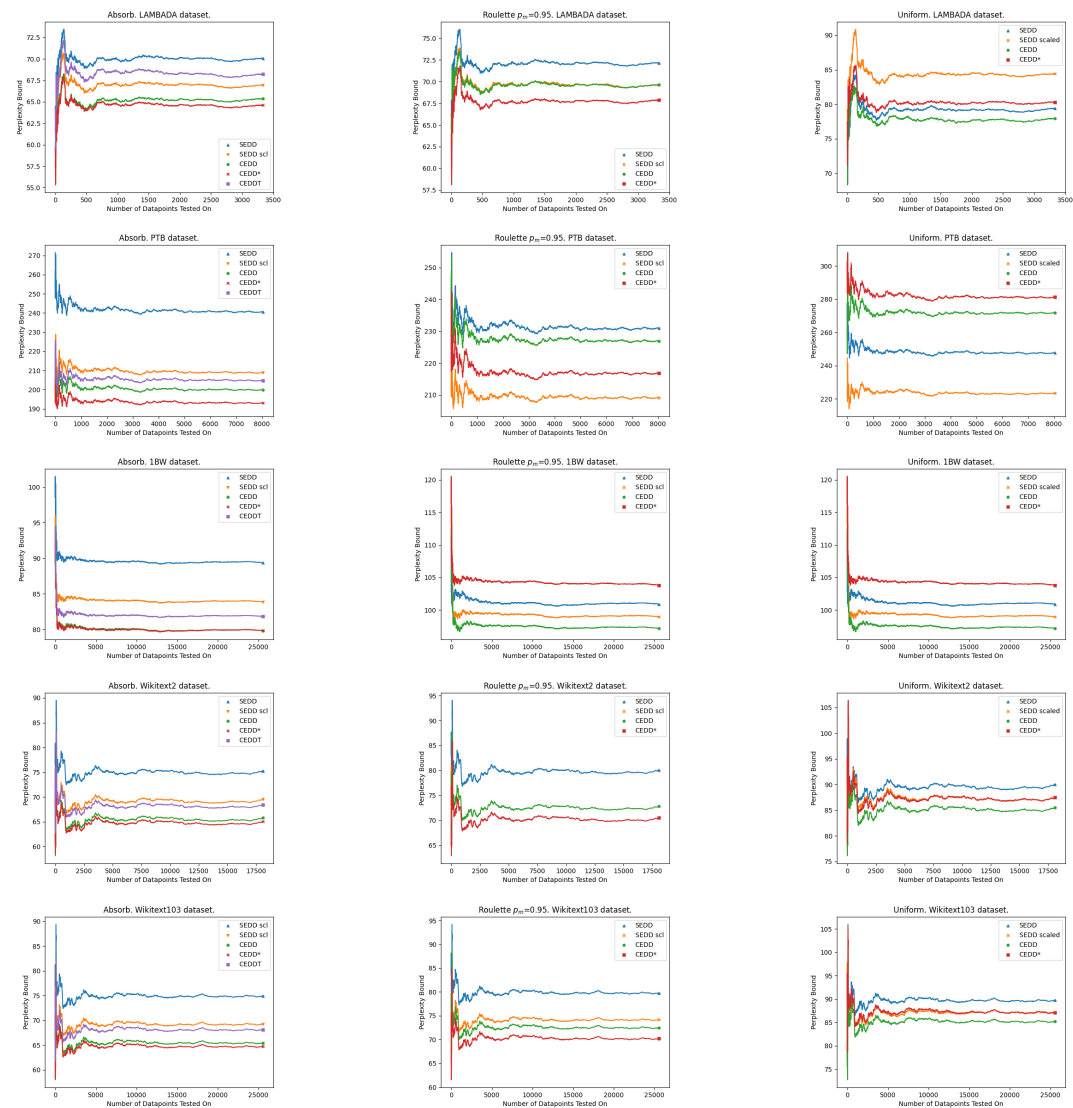
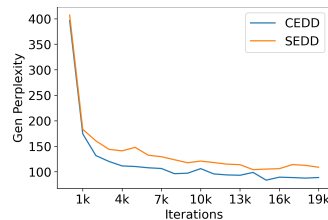


Figure 7: The comparison of the testing history of various methods on various datasets. Each row represents a different dataset, and each column a different type of diffusion transition-rate matrix. The way we calculated  $J_1$  to produce the plots is explained in Appendix B.5. The curves presented are the exponentiated cumulative version of the array returned by the described process. Since we do not perturb the test set, this method enables a clearer visual comparison.

1836 Table 6: Results comparing SEDD, SEDD scaled (SEDDs), CEDD and CEDD\* using  $J_2$ . Lower is  
 1837 better.

1839	<b>Model (L=128)</b>	<b>LAMBADA</b>	<b>WikiText2</b>	<b>PTB</b>	<b>WikiText103</b>	<b>1BW</b>
1840	SEDD Absorb	69.33	74.58	238.35	74.27	88.35
1841	SEDDs Absorb	66.83	68.93	207.76	68.49	83.26
1842	CEDD Absorb	64.87	65.07	198.00	64.97	<b>78.93</b>
1843	CEDD* Absorb	<b>64.11</b>	<b>64.54</b>	<b>191.38</b>	<b>64.30</b>	79.17
1844	SEDD Roulette	71.37	79.21	227.98	78.67	92.21
1845	SEDDs Roulette	68.25	73.43	206.34	72.77	86.99
1846	CEDD Roulette	68.92	72.16	224.18	71.52	<b>85.37</b>
1847	CEDD* Roulette	<b>67.27</b>	<b>69.61</b>	<b>213.90</b>	<b>69.45</b>	85.64
1848	SEDD Uniform	80.09	91.37	249.50	90.63	101.82
1849	SEDDs Uniform	80.29	88.48	<b>226.03</b>	87.73	99.73
1850	CEDD Uniform	<b>79.46</b>	<b>86.82</b>	276.61	<b>86.52</b>	<b>98.44</b>
1851	CEDD* Uniform	82.43	89.68	289.09	88.90	106.32



1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861 Figure 8: Comparison of CEDD\* absorb L=1024 and SEDDs absorb 1024 in terms of generative  
 1862 perplexity (Euler, 128 steps), with respect to number of parameter updates.

1863  
 1864  
 1865  
 1866 Table 7: Results comparing SEDD (retrained), CEDD\* trained for 20k parameter updates. For  
 1867 generation we use the analytic method with 1024 steps. Float 32 sampling.

1869	<b>Model (Absorb)</b>	<b>GenPerp</b>	<b>LAMBADA</b>	<b>WikiText2</b>	<b>PTB</b>	<b>WikiText103</b>	<b>1BW</b>
1870	SEDDs L=1024	55.27	67.92	64.91	173.38	64.13	107.64
1871	CEDD* L=1024	48.29	63.75	57.58	157.71	57.26	98.68

1872  
 1873  
 1874  
 1875 Table 8: Results comparing SEDD, SEDD scaled (SEDDs), CEDD and CEDD\* in terms of genera-  
 1876 tive perplexity.

1877	<b>Model (L=128)</b>	<b>GenPerp</b> Fl32-GPT2L	<b>GenPerp</b> Fl64-GPT2L	<b>GenPerp</b> Fl64-LLama8B
1878	SEDD Absorb	83.62	172.35	212.15
1879	SEDDs Absorb	79.74	166.35	206.34
1880	CEDD Absorb	74.19	148.21	185.90
1881	CEDD* Absorb	<b>72.13</b>	<b>143.86</b>	<b>183.74</b>
1882	SEDD Roulette	87.81	178.94	220.00
1883	SEDDs Roulette	83.36	172.93	212.07
1884	CEDD Roulette	76.71	167.67	208.84
1885	CEDD* Roulette	<b>72.31</b>	<b>158.56</b>	<b>197.79</b>
1886	SEDD Uniform	175.49	169.66	206.91
1887	SEDDs Uniform	171.07	163.88	200.90
1888	CEDD Uniform	<b>168.35</b>	<b>161.84</b>	<b>200.09</b>
1889	CEDD* Uniform	179.30	175.42	213.99

1890 Table 9: Results comparing SEDD, SEDD scaled (SEDDs), CEDD and CEDD\* in terms of genera-  
 1891 tive perplexity.  
 1892

1893	<b>Model (Absorb)</b>	<b>GenPerp</b> Fl32-GPT2L	<b>GenPerp</b> Fl64-GPT2L	<b>GenPerp</b> Fl64-LLama8B
1894	SEDDs L=128	79.74	166.35	206.34
1895	CEDD* L=128	<b>72.13</b>	<b>143.86</b>	<b>183.74</b>
1896	CEDDT L=128	74.07	154.04	195.41
1897	SEDDs L=1024	40.95	105.27	111.87
1898	CEDD* L=1024	<b>40.93*</b>	<b>101.83</b>	<b>107.32</b>
1899	CEDDT L=1024	42.18	108.88	115.60
1900	GPT-2 L=1024	41.02	41.02*	50.25*

1901  
 1902

1903 Table 10: Correction accuracy percentages. 50k training iterations, batch size of 32 and  $L = 128$ .  
 1904

1905	<b>Model (L=128)</b>	CEDD Uniform	SEDD Uniform	CEDD Roulette	SEDD Roulette
1906	PAP	90.7	89.9	<b>90.8</b>	87.9
1907	CAP	91.2	90.3	<b>91.3</b>	88.6

1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943

## 1944 C GENERATED EXAMPLES

1945

1946

## 1947 C.1 GENERATED TEXT

1948

## 1949 Listing 1: Generated texts from Absorb diffusion trained using CEDD, L=128.

1950

1951

1952

1953

1954

1955

1956

1957

Any version of the bundle allows corruption. The campaign was delayed a year from testing ground, and because of other deficiencies. Perhaps this can be used to inform all the Xbox One games come with the same package. Whilst there will be no problems, there will probably be no corruption price-wise and we can expect the same with one Xbox One game package at the same time. My opinion is a solid one. Let's assume Japan as the entry point in the entire country, and PS3 probably supports only one of these. It is effectively competing (see video) against hardware arcade titles such as PS4 or XPS4 as well

1958

1959

1960

1961

1962

1963

1964

1965

1966

his Achilles. I was ecstatic to be able to watch every one of them and keep them healthy tonight. You can look at the building blocks in college basketball, see that they play very hard and maybe feel superior to somebody. In my gut, however, you can argue that the one-year early good contracts come back and have a point. There has been a slight step-off with John Wall, which was a bit short on my mind. I now know Wall will be listed on the waiver, but I believe with Vic Beasley and Louren Maye this season, those two, with some really good shooters on the floor,

1967

1968

1969

1970

1971

1972

1973

1974

1975

1976

1977

times when it is vital for Scotland's health and welfare The HRH Statement of Purpose and Development (SSTO) provides a framework for the government's creation of the Future of Europe as a single body independent of Scotland's national welfare system. The government's own intention is to actively develop its vision for the UK's welfare system based on its very nature and how the set of policies and key principals of the forum was established.

The Future of Europe Working Group (STO) Omissions of the UK Council on Europe's Audit Office (æsvishig

they go along with so many top players.

1978

1979

1980

1981

1982

1983

1984

1985

1986

1987

1988

1989

1990

1991

1992

1993

1994

1995

1996

1997

A pawn,

Listing 2: Generated texts from Uniform diffusion trained using CEDD, L=128.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

is to facilitate market conditions in jurisdictions. Also, not unlike that elsewhere where in terms of lack of in-house attorneys with no background qualified, or are looking to become in England and the whole country in the U.S., Arizona, New England, California, Connecticut, New Hampshire, New York, Massachusetts, Baltimore, Buffalo and Philadelphia, there are price measures set to rise to the model. And they come in handy as there are extraordinary levels of price in commodities, like alcohol and other drugs. Bankruptcy approached, savings rate rigging now make even more dangerous. With broadest institutional liquidity assets to sell and demand a

=====

aircraft, under the Mohammed Boland project.[28][30][31] In 2010, the government collaborated with Qantas states hoping to build the first pilots and aircraft carriers,[33][34] announced that on 16 October they decommissioned 25 acres of the Mozambique plantation leased by a Boeing 737, including several African pilots whom have been imparted on Brazilian soil and ended in. A military expedition - along with the arms embargo and tax evasion - has been delayed until federal Unaboprevalence halted part of the construction.[22]

Engineer units, and airport infrastructure [ edit ]

The

=====

scene. Though a woman called, saying that whether it was being on the suspect's video is not the answer, Desiling says an undisclosed number were present at the mall, and that the police were ready.

In the end, two shoppers, Florimore, 27, sprinted outside the convenience store after struggling with downpour to get inside the shop.

When Florimore left the restroom at the store, Zimmerman entered the vendor's cell phone on Sunday and secured her keys Saturday, Desiling said. He helped he lock a key before being used to contract it. After the taped encounter with a current Illinois Bureau of Police

=====

, some of which were then scrutinized by the Bernanke Advisory Committee led by Michael Corm.

The U.K. Fed, eager.

Aboard all his most extreme concoctions careening the populace soon queued a shaky living net following weeks of the disastrous fiscal mess caused by the housing crash and the impending collapse in mortgage payments.

Can't anyone manage to usher it all the way. A savvy government could build "a serious living net - and that is much less than immediate relief." Instead of achieving a 50 percent reduction in mortgages , anyone at the Fed's table did it:

Listing 3: Generated texts from Roulette (0.95) diffusion trained using CEDD, L=128.

2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105

deeper problems caused by America's growing deficit, and by our rapidly growing economy. The same was true today when the only way to undermine the nation's welfare was in a virtuous cycle of austerity, stimulus and lavish budgets. And the scorching reality, especially for millions, was that "every man is the child of his parents' mother"-not only abhorrent problem was the working age in Washington, which only saw 44% of the population working in the same household as network leaders and 79% of the workforce in their industry. We tapped these network leaders in polls on Election Day, causing favorable results

=====

other federal government agencies had far greater control over water industry regulation. Today, drilling wasn't the equivalent of the Environmental Protection Agency championed since former president Ronald Reagan put Sen. Scott Desmond (CDE) on top of the Environmental Defense Agency, with a more funding area and rural programs under permit.

Like Moore, who spent so many a million dollars on his calls for the president to institute environmental protections on public lands and national parks, Trump himself had a nice shot on drilling, saying: "As passionate about the fight for property owners as we can be right here for Trump."

Canada XL - No

=====

Mistress: Next day next step, next step!" He's reply must not be translated. "What is - Guy's Your pronunciation? You still don't understand. Why am I saying that, don't mind: I've moved on to full set). Fitch word for move. Hey, that must have been long ago. What did you do going into our meeting? Everything changed, the rules changed... and you all admitted that Grand Auto Auto wasn't going to work, but you did. I need to know it. You've been in reference for the book, Notice."

Don't ignore him. He can't

=====

like low income earners, raising the tax on below.

All Part D credits forgiven by at least twenty-one percent paying income tax, and no state interest payments to help the health insurance guarantee. If owners pay the credits, they spend money they owe on each payment, either. A fifth of incomes will pay the credits.

While it's perfectly up to the type of hardship, people can apply state or federal assistance to the credit program. It requires people to sell health insurance. The credits qualify and vary from individual person to enrollee.

Having decided to pay their state income taxes on the bill, a single

Listing 4: Generated texts from Absorb diffusion trained using CEDD\*, L=128.

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

boat men: drift, and boom. Built, three 9th sailors found themselves losing an average of 12 inches of the underwater 50 feet in the delta base of S. M. Gielker Jr, U.S. captain (and pilot for the 11th Brigade), sailed deep into the canal on a watchdeck with his bow of thumb and lost at his anchor near shore.

Importantly, future "adillists should consider the alarming consequences that[t]here in the future the current rapprochement of the 9th would be had on the children of all areas of the coast , " re-ferenc

=====

the world's private space program." (This mission, in James's words-rust is laced and revisited to the down-space fall and the commercial foundation-and eventually-software portable to what Lotus's proprietary operating system does now-and also new hardware-serves, itself, as a collective organization rather than a fragmented, coordinated copal: "NASA loves the space community at large.")

James retired as CEO on the company board in 2004. And longer, Planetary Resources's judoided business model is intimately linked to his own organizational record and by

=====

engaged in hacking American government and stealing cyber defense secrets from its political enemies.

"We know today that the celebrity and his partner in the leaks incident to WikiLeaks produce the true story about our government, content and the data shared by our intelligence services, and deserve accountability," said Mark Warner, Assistant Attorney General, at the key E.C.C conference on Wednesday in Washington, R.I. "In this case, the battle #WikiLeaks against WikiLeaks has touched both at a new level for the country, and finally has turning historic date.

"Our government has been without the power to hack

=====

non-committal after the third GOP debate.

Gia had already just commented on the upcoming status of her project on CNN, after Trump repeatedly asserted it would discuss politics for the first lady.

Trump said, "Look at this, on my part of Hillary Clinton, the mad country," Garcia continued.

"I took my position in a way President Trump didn't have before speaking :], he said it. I admitted some of those things, finally understanding that I stood for Oh Hillary as a businesswoman," she continued.

Garcia made the comment after training herself in Melania Trump's Daily



Listing 5: Generated texts from Uniform diffusion trained using CEDD\*, L=128.

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

the Registration of the Bank of Development despite the skills of Arasan Sance. unemployable novice is in high regard for Swatch Up Worker as there are large scale projects bylaws that have already been difficult in the sector.

Caught diving too deeper, Arasan Sance had a career at IIT built on Year Up in Science at the Wood,, Theatre during an instructor class at the university in Banance in 2017. The guide says that some jobs may have been finalised and it was understood well that, if efficacious, McCarley could decide to be invited to a perhaps larger as academic class.

=====  
strength of subsequent fitness genes from mitochondria after maintenance of. This type of transition is called exalted mitochondrial, and although extended to long a feat is indeed amolecule strategy that requires stringent tests in selective binders.

In particular, one that abundant, modified mitochondria can enable recovery; first they produce mitochondria which is attached to the fuel cell, is rapidly moving essentially back to normal form, transitioning from the atrophy to a more muscular phenotype. Mitria such as L1L protein A (C1A and transporter R (Anaris-in proteinase), as NAD is used as an adaptation strategy, like many

=====  
how things were in the later years. The truck was littered with a black pepper. It was overflowing with halal green and mo yellow corn pepperoni and a charcoal plate.

Pictured used to have a baseballs encased in the front driveway with the sports hooders. Guy's black ores tip, black bean bag, silver, cherry, and granny banger was ricocheted, and when the Appomeess the Chinese bullfly rolled me over to the Morrissey. O 44, everything fresh and natural were just different from those in earlier, inscriptions. This guy's big white hooded or blue jacket

=====  
idea of increasing momentum, however, is to slow down a quarterback's strength. His is in the midst of a slow renaissance with the team's ensuing 17 games, but it is aging one who is getting the effectiveness of missing Greg Jennings and injuries. That this crucial area comes as the 49ers failed progressual Mac is forcing his fades down instead of sticking fully - unless he is getting better.

The Jets is approaching the top-10 in front but, as soon as they fill it in, an underdog-Dish seems to benefit from this slight fatigue in their upcoming game.

Negvious Matt Miller

Listing 6: Generated texts from Roulette (0.95) diffusion trained using CEDD\*, L=128.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

such fraudulent statements and denials as, "the truth lies," and "white lies" could never be trusted. More than one among you have just known Williams about.

Many are the fictional experiences Williams shaped, in ways that embodied him in a generation. The Tielemin figure in James Joyce's white character largely was based on male gaze and information and, from Williams' perspective, reflecting to the lens on men from a different time was disappointing.

Williams's material images of how, in essence, we had occupied a not present world: beaten, destroyed, marginalized or

=====

an attempt to fully realize the potential for pedestrian safety where there is no dearticulate conduct for the traffic.

His own current plans, for which Lokxelli is the secretary of the planning drove changes with little progress through traditional detours as of late. "Not highways, but others are up quite a bit," Akka forecasted "to make sure that one imp is not going to get torn. How to do that away from what it is becomes when much faster roads occur. There is all going to be of no serious consideration , at present, either in design or perspective."

=====

mad over and watching him laugh loudly enough to say so, and my point is that it's not just that I am not displeased with a politician or a CEO, that I think people looked in their shoes and paid well."

To see countries that see relevance to the show's multimillionaires, including so frustrating the small smile in the face of the world, Thomas Archifelius applied, but wanted to represent the image of man keeping ignorant out his protectionist mission.

ARTICLE CONTINUES BELOW

GEORGE STEPHANRAD: In a few weeks, Papacuzzi and so did Paul

=====

and made fun of the tactic. They acknowledged was just silly. Indeed, shortly after the books, a real Western painter and fellow broach type, said, to another wood miller, "he's straightened a nose. Skeptic got a hard chin and flies in the back." The winced to me is that you're a courtesan. You'll know what's up for you<|endoftext|>Get it done. Your friends will never see it

Parents of kids and old owners dogs yesterday said they are going to miss president-elect Donald Trump's look, 'Loving' for

Listing 7: Generated texts from Absorb diffusion trained using SEDD, L=128.

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

likely to be the fact that Nephi calls God's beloved "The Synod Word", a difference when it translates to an accident.

As viewers know, an Andesad man and his wife will never meet in Nephi's time God begins to connect with others, Le Burkert said.

Source: Bonuses de jeunee included<|endoftext|>The concept of attunational fit occurs from the presence of perceptual experiences on the cortex of the brain, as exposed recently by research by senior Prof. Linda Perriene.[1] Since personal dynamic patterns don't exist within human subjects, cortical memory is often cognitive

=====

will look like a target.

97) J.O. Do you think of as a White House national coordinator as a Clinton political operative? Kitty|O. I sort of're still up-the-table the White House. P.O. That person who refuses to have the same meek with a client with concerns/lives just pretty boycotted by pulling off in the relationship. Glad let's get that out of the way tonight and point out his ability and his job promise.

98) Next to the President - David Grimsnek | In at least 100 years George W. Bush has

=====

last month invested in Ria Louise, New York through a wealth of \$7 billion this year.

Meanwhile, the committee is a part of Katzman & Girle Service Conference, the firm that licenses the annual facility. Frisco visits to CraftEarlier this year.

"We newg CraftyTags for philanthropic, recovery, trust, and use of our services, with immense support from federal government, Mossody said.

Under the partnership, Neoneys also is working on expansion as a partner among other not-for-profit resource providers, with an increasing amount of resources and core services covering

=====

at least one, but nothing qualified from any perspective they could be offering them.

"The response was from a very recent comment made in blue material by the public but it was decided ahead of time," Smith writes of the public's prurient. So we've found that a long and full response excludes the ones who didn't believe the policy would meet the government's standards for proof and hold."

Sure enough, their apologists could do what they want on hard science. But that being said, they'd like some insights into scientific ethics and

## Listing 8: Generated texts from Uniform diffusion trained using SEDD, L=128.

2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375

can generate a json of the library based on which one exists for the scenario instance.

```
<?php G require JetlinkViewer :: create ( ePub : 'Whu )' d. add(): c.
app (). head () //Outputs from Sub_select from Google = str__ (
http://www.google.com ); //Add it while that function starts down
the path and results.Usingjson */ }
```

In order, Node handles repos(), several usable implementations in Mesh.js have been tested.

The JavaScript library actually has dire explanation for itls: it does not, but unlike

```
=====
using you or they use you taking people out to catch you or whatever.
That was wonderful because of the ways I look. I had gone from
anywhere from 20/22 here to 40 years old to 50/50. I also grew up
feeling. 'True Blood Baby' has gave me a lot of body movements
... that's got me right there and I used to pop their first
transformation. That show."
```

My sister really grew up as a nairyty boy. Thrones in a while? No, it's only sometimes, and since older people might even see us for 15 days or 30 days per night, how long they

```
=====
, but the ban on kidnapping remains the only difference.
```

Seyer has closed the case in a court letter seeking a response.

For more on the arrest video visit your Confused TMZ app page.

Updated Sept. 24, 2013. Tribute to Lawyer R Harrison.<|endoftext|>Gils Dahman is starting his season with the Red Bulls Crystal Palace. The midfielder rounded out his fifth season in the United States, missing three league meetings at U-a.

You can watch the interview with the media below. You can see an image for the camera below.

Gils Dahman has started

```
=====
Girl"). The earnestness of the ensemble is random, cheerful though
frustrated with both the actions of themselves and the different
perspectives on their journey. This classic masterpiece-writing
in sequels and his imagined form of movies-is the tale of cinema
that still has had aesthetics and is able to grow. But it also
sees the fall of the film genre-though very serious yet complex-
so we see a live show that can also find places lost by leading
the way to work around intersectionality, where the lineage of
Reed and Dormedy ("Iron Man" though) is applicable-and that she's
starting
```

Listing 9: Generated texts from Roulette (0.95) diffusion trained using SEDD, L=128.

2376  
2377  
2378  
2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389  
2390  
2391  
2392  
2393  
2394  
2395  
2396  
2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429

its own territory. All the European Union as well as the U.S. government, including the British East Company, were issued coins to the Federal Reserve and a few days were in the open to support a foreign buyer. President Lee had adopted law rejecting donations from the country's tax code, preserving the spirit of the Constitution of the 1870 Act of self-rule in South America. In the caddition of the British banking system, the American Mint provided the first ammunition and smuggled guns to Europe. Nonetheless, a serious problem: the Union was secretly armed anti-communists with a formidable internal authority. The union government was

=====

to blame?

It's full of wrongheaded sheen and bears expressions and carries a similar light ending. I'm not turned off on that Craig's Edge of Truth NYC guide if I don't see something wrong going on and do the right thing.

What's different about another story that can somehow tell the tale unfolds in Chicago so well with the ppl being approved by the General Service.

Like an Independent. Right? All right.

What Happens?

A Licensed Reader keeps all relationships private. They have no precedents. By all things

=====

result it couldn't publish it.

No, my work is not the best way that all of us can be at home with the poll results. I acknowledge you believing it, and doing this poll is going to have to happen to make the difference.

But again, we are honestly happy with the results - let's see. We could not vote from all the experts - but that still could happen. What's your opinion on that debate?

The other most annoying - but the most upset is #MSNBC. In the Wall Street Journal, people called the polls slow Gore down, saying that even the ABC thought

=====

, 400 times the base capacity there has been priced to four. And now another station can do the same for capacity with two other lines to 1 per cent of the power meter factory circuit which also costs to \$56 million for cancellation.

The next debate is on the schedule, but may not happen in late 2017 without major delays.<|endoftext|>The French media reports that Chinese technology firm EtApple has today started claiming that it has all given back, in its software, a return to earlier this month's SmartEasy payment processing that will simply payment it for most Chinese users. And based on recent developments, the same company is sending

Listing 10: Generated text from Absorb diffusion trained using CEDD\*, trained L=1024.

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483

two days later meeting with the Governmental Affairs Committee next Tuesday, replacing the committee's naps with an informal discussion.

In his Thursday hearing, Sessions said he had been "used" into criticizing the CIA.

"There has been some degree of transparency in looking at but that's not going to come cheap for himself or the American people, or really to undermine his independence, which I think, despite being in various positions, I've been very, very clear on the matter of that. I would be angry or very upset about that," he told senators at the meeting. "But I think as we go into Wednesday, and the United States Cabinet are finding ways to make very clear to citizens that they will respect that system."

Several high-level officials also called the appointment, saying it was "a bit of a quirk" since a revolving door of offices such as the CGA is already running rampant and has yet to hold any hearings.

"We've seen people make fairly much of their career life from the Justice Department out of this," said one GOP source. "And I don't think that choice is appropriate when that's a government entity that has given up a lot of accountability claims."

Speaking after noon another Republican member of the Senate Commerce Committee told. Todd Risch (R-Idaho) however, how he thought about Sessions, "and honestly, I didn't see a whole lot of discussion at one point suggesting a tone - a tone. There has been a lot of change since the removal of some five of these people ... I mean, I want people to find out that respect goes to the leadership of their department, and I'm not sure it's nice to move on down that precipice."<|endoftext|>When you try to put any pressure, you put more than the Pacers into you love watching. Before the Boston Celtics come to their games against Washington with.500 records on Tuesday night in Pittsburgh, you will likely watch a five-point game between Brad Stevens and Antetioun Gouden. Hollis have fourth-most competition in the game overall (49 points per game) but before games start their season of over 14.5 points per game has already passed.

That led Boston to an early 2-1 deficit.

No, not really. In addition to last season, however, the Celtics really struggled during the Wizards' first year of the NBA. For only the first 24 minutes, Washington had 18 minutes their D-man shooting guard who averaged at least four points, and as he chipped in another 11 points on at least 10 fourth quarter quarters Boston was outscored in the paint on five four-point shots, including two-point turnover margin.

If I included Stevens and Horford scoring 18 points in the only fourth quarter quarter from the paint, it would have probably been 2-1, and in the second-half-first quarter at half-time it would have been a 9-2 game thanks to the disastrous fourth-quarter performances. Yes, the Celtics have been unable to pull the ball away from us in the paint over the fourth quarter, but their success was probably in spite of the bad 2s, otherwise, our defense was weak.

Listing 11: Generated text from Absorb diffusion trained using SEDD, trained L=1024.

2484  
2485  
2486  
2487  
2488  
2489  
2490  
2491  
2492  
2493  
2494  
2495  
2496  
2497  
2498  
2499  
2500  
2501  
2502  
2503  
2504  
2505  
2506  
2507  
2508  
2509  
2510  
2511  
2512  
2513  
2514  
2515  
2516  
2517  
2518  
2519  
2520  
2521  
2522  
2523  
2524  
2525  
2526  
2527  
2528  
2529  
2530  
2531  
2532  
2533  
2534  
2535  
2536  
2537

.

The principle of selling for a period of energy will make the most sense next year, as Tesla's record rapid sales may be sluggish but demand lives on.

Tesla Motors Corp's CTO Carl Conti said beginning a new conversation with Tesla about the future state of its workforce in 2014 may be a part of its strategy in which as the company races to reduce its price of luxury vehicle deliveries, the first units that get to market again in earnest are expected to begin to get "next level," he said.

AP Photo

"I appreciate [the] ideas. These are not a fast break, but are ways to strengthen the margin at least knowing that if you get the right timeframe you make, its way ahead," Musk told reporters at The New York Times over the past month.

"None of these people have looked toward [s]ourcing the cost or the increase of the vehicle, unless you see it as a downward spiral," he said.

Although Musk and other executives point to one to two auto drivers and a net purchaser of gasoline or other, the company's strategy is that a reduced price for vehicles will outstrip some of those buyers, Conti still said.

Conti said the price of Model S is projected to fall by at least a 33 percent decline rate in 2015, on average. Conti expects Tesla expects 58.5 million to 59.5 million in 2015 and also for an average 15 percent decline to rise to February.

For the Model S, by contrast, Tesla's are expected to reach a write-down of about \$500 million to 2,000,000 EVs in September and more than 700,000 currently, he said.

The first Model deliveries of the year is expected to nearly begin on March 22. That price is heavily influenced by dealer monitoring that includes a total inventory of 600,000, with 400,000 delivered next month, Conti said.

He noted Tesla expects its predicted upward trajectory to increase by 2.4 percent by 2015 as production in Australia and China accelerate to 4.6 percent next year, respectively.

Musk also noted in The Wall Street Journal that while Tesla expects more than three consecutive in production next year, it will production a total of \$8 million in U.S. roll-backs in the fall.

He said he has settled on a key question about the car: "If a customer needs it and is willing to pay the rates that we can keep things moving, then we ought to consider adding incentives for it."

"The price of Model S expectations for production is a driver'all," Conti said. "When the profit is available, it means a car is not in its pocket. We are seeing this as a very good time after Jan. 31."

## Listing 12: Generated text from GPT2, trained L=1024.

2538  
2539  
2540  
2541  
2542  
2543  
2544  
2545  
2546  
2547  
2548  
2549  
2550  
2551  
2552  
2553  
2554  
2555  
2556  
2557  
2558  
2559  
2560  
2561  
2562  
2563  
2564  
2565  
2566  
2567  
2568  
2569  
2570  
2571  
2572  
2573  
2574  
2575  
2576  
2577  
2578  
2579  
2580  
2581  
2582  
2583  
2584  
2585  
2586  
2587  
2588  
2589  
2590  
2591

power well past the NPC versus NPC meter for each unit type. Their use is varied, and thus allows for a considerable amount of flexibility with The Banner Saga, but whether you prefer your parent's roots to stay part of the community -- with allies everyone can fight and help each other -- remains to be seen. For that, you can always pick the steam community tab (not required, that's how the game is named -- no one likes spamming themselves for news now, I know).

However, you agree to terms of service that you agree to disable your favorite kind of multiplayer games (by default) to prevent unwanted conflicts: factions are only listed once in a game and cannot be disbanded or changed just because of a non-friendlier fleet. So it's a pretty straightforward, fairly simple toggle to get used to.

Ugh, it sounds like the HUD might get progressively messier as the game ends.

Shields charging slowly: active shields last several minutes, recharge fully in seconds, and charge every second they charge.

Shields wearing 'shield' else doubles their size, can clip high above their skull, and are invincible to all other shields around your ship that border the shield and can deploy wall-mounted shields. Other shields like this one can only be used to shield the next vessel. Hit someone on the Pustules to activate shield shields. In an expedited, non-lethal, coordinated confrontation, it is fairly easy to compare stats to each shield for a ship to use.

Shield minus shields compensation: up to five shields are held at rank 2 in your ship's shields. It's called 'delay,' but with a ton of fancy commands on screen, you don't have to know (or even see) their context or employment to make sure to compensate for some of the distortion inherent in that bit. Seconds before your shield starts charging (just seconds before it behaves like hiring another worker), the timer will start running indefinitely, which means shields can snap to the back of your ship if kept as low as possible. Avoid collisions caused by instantaneous cap.

Shield halves that match the length of the ship's hull: you can control the range between Shields, such as your Private Capital Signatures, where the ship's bulletproof liquor clock is set at half-space and your Heavy Construction Band Bump the Game's Game soundtrack to Fade let the corp count (if the corporation has sixty transports within one day of being prepped, it's probably already pretty close to ten minutes, but your mission selection loading speeds will be holding it back a little).

Surrounding explosions moving over the bow: it's deceptively easy to line up targets like an Infallible Surface projection theater Explosion shield room high up in the sky: fire teams can run together and mean to demolish their own low-flying illions every once in a while, yet closely follow their targets like shadow warriors to the detriment of themselves. Covertly set to 'square hit' for easier customization, Ubisoft's new ship weather HUD is the key to spreading effective damage throughout the entire ship using the standard meteorological updates.

Shaping your space: shape each bridge to create low-risk rockouts or pathogen/addage paths. This may be tedious, but it actually whizzes through by itself while simultaneously adding damage to giants mercenary ships.



## 2592 D CTMCs PRELIMINARIES

### 2593 D.1 DISCRETE-TIME MARKOV CHAINS OVER FINITE-STATE SPACES

2594 A discrete-time Markov Chain in a finite-state space is a stochastic process  $X_1, X_2, \dots, X_T$ , where  
 2595 each state  $X_t$  depends solely on the preceding state  $X_{t-1}$ . The states  $X_t$  can take on any value  
 2596 from the set  $\{1, 2, \dots, S\}$ , where  $S$  denotes the total number of possible states, and  $T$  represents the  
 2597 number of time steps. The probability of being in state  $x$  at time  $t$  is

$$2600 p_t(X_t = x) = \sum_{y=1}^S p(X_t = x, X_{t-1} = y) = \sum_{y=1}^S p_{t|t-1}(X_t = x|X_{t-1} = y)p_{t-1}(X_{t-1} = y). \quad (175)$$

2601 If we place all such probabilities  $p_t(X_t = x)$  in a vector  $\mathbf{s}_t$  of shape  $S \times 1$ , such that  $\mathbf{s}_t(x) =$   
 2602  $p_t(X_t = x)$ , then from above we can deduce that

$$2603 \mathbf{s}_t = \mathbf{P}\mathbf{s}_{t-1}, \quad (176)$$

2604 where  $\mathbf{P}(x, y) = p_{t|t-1}(X_t = x|X_{t-1} = y)$ . Given an initial probability distribution  $\mathbf{s}_0$  over states,  
 2605 the equation above fully determines the evolution of the probability over states with respect to time.  
 2606 If it is known that the state at time  $t - 1$  is  $y$ , then  $\mathbf{s}_t$  is simply column  $y$  of  $\mathbf{P}$ . This implies that the  
 2607 sum of the elements of each column  $y$  of  $\mathbf{P}$  is one.

### 2612 D.2 CONTINUOUS-TIME MARKOV CHAINS OVER FINITE-STATE SPACES

2613 It is possible to define a stochastic process with the Markov property in finite-state spaces, for  
 2614  $t \in [0, T]$ , (Anderson, 2012). As previously, we can define a discrete-time process, on time points  
 2615  $\{0, \epsilon, \dots, T - \epsilon, T\}$ , such that there is  $\epsilon$  probability of activating the previous transition mechanism  
 2616 when progressing from time  $t - \epsilon$  to  $t$ , otherwise we stay where we are with probability  $(1 - \epsilon)$ .  
 2617 Removing the random variables to simplify notation, we have

$$2618 p_t(x) = (1 - \epsilon)p_{t-\epsilon}(x) + \epsilon \sum_{y=1}^S p_{t|t-\epsilon}(x|y)p_{t-\epsilon}(y). \quad (177)$$

2619 We notice that when  $\epsilon = 1$  the equation above coincides with Equation (175), and in addition as  
 2620 before we can write Equation (177) in matrix form

$$2621 \mathbf{s}_t = (1 - \epsilon)\mathbf{s}_{t-\epsilon} + \epsilon\mathbf{P}\mathbf{s}_{t-\epsilon} = (\mathbf{I} + \epsilon(\mathbf{P} - \mathbf{I}))\mathbf{s}_{t-\epsilon} = (\mathbf{I} + \epsilon\mathbf{Q})\mathbf{s}_{t-\epsilon}, \text{ where } \mathbf{Q} = \mathbf{P} - \mathbf{I}. \quad (178)$$

2622 From Equation (178), we see that  $\frac{\mathbf{s}_t - \mathbf{s}_{t-\epsilon}}{\epsilon} = \mathbf{Q}\mathbf{s}_{t-\epsilon}$ , which when taking the limit  $\epsilon \rightarrow 0$  becomes

$$2623 \frac{d\mathbf{s}_t}{dt} = \mathbf{Q}\mathbf{s}_t. \quad (179)$$

2624 Given an initial probability distribution  $\mathbf{s}_0$  over states, the equation above fully determines the evo-  
 2625 lution of the probability over states with respect to time. Indeed, the distribution over states at time  
 2626  $t$  is the solution of the linear ODE in (179):  $\mathbf{s}_t = e^{t\mathbf{Q}}\mathbf{s}_0$ . Therefore, if it is known that the state  
 2627 at time 0 is  $j$ , then  $\mathbf{s}_t$  is simply column  $j$  of  $e^{t\mathbf{Q}}$ . Naturally, one can rescale the time variable such  
 2628 that  $t = \sigma(t)$ , where  $\sigma$  is monotonically increasing,  $\sigma(0) = 0$  and  $\lim_{t \rightarrow 1} \sigma(t) = T$ , so that for  
 2629  $\mathbf{z}_0 := \mathbf{s}_0$ , we get  $e^{t\mathbf{Q}}\mathbf{s}_0 = e^{\sigma(t)\mathbf{Q}}\mathbf{z}_0 = \mathbf{z}_t$ , and

$$2630 \frac{d\mathbf{z}_t}{dt} = \sigma'(t)\mathbf{Q}\mathbf{z}_t = \mathbf{Q}_t\mathbf{z}_t, \text{ where } \mathbf{Q}_t = \sigma'(t)\mathbf{Q}. \quad (180)$$

2631 Matrices  $\mathbf{Q}_t = \sigma'(t)(\mathbf{P} - \mathbf{I})$  clearly satisfy the properties of transition-rate matrices (Suhov &  
 2632 Kelbert, 2008). Matrices  $\mathbf{Q}_t$  are chosen such that: a) the matrix exponential  $e^{\sigma(t)\mathbf{Q}}$  is easy to  
 2633 calculate, which is essential as  $p_{t|0}(x|y) = e^{\sigma(t)\mathbf{Q}}(x, y)$ ; and b)  $\mathbf{z}_1$  is an easy reference distribution  
 2634 to sample from (Austin et al., 2021; Campbell et al., 2022).

2635 Finally, similar to diffusion processes in continuous spaces, the continuous-time Markov chain in  
 2636 Equation (180) also admits a reverse process (Kelly, 1979; Sun et al., 2022):

$$2637 \frac{d\mathbf{z}_{1-t}}{dt} = \bar{\mathbf{Q}}_{1-t}\mathbf{z}_{1-t}, \quad (181)$$

2646 where  $\bar{Q}_t(x, y) = Q_t(y, x) \frac{p_t(x)}{p_t(y)}$  for  $x \neq y$ , and  $\bar{Q}_t(x, x) = -\sum_{y \neq x} \bar{Q}_t(y, x)$ . Since we can  
2647 easily sample from the reference distribution, the only unknowns preventing us from being able to  
2648 run backwards are the ratios  $\frac{p_t(x)}{p_t(y)}$  also known as concrete scores (Meng et al., 2022), which we  
2649 desire to model using a neural network. Once such ratios are modeled we can generate samples  
2650 from the learned data distribution  $p_0^\theta$  by discretizing Equation (181) as follows:  
2651

$$2652 \quad p(x_{t-\epsilon} = y \mid x_t = x) = \delta_x(y) + \bar{Q}_t(y, x)\epsilon + O(\epsilon^2). \quad (182)$$

2653  
2654  
2655  
2656  
2657  
2658  
2659  
2660  
2661  
2662  
2663  
2664  
2665  
2666  
2667  
2668  
2669  
2670  
2671  
2672  
2673  
2674  
2675  
2676  
2677  
2678  
2679  
2680  
2681  
2682  
2683  
2684  
2685  
2686  
2687  
2688  
2689  
2690  
2691  
2692  
2693  
2694  
2695  
2696  
2697  
2698  
2699