

A Additional details regarding D3PMs

A.1 Doubly-stochastic matrices

As discussed in Section 3.1, there are two constraints on \mathbf{Q}_t that allow it to be used within a D3PM: the rows of \mathbf{Q}_t must sum to one to conserve probability mass, and the rows of $\bar{\mathbf{Q}}_t = \mathbf{Q}_1 \mathbf{Q}_2 \dots \mathbf{Q}_t$ must converge to a known stationary distribution as t becomes large. Technically, it is also possible to use a learned prior $p_\theta(\mathbf{x}_T)$, but assuming this is still modeled under a conditional independence assumption, $q(\mathbf{x}_T | \mathbf{x}_0)$ must still be close to a stationary distribution for the L_T loss term to be small.

One way to ensure that this occurs is to choose \mathbf{Q}_t as increasing powers of a doubly stochastic base matrix \mathbf{Q} (rows and columns sum to 1) with strictly positive entries. This is enough to ensure that \mathbf{Q} is irreducible and aperiodic and that product $\bar{\mathbf{Q}}_t$ converges as $t \rightarrow \infty$ to a uniform distribution over all states. To show this, consider $\pi_i = 1/K$ for $i = 1, \dots, K$, and $\sum_{i=1}^K \mathbf{Q}_{i,:} = \mathbf{1}$ and $\sum_{j=1}^K \mathbf{Q}_{:,j} = \mathbf{1}$, then $[\mathbf{Q}\pi]_i = \sum_{j=1}^K \mathbf{Q}_{i,j} \pi_j = 1/K \sum_{j=1}^K \mathbf{Q}_{i,j} = 1/K = \pi_i$, thus the uniform distribution is an eigenvector of the transition matrix with eigenvalue 1. Convergence to this distribution follows from the Perron-Frobenius theorem for positive square matrices.

More generally, a similar argument shows that even for \mathbf{Q}_t that are not powers of the same base matrix, as long as each \mathbf{Q}_t is doubly stochastic, irreducible, and aperiodic, the uniform distribution is the only possible stationary distribution, and as long as the second largest eigenvalue of \mathbf{Q}_t is bounded below, the cumulative product $\bar{\mathbf{Q}}_t$ will converge to the uniform distribution. In practice, we choose \mathbf{Q}_t to add more noise as t increases, which ensures that $\bar{\mathbf{Q}}_T$ is very close to reaching a uniform stationary distribution.

A.2 More details on possible choices of Markov transition matrices

A.2.1 Uniform diffusion

The transition matrix described by Sohl-Dickstein et al. [17] for the binary case, and extended by Hooeboom et al. [9], to the categorical case, can be represented using the following $K \times K$ transition matrix

$$[\mathbf{Q}_t]_{ij} = \begin{cases} 1 - \frac{K-1}{K} \beta_t & \text{if } i = j \\ \frac{1}{K} \beta_t & \text{if } i \neq j \end{cases}, \quad (6)$$

This transition matrix can also be written as $(1 - \beta_t)I + \beta_t \mathbf{1} \mathbf{1}^T / K$, where $\mathbf{1}$ is a column vector of all ones.

A.2.2 Diffusion with an absorbing state

For our diffusion models with an absorbing state m , we use the following matrix:

$$[\mathbf{Q}_t]_{ij} = \begin{cases} 1 & \text{if } i = j = m \\ 1 - \beta_t & \text{if } i = j \neq m \\ \beta_t & \text{if } j = m, i \neq m \end{cases} \quad (7)$$

The transition matrix can also be written as $(1 - \beta_t)I + \beta_t \mathbf{1} e_m^T$, where e_m is a vector with a one on the absorbing state m and zeros elsewhere. Since m is an absorbing state, the corruption process converges not to a uniform distribution but to the point-mass distribution on m .

For text generation, we let m be the [MASK] token at index $K - 1$; this leads to a BERT-like training objective, which masks tokens according to some schedule and learns to denoise them iteratively (see Section 4). For image generation, we set m to the gray RGB pixel (128, 128, 128) at index $K/2$.

A.2.3 Discretized Gaussian transition matrices

For our D3PM models applied to ordinal data, inspired by continuous-space diffusion models, we use the following $K \times K$ matrix:

$$[\mathbf{Q}_t]_{ij} = \begin{cases} \frac{\exp\left(-\frac{4|i-j|^2}{(K-1)^2\beta_t}\right)}{\sum_{n=-(K-1)}^{K-1} \exp\left(-\frac{4n^2}{(K-1)^2\beta_t}\right)} & \text{if } i \neq j \\ 1 - \sum_{l=0, l \neq i}^{K-1} [\mathbf{Q}_t]_{il} & \text{if } i = j \end{cases} \quad (8)$$

Normalization is ensured by assigning the diagonal values to one minus the sum of each row (not including the diagonal entry). Note that due to the normalization of the off-diagonal values over the range $\{-K+1, \dots, K-1\}$ the sum of each row excluding the diagonal entry is always smaller than 1. The result yields an irreducible doubly stochastic matrix and a forward process with a uniform stationary distribution. Similar to the continuous Gaussian diffusion model, the parameters β_t influence the variance of the forward process distributions.

A.2.4 Structured diffusion in text: using word-embedding distance to introduce locality

For text, we construct a k -nearest neighbor adjacency matrix

$$[\mathbf{G}]_{ij} = 1 \text{ if } w_i \text{ is a } k\text{-nearest neighbor of } w_j \text{ else } 0$$

constructed from a pre-trained embedding space over the vocabulary. Then we consider a symmetrized adjacency matrix of the form $\mathbf{A} = (\mathbf{G} + \mathbf{G}^T)/(2k)$ where k is the number of nearest neighbors of each node, and finally construct a doubly stochastic rate matrix with

$$[\mathbf{R}]_{ij} = \begin{cases} -\sum_{l \neq i} A_{il} & \text{if } i = j \\ A_{ij} & \text{otherwise} \end{cases} \quad (9)$$

Our final transition matrix is constructed as a matrix exponential of this rate matrix:

$$\mathbf{Q}_t = \exp(\alpha_t \mathbf{R}) = \sum_{n=0}^{\infty} \frac{\alpha_t^n}{n!} \mathbf{R}^n$$

Since \mathbf{R} is symmetric and sums to zero along each row, \mathbf{Q}_t is doubly stochastic, which ensures we have a uniform stationary distribution (as long as G is connected). Increasing α_t over time allows us to add more noise for larger values of t .

Assuming word embeddings are some metric for syntactic or semantic similarity, this results in a corruption process that gradually moves away from the ground-truth sentence, swapping words with nearest-neighbors in embedding space. For character level modeling, this is a graph over characters, which more often transitions for instance from vowels to other vowels than from vowels to consonants. For words, this could transition between semantically similar words.

For example, in Figure 4, we construct the forward process to diffuse from "dog" to "cat" or "cow", which are nearby in embedding space, but not to more distant words. We can either bootstrap this process by updating the transition matrix \mathbf{Q} dynamically during training, or use pretrained embeddings; we use pretrained embeddings for all of our experiments. Specifically, we train an autoregressive language model on the dataset in question (either text8 or LM1B) with randomly initialized word embeddings (768 dimensional in most cases), and then use L^2 or cosine similarity to compute the k -nearest neighbors of each token. We transition preferentially to these tokens, although the matrix exponential in theory allows transitions to any other token. We choose k large enough so the resulting graph is connected.

A.2.5 Band-diagonal transitions

A class of transition matrices that introduce local, ordinal inductive biases for structured data are band-diagonal transition matrices which only allow the corruption process to transition locally between states and biases the reverse process towards local iterative refinement. For example, in images, this can be used to allow transitions only between adjacent pixel values.

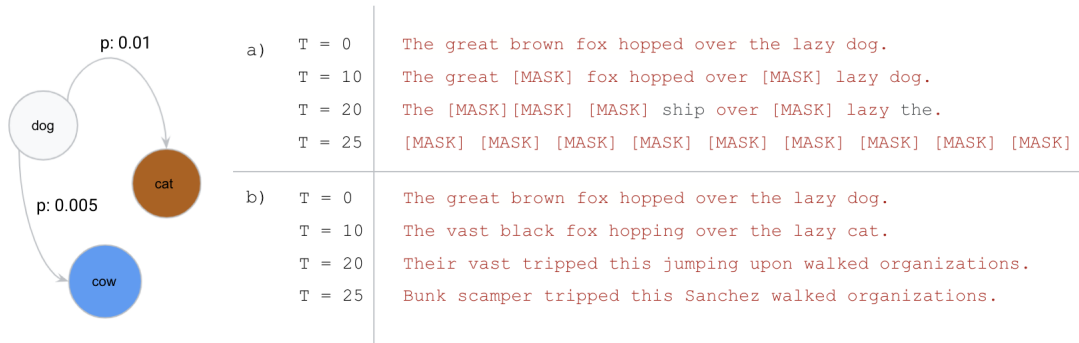


Figure 4: Two examples of noise schedules transforming text data. The top is a BERT-like absorbing + uniform diffusion which replaces tokens with [MASK] tokens (and occasionally with any other token, in black). The bottom is nearest-neighbor diffusion in embedding space. At left represents a possible column in the transition matrix.

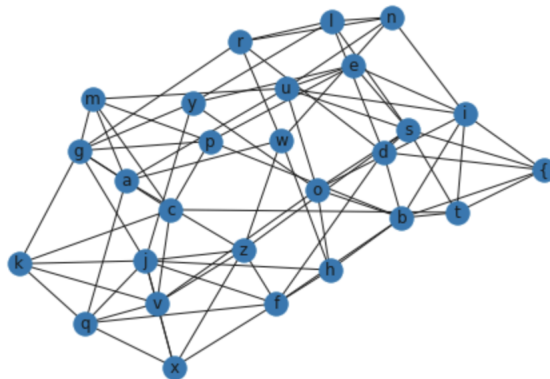


Figure 5: The character-level symmetrized 5-NN graph.

$$[Q_t]_{ij} = \begin{cases} \frac{1}{K}\beta_t & \text{if } 0 < |i - j| \leq v \\ 1 - \sum_{l \neq i} Q_{il} & \text{if } i = j \end{cases} \quad (10)$$

where v is the number of nonzero off-diagonal elements of Q above (and below) the main diagonal. Note that this is a doubly stochastic matrix, so the stationary distribution is uniform. We do not use these in our experiments.

A.2.6 Combinations of absorbing diffusion and other diffusion

A few ablations in Appendix B.2.1 consider transition matrices that combine absorbing-state or nearest-neighbor and uniform D3PM models. For instance, an absorbing-uniform transition matrix can be constructed $Q = \alpha \mathbb{1}e_m^T + \beta \mathbb{1}\mathbb{1}^T/K + (1 - \alpha - \beta)I$, where e_m is a one-hot vector on the [MASK] token.

A.3 Generative Masked Language Models are Diffusion Models

Generative Masked Language Models [5, 21] are generative models that generate text from a sequence of [MASK] tokens. These are usually trained by sampling a sequence x_0 , masking tokens according to some schedule, and learning to predict the masked tokens given context. The actual masking procedure can either be done independently, i.e. by masking each token with probability $p = k/T$,

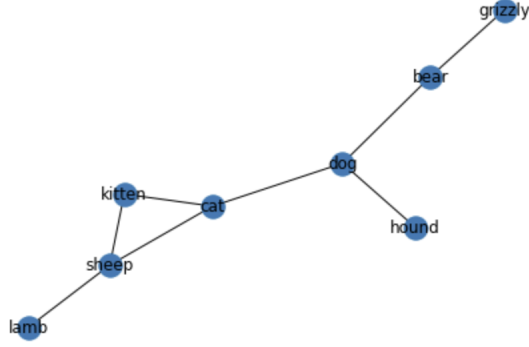


Figure 6: Subgraph of a word-level NN graph.

like Devlin et al. [3], or by sampling exactly k tokens. The usual objective is⁷:

$$\min -\mathbb{E}_{q(\mathbf{x}_0)} \left[\mathbb{E}_{k \in [1 \dots |\mathbf{x}_0|]} \left[\frac{1}{k} \mathbb{E}_{\mathbf{x}_k \text{ with } k \text{ masked tokens}} \left[\sum_{i \text{ with } [\mathbf{x}_k]_i = m} \log p_\theta([\mathbf{x}_0]_i | \mathbf{x}_k) \right] \right] \right] \quad (11)$$

where we first sample a datapoint \mathbf{x}_0 , sample a number of tokens to mask k (either uniformly or according to some schedule), then mask that many tokens at random and compute a cross entropy loss over those masked tokens. We claim that this training objective is a (reweighted) absorbing-state D3PM objective with a particular noise schedule and the \mathbf{x}_0 -parameterization from 3.3 (and indeed, that any absorbing-state D3PM model with [MASK] as the absorbing state will be a reweighted version of this loss with different weights assigned to different numbers of masked tokens k).

Consider a D3PM with a schedule that masks tokens with probability β_t . The reverse process predicts $\tilde{p}_\theta(\tilde{\mathbf{x}}_0 | \mathbf{x}_t)$, then uses the forward process to compute $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) \propto \sum q(\mathbf{x}_{t-1}, \mathbf{x}_t | \tilde{\mathbf{x}}_0) \tilde{p}_\theta(\tilde{\mathbf{x}}_0 | \mathbf{x}_t)$. In the particular case of absorbing-state diffusion, for each masked token $[\mathbf{x}_t]_i = m$ in \mathbf{x}_t , we thus have

$$p_\theta([\mathbf{x}_{t-1}]_i | \mathbf{x}_t) \propto \begin{cases} [\beta_t \prod_{s < t} (1 - \beta_s)] \tilde{p}_\theta([\tilde{\mathbf{x}}_0]_i = [\mathbf{x}_0]_i | \mathbf{x}_t) & \text{for } [\mathbf{x}_{t-1}]_i = [\mathbf{x}_0]_i \neq m \\ 1 - \prod_{s \leq t} (1 - \beta_s) & \text{for } [\mathbf{x}_{t-1}]_i = m \end{cases}$$

We note that for each unmasked token $[\mathbf{x}_t]_i = [\mathbf{x}_0]_i$, the KL-divergence is zero since unmasked tokens cannot make any other type of transition other than becoming masked. Also, the term in the KL divergence due to the probability of mask transitions is a constant, since mask transitions are independent of the model parameters θ . Our L_t term is then

$$D_{\text{KL}}[q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)] = - \left[\beta_t \prod_{s < t} (1 - \beta_s) \right] \sum_{i \text{ with } [\mathbf{x}_t]_i = m} \log \tilde{p}_\theta([\mathbf{x}_0]_i | \mathbf{x}_t) + C$$

where C is independent of θ and the sum is taken over the masked tokens in \mathbf{x}_t . For example, if we use $\beta(t) = 1/(T - t + 1)$ from Sohl-Dickstein et al. [17], $\beta_t \prod_{i=0}^{t-1} (1 - \beta_i) = 1/T$ and $1 - \prod_{i=0}^t (1 - \beta_i) = (t - 1)/T$, so $q([\mathbf{x}_{t-1}]_i = [\mathbf{x}_0]_i | [\mathbf{x}_t]_i = m, \mathbf{x}_0) = 1/t$ for non-mask tokens and we can simplify our L_t objective to

$$D_{\text{KL}}[q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)] = - \left[\frac{1}{t} \sum_{i \text{ with } [\mathbf{x}_t]_i = m} \log \tilde{p}_\theta([\mathbf{x}_0]_i | \mathbf{x}_t) \right] + C$$

where \mathbf{x}_t masks tokens independently and uniformly with probability t/T . The L_T term in our ELBO is 0 for the $1/(T - t + 1)$ schedule, so the full objective (up to a constant) reduces to

⁷Sometimes the loss is un-normalized or normalized by the full sequence length.

$$\begin{aligned}
& \mathbb{E}_{q(\mathbf{x}_0)} \left[- \sum_{t=2}^T \frac{1}{t} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\sum_{i \text{ with } [\mathbf{x}_t]_i=m} \log p_\theta([\mathbf{x}_0]_i|\mathbf{x}_t) \right] \right. \\
& \quad \left. - \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} \left[\sum_{i \text{ with } [\mathbf{x}_1]_i=m} \log p_\theta([\mathbf{x}_0]_i|\mathbf{x}_1) \right] \right] \\
&= - \mathbb{E}_{q(\mathbf{x}_0)} \left[\sum_{t=1}^T \frac{1}{t} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[\sum_{i \text{ with } [\mathbf{x}_t]_i=m} \log p_\theta([\mathbf{x}_0]_i|\mathbf{x}_t) \right] \right] \tag{12}
\end{aligned}$$

Note that while this looks very similar to Equation [11](#) (with each term reweighted by $1/t$, the expected number of masked tokens) it is not exactly identical since masking is computed independently per-token position (instead of choosing exactly k tokens to mask). This is an entirely practical way to do masking (and indeed some methods implement it this way).

Furthermore, since the masking probability varies linearly as $1 - \prod(1 - \beta_t) = t/T$, this is very close to uniformly sampling the number of masked tokens k , but k is actually drawn from a mixture of binomial distributions, i.e.

$$= - \mathbb{E}_{q(\mathbf{x}_0)} \left[\mathbb{E}_{k \in [1 \dots |X|]} \left[\mathbb{E}_{\mathbf{x}_k \text{ with } k \text{ masked tokens}} \left[\alpha(k) \sum_{i \text{ with } [\mathbf{x}_k]_i=m} \log p_\theta([\mathbf{x}_0]_i|\mathbf{x}_k) \right] \right] \right] \tag{13}$$

$$\alpha(k) = q(\mathbf{x}_t \text{ has } k \text{ masked tokens} | \mathbf{x}_0 \text{ has } n \text{ tokens}) = \frac{1}{T} \sum_{t=1}^T \binom{n}{k} \left(\frac{t}{T}\right)^{n-1} \left(1 - \frac{t}{T}\right)^{n-k} \tag{14}$$

which is very close to uniform weight over terms, but slightly downweights terms near 0 and T . By upweighting terms near the boundary, you could in theory make this exactly uniform and thus exactly recover Equation [11](#). For instance, for 50 categories, absorbing-state diffusion produces the weighting shown in Figure [7](#).

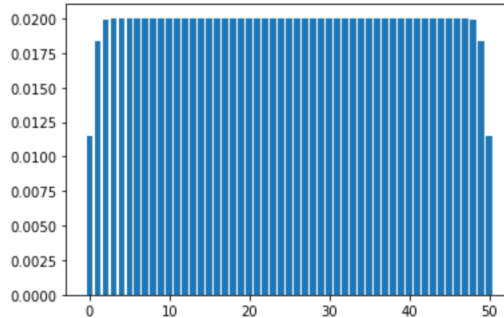


Figure 7: Plot of the probabilities of having k tokens masked out of a length-50 sequence under a D3PM absorbing schedule with $T = 50$ steps, which is very similar to the uniform weighting used by Ghazvininejad et al. [\[5\]](#).

A.4 Scaling to a large number of categories

When the number of categories K is large, it can quickly become impractical to store all of the transition matrices \mathbf{Q}_t in memory, as the memory usage grows like $O(K^2T)$. And even if there is an algorithm to compute individual step matrices \mathbf{Q}_t on demand, it may or may not be possible to do the same for the cumulative products \mathbf{Q}_t . We propose two approaches to scaling D3PMs to large numbers of categories that ensure cumulative products are efficient: using low-rank corruption and using matrix exponentials.

A.4.1 Low-rank corruption

In the low-rank case, we consider structuring our transition matrices as

$$\mathbf{Q}_t = \beta_t \mathbf{A}_t + (1 - \beta_t) \mathbf{I}, \quad (15)$$

where each \mathbf{A}_t is a diagonalizable low-rank matrix with the same nonzero eigenvectors. In particular, recall that both absorbing-state diffusion and uniform diffusion have this form: for uniform diffusion, $\mathbf{A}_t^{\text{uniform}} = \mathbb{1}\mathbb{1}^T/K$, and for absorbing-state diffusion $\mathbf{A}_t^{\text{abs}} = \mathbb{1}e_m^T$ where e_m is a one-hot vector on the absorbing state. Since products of \mathbf{A}_t 's are also low rank, the cumulative products $\overline{\mathbf{Q}}_t$ can be efficiently precomputed and stored using a much smaller amount of memory $O(r^2T)$ where $r = \text{rank}(\mathbf{A}_t)$.

As an illustrative example, we describe in more detail how to efficiently represent uniform and absorbing-state transition matrices using the low-rank structure.

To compute products of uniform transition matrices (i.e. $\prod_i (1 - \beta_i) \mathbf{I} + \beta_i \mathbb{1}\mathbb{1}^T/K$), we can take advantage of the useful fact that products of matrices of the form $\alpha \mathbf{I} + \beta \mathbb{1}\mathbb{1}^T$ also have this same form: $I^2 = I$ and $(\beta \mathbb{1}\mathbb{1}^T)^2 = \beta^2 K \mathbb{1}\mathbb{1}^T$. We can thus treat this as a formal polynomial in one variable $X = (\mathbb{1}\mathbb{1}^T/K)$. Then products can be computed as $\prod_i [(1 - \beta_i) + \beta_i X]$ over the quotient ring $\mathbb{R}[X]/(X^2 - X)$, since $X^2 = X$. Functionally, this means you can instantiate a polynomial $(1 - \beta_i) + \beta_i X$ and repeatedly perform ordinary polynomial multiplication over $\mathbb{R}[X]$ for the $t < T$ timesteps. After each multiplication, the higher-order terms are reduced by $X^2 = X$, leaving a polynomial of degree 1 where the X term has coefficient given by the sum of all higher-order terms. This can be computed with the convenient `np.polynomial` module.

Similarly, the transition matrices for D3PM absorbing can be computed in closed form. Fundamentally, in each step, we transition to a [MASK] token with probability β_t and stay the same with probability $1 - \beta_t$. Since the [MASK] state is absorbing, after t steps, the only operative quantity is the probability of not yet having transitioned to the [MASK] state, given by $\tilde{\alpha}_t = \prod_{i=0}^t (1 - \beta_i)$. Hence for D3PM absorbing, $\overline{\mathbf{Q}} = \tilde{\alpha}_t \mathbf{I} + (1 - \tilde{\alpha}_t) \mathbb{1}e_m^T$ where e_m is a one-hot vector on the [MASK] token.

A.4.2 Matrix exponentials

In the matrix exponential case, we specify our transition matrices as

$$\mathbf{Q}_t = \exp(\alpha_t \mathbf{R}) = \sum_{n=0}^{\infty} \frac{\alpha_t^n}{n!} \mathbf{R}^n, \quad \overline{\mathbf{Q}}_t = \exp\left(\left(\sum_{s \leq t} \alpha_s\right) \mathbf{R}\right), \quad (16)$$

where \mathbf{R} is a *transition rate matrix* and \exp denotes the matrix exponential operation; the similar form for \mathbf{Q}_t and $\overline{\mathbf{Q}}_t$ is a consequence of the ‘‘exponential of sums’’ property for commuting matrices. For efficiency, we further assume that each of the α_t is an integer multiple $n_t \alpha_*$ of some common factor α_* , and precompute matrices $\exp(2^k \alpha_* \mathbf{R})$ for $0 \leq k \leq \log_2(\bar{\alpha}_T/\alpha_*)$, where $\bar{\alpha}_T = \sum_{t < T} \alpha_t$, taking space $O(K^2 \log(\bar{\alpha}_T/\alpha_*))$. Then, to compute matrix-vector products with \mathbf{Q}_t or $\overline{\mathbf{Q}}_t$, we can iteratively take products with a subset of these precomputed matrices based on the digits of a binary expansion of the desired multiple n_t in time $O(K^2 \log(\bar{\alpha}_T/\alpha_*))$.⁸

As long as \mathbf{R} has non-positive off-diagonal entries and sums to zero along each row, the matrix exponential produces a valid transition matrix \mathbf{Q}_t ; convergence to a specific stationary distribution can also be ensured by controlling the eigenvectors. In particular, if every column also sums to zero, the resulting \mathbf{Q}_t will be doubly stochastic and will thus have a uniform stationary distribution.

We note that this parameterization can be viewed as a discretization of a continuous-time discrete-space Markov processes; we describe this connection in more detail in the following section.

A.5 Continuous-time Markov process transition rates

Following Feller [4], we define a continuous-time discrete-space Markov process as a collection of random variables $\{\mathbf{x}_t\}_{t > 0}$ parameterized by $t \in \mathbb{R}^+$ and characterized by a Markov property

⁸This is closely related to the well-known ‘‘exponentiation-by-squaring’’ technique.

$(\mathbf{x}_t \perp \mathbf{x}_s \mid \mathbf{x}_\tau$ if $t < \tau < s$), a transition probability matrix $\Pi(t) \in \mathbb{R}^{N \times N}$ where N is the cardinality of \mathbf{x}_t , and a set of transition rates $\gamma_i(t)$.

A conceptual way to understand these processes is to imagine a continuous Poisson process occurring in each state i at rate $\gamma_i(t)$ determining when a transition between states occurs. When a transition occurs (at time t), a Markov transition occurs between states i and j with probability $\Pi_{ij}(t)$. Many common stochastic processes fall into this family, including Poisson processes. Like in the case of stochastic differential equations (Song et al. [18]), we can derive a set of Kolmogorov equations (or Fokker-Planck equations in the continuous-state space case) that determine the marginal probability $\partial q_{ij}(\tau, t)$ of ending up in state j at time t having started in state i at time s . The general form of the Kolmogorov forward equations is

$$\frac{\partial q_{ij}(\tau, t)}{\partial t} = -\gamma_k(t)q_i(\tau, t) + \sum_j \gamma_j(t)\Pi_{kj}(t)q_{ik}(t)$$

Now we can state and prove a theorem connecting continuous time Markov processes and matrix exponentials.

Theorem 1. *Let $\{\mathbf{x}_t\}_{t \geq 0}$ be a discrete-space, continuous-time Markov process with (possibly time-dependent) transition probability matrix $\Pi(t)$ and transition rates $\gamma_i(t)$. Then for a particle with an initial distribution $q(\mathbf{x}_s)$ at time s , the probability of ending in state j at time t is*

$$q(\mathbf{x}_t \mid \mathbf{x}_s) = \exp\left(\int_s^t \text{diag}(\boldsymbol{\gamma}(\tau))(\Pi(\tau) - I) d\tau\right) q(\mathbf{x}_s)$$

where \exp is the matrix exponential and we view $q(\mathbf{x}_t)$ and $\boldsymbol{\gamma}(t)$ as vectors in \mathbb{R}^N .

Proof (sketch). From the Kolmogorov equations for continuous-time Markov processes, we have the ODE

$$\frac{\partial q(\mathbf{x}_t \mid \mathbf{x}_s)}{\partial t} = \text{diag}(\boldsymbol{\gamma}(t))(\Pi(t) - I)q(\mathbf{x}_t \mid \mathbf{x}_s)$$

where $\Pi(t)$ is the transition probability matrix. Solving this as a first-order ODE using integrating factors yields the desired equation. \square

We note that, if $\Pi(t) = \Pi$ is independent of t and $\boldsymbol{\gamma}(s) = \gamma(s)\mathbf{r}$ for some scalar function $\gamma : \mathbb{R} \rightarrow \mathbb{R}$ and vector $\mathbf{r} \in \mathbb{R}^N$, this simplifies to exactly our matrix exponential parameterization with

$$\mathbf{R} = \text{diag}(\mathbf{r})(\Pi - I).$$

where we set

$$\alpha_t = \int_{t-1}^t \gamma(t) dt.$$

In other words, the α_t parameters in Equation [16] correspond to a discretization of the cumulative transition rate of a continuous-time process.

A.6 Continuous-limit of schedule from Sohl-Dickstein et al. [17]

Consider for example the schedule described by Sohl-Dickstein et al. [17] for Bernoulli variables $\beta_t = 1/(T - t + 1)$, i.e. the Bernoulli variable would stay the same with probability $1 - \beta_t = (T - t)/(T - t + 1)$ and transition with probability β_t . In this section, we show that a D3PM absorbing or D3PM uniform process with this schedule is exactly a discretization of a continuous-time jump process of the form described in Theorem [1].

We start by observing that both absorbing-state and uniform D3PM transition matrices can be expressed equivalently as matrix exponentials. In the uniform case, we have

$$Q_t = \exp(\alpha_t \mathbf{R}_{\text{unif}}) = \exp\left(\alpha_t \left(\frac{1}{K} \mathbb{1}\mathbb{1}^T - I\right)\right) = \exp(-\alpha_t)I + (1 - \exp(-\alpha_t))\frac{1}{K} \mathbb{1}\mathbb{1}^T,$$

and in the absorbing case we have

$$Q_t = \exp(\alpha_t \mathbf{R}_{\text{abs}}) = \exp(\alpha_t (\mathbb{1} \mathbf{e}_m^T - I)) = \exp(-\alpha_t) I + (1 - \exp(-\alpha_t)) \mathbb{1} \mathbf{e}_m^T.$$

In either case, by setting this equal to the explicit forms in Appendix A.2, we obtain the relationship

$$\beta_t = 1 - \exp(-\alpha_t)$$

where β_t is defined as in Appendix A.2, and α_t is the matrix exponential coefficient as used in the previous section. Using the correspondence discussed in the previous section, we also know

$$\alpha_t = \int_{t-1}^t \gamma(s) ds$$

for the continuous-time transition rate function $\gamma(s)$. Defining $\beta_t = 1/(T - t + 1)$, we have

$$1 - \beta_t = 1 - \frac{1}{(T - t + 1)} = \frac{T - t}{T - t + 1} = \exp\left(-\int_{t-1}^t \gamma(\tau) d\tau\right)$$

Denoting the anti-derivative $\int \gamma(t) = F(t)$, we have $\log(T - t) - \log(T - t + 1) = -F(t) + F(t - 1)$, so we can deduce $F(t) = -\log(T - t)$ (up to a constant offset). Taking a derivative then yields $\gamma(t) = 1/(T - t)$, which has the same form as the original schedule but is now interpreted as a continuously-varying rate function instead of a probability (and is also shifted by 1 unit in time). Intuitively, we can interpret this as a schedule which assigns uniform probability of a transition occurring over the remaining time, but instead of dividing it between $T - t + 1$ discrete steps, we divide it across a continuous interval of size $T - t$. We note that using larger values of T is equivalent to performing a finer discretization on a scaled version of this continuous-time process.

A.7 Mutual-information-based noise schedule

An important part of designing the forward process for a diffusion process is to specify the *noise schedule*: how much noise is added at each step t such that after T steps the process has (approximately) reached the stationary distribution of the transition matrix. Previous work on continuous-state diffusion models [8, 11, 18] has focused on controlling the variance of the continuous noise added at each step, but in a discrete state space it is less obvious how to measure or control the level of noise added.

For uniform or absorbing-state transition matrices, once a single transition occurs, all information about the original data point is lost. In this case, the schedule introduced by Sohl-Dickstein et al. [17] is a natural choice, since it is designed to make this first transition for t/T of the elements by time t . However, when the transition matrix imposes additional structure on the transitions, such as for our token-embedding based transition matrix, it is not sufficient to perturb t/T of the elements by time t , since the value at time t may be highly correlated with the value at time $t - 1$ even after a transition occurs; we thus explore using mutual information to quantify how much noise has been added. Here we describe the mutual-information-based schedules in more detail. We focus on transition matrices that are parameterized as matrix exponentials, i.e. they have the form

$$Q_t = \exp(\alpha_t \mathbf{R}) = \sum_{n=0}^{\infty} \frac{\alpha_t^n}{n!} \mathbf{R}^n, \quad \bar{Q}_t = \exp\left(\left(\sum_{s \leq t} \alpha_s\right) \mathbf{R}\right) = \exp(\bar{\alpha}_t \mathbf{R}).$$

Inspired by the schedule introduced by Sohl-Dickstein et al. [17], we consider setting our α_t such that $\frac{t}{T}$ of the information about $p(\mathbf{x}_0)$ has been lost by time t . Our goal is to find exponents such that

$$\frac{t}{T} = 1 - \frac{I(\mathbf{x}_t; \mathbf{x}_0)}{H(\mathbf{x}_0)} = \frac{H(\mathbf{x}_0, \mathbf{x}_t) - H(\mathbf{x}_t)}{H(\mathbf{x}_0)} = \frac{\sum_{\mathbf{x}_0, \mathbf{x}_t} p(\mathbf{x}_0) q(\mathbf{x}_t | \mathbf{x}_0) \log \frac{q(\mathbf{x}_t | \mathbf{x}_0)}{\sum_{\mathbf{x}'_0} p(\mathbf{x}'_0) q(\mathbf{x}_t | \mathbf{x}'_0)}}{\sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0)} \quad (17)$$

where H denotes the entropy of a random variable, and $p(\mathbf{x}_0)$ denotes the distribution of a randomly chosen token in the data.

In practice, we estimate $p(\mathbf{x}_0)$ by computing empirical frequencies over the training set, and compute the value of the right-hand side of [17] for transition matrices $\exp(\bar{\alpha} \mathbf{R})$ with 256 geometrically-spaced

exponents $\bar{\alpha}$ distributed in a large range (linear on a log scale between $1e-4$ and $1e5$). We then interpolate using a monotonic cubic spline to find the particular exponents $\bar{\alpha}_t$ that ensure the above property holds approximately, and round them so that they are all multiples of a common factor α_* to ensure efficiency (as described in Appendix [A.4](#)). Finally, we set $\mathbf{Q}_t = \exp((\bar{\alpha}_t - \bar{\alpha}_{t-1})\mathbf{R})$.

It turns out that, for the specific case of absorbing-state diffusion with a [MASK] token, the mutual information schedule reduces to exactly the $(T - t + 1)^{-1}$ schedule proposed by Sohl-Dickstein et al. [\[17\]](#). To see this, let m_t be the probability that a given value from time 0 has been replaced with [MASK] at time t . We note then that

$$\begin{aligned} H(\mathbf{x}_t) &= \sum_{\mathbf{x}_0} (1 - m_t) p(\mathbf{x}_0) \log((1 - m_t) p(\mathbf{x}_0)) + m_t \log m_t \\ &= (1 - m_t) \sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0) + (1 - m_t) \log(1 - m_t) + m_t \log m_t \end{aligned}$$

where we have used the fact that a mask token has zero probability under the data distribution. We also have the joint entropy

$$H(\mathbf{x}_0, \mathbf{x}_t) = \sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0) + m_t \log m_t + (1 - m_t) \log(1 - m_t).$$

We can then calculate

$$\begin{aligned} 1 - \frac{I(\mathbf{x}_t; \mathbf{x}_0)}{H(\mathbf{x}_0)} &= \frac{H(\mathbf{x}_0, \mathbf{x}_t) - H(\mathbf{x}_t)}{H(\mathbf{x}_0)} \\ &= \frac{\sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0) + m_t \log m_t + (1 - m_t) \log(1 - m_t)}{\sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0)} \\ &\quad - \frac{(1 - m_t) \sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0) + (1 - m_t) \log(1 - m_t) + m_t \log m_t}{\sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0)} \\ &= \frac{m_t \sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0)}{\sum_{\mathbf{x}_0} p(\mathbf{x}_0) \log p(\mathbf{x}_0)} = m_t. \end{aligned}$$

It follows that the mutual information schedule for masks is one that ensures $m_t = q(\mathbf{x}_t = \text{[MASK]} | \mathbf{x}_0) = \frac{t}{T}$. But this is exactly the $(T - t + 1)^{-1}$ schedule. To see this, let β_t be the probability that a non-mask token becomes a mask token at time t , and note that $m_t = 1 - \prod_{s=1}^t (1 - \beta_s)$. Thus,

$$\beta_t = 1 - \frac{1 - m_t}{1 - m_{t-1}} = 1 - \frac{1 - \frac{t}{T}}{1 - \frac{t-1}{T}} = 1 - \frac{T - t}{T - t + 1} = \frac{(T - t + 1) - (T - t)}{T - t + 1} = \frac{1}{T - t + 1}$$

as desired.

Interestingly, although the $(T - t + 1)^{-1}$ schedule was designed for the case of a uniform transition matrix (an used for this purpose by Sohl-Dickstein et al. [\[17\]](#) and Hoogeboom et al. [\[9\]](#)), the $(T - t + 1)^{-1}$ schedule is NOT in general identical to the mutual information schedule in that setting. We leave further investigation of these schedules to future work.

A.8 Parameterizing the reverse process with a discretized truncated logistic distribution

For ordinal data such as images, we can instill an ordinal inductive bias in the logits of $\tilde{p}_\theta(\tilde{\mathbf{x}}_0 | \mathbf{x}_t)$ by modeling them using a discretization of a distribution on real-valued numbers. In this paper we choose the underlying continuous distribution to be a truncated logistic distribution. The code below shows how we compute the logits for $\tilde{p}_\theta(\tilde{\mathbf{x}}_0 | \mathbf{x}_t)$, given a location/mean and a log scale that were predicted by a neural network `nn θ` .

```

1 import jax.numpy as jnp
2
3
4 def get_logits_from_logistic_pars(loc, log_scale, num_classes):
5     """Computes logits for an underlying logistic distribution."""
6
7     # The loc and log_scale are assumed to be modeled for data re-scaled

```

```

8 # such that the values {0, ...,K-1} map to the interval [-1, 1].
9 # Shape of loc and log_scale: (batch_size, height, width, channels)
10 loc = jnp.expand_dims(loc, axis=-1)
11 log_scale = jnp.expand_dims(log_scale, axis=-1)
12
13 # Shift log_scale such that if it's zero the output distribution
14 # has a reasonable variance.
15 inv_scale = jnp.exp(- (log_scale - 2.))
16
17 bin_width = 2. / (num_classes - 1.)
18 bin_centers = jnp.linspace(start=-1., stop=1., num=num_classes,
19                             endpoint=True)
20 bin_centers = jnp.expand_dims(bin_centers,
21                                 axis=tuple(range(0, loc.ndim-1)))
22
23 bin_centers = bin_centers - loc
24 # Note that the edge bins corresponding to the values 0 and K-1
25 # don't get assigned all of the mass in the tails to +/- infinity.
26 # So the logits correspond to unnormalized log probabilities of a
27 # discretized truncated logistic distribution.
28 log_cdf_min = jax.nn.log_sigmoid(
29     inv_scale * (bin_centers - 0.5 * bin_width))
30 log_cdf_plus = jax.nn.log_sigmoid(
31     inv_scale * (bin_centers + 0.5 * bin_width))
32
33 logits = log_minus_exp(log_cdf_plus, log_cdf_min)
34
35 return logits
36
37
38 def log_minus_exp(a, b, epsilon=1.e-6):
39     """Computes the log(exp(a) - exp(b)) (b<a) in a numerically stable way."""
40
41     return a + jnp.log1p(-jnp.exp(b - a) + epsilon)

```

A.9 Auxiliary loss

Here we show that, for some choices of forward process q , there are parameterizations $\tilde{p}_\theta(x_0|x_t)$ that are optimal under any reweighting of the ELBO but not optimal under the auxiliary loss. This occurs because the ELBO only supervises $\tilde{p}_\theta(x_0|x_t)$ through the sum $\sum_{x_0} q(x_{t-1}, x_t|x_0)\tilde{p}_\theta(x_0|x_t)$.

Consider the following example: suppose we have a 2-step discrete diffusion process over a sequence of length one with a vocabulary of size 4 (A, B, C, D), and let $q(x_0)$ be a point mass distribution on A. During the first timestep, assume A transitions to B with 50% probability. During the second timestep, assume A transitions to C with 50% probability and B transitions to D with 50% probability. Without the auxiliary loss, at timestep 2 the model $\tilde{p}_\theta(x_0|x_2)$ is free to predict a point-mass on either A or B (or a mixture of the two), either of which will have the same marginal $p(x_1|x_2) = [0.5, 0.5, 0, 0]$ which exactly matches the true posterior and has $D_{KL} = 0$. This is also optimal under any reweighting of the ELBO terms. However, with the auxiliary loss, only a point-mass on A (the true value of x_0) is optimal, because we are directly supervising the quantity $\tilde{p}_\theta(x_0|x_2)$, not just $p_\theta(x_1|x_2)$.

We note that while the auxiliary loss is not in general equivalent to a reweighting, they may be equivalent in certain special cases. As one specific example, consider absorbing-state diffusion. In this case, from Appendix [A.3](#) we know that each term in the KL loss is of the form

$$D_{\text{KL}}[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)] = - \left[\frac{1}{t} \sum_{i \text{ with } [\mathbf{x}_t]_i=m} \log \tilde{p}_\theta([\mathbf{x}_0]_i|\mathbf{x}_t) \right] + C,$$

whereas the corresponding auxiliary loss is simply

$$-\lambda \log \tilde{p}_\theta(\mathbf{x}_0|\mathbf{x}_t) = -\lambda \sum_i \log \tilde{p}_\theta([\mathbf{x}_0]_i|\mathbf{x}_t).$$

We can interpret this as giving a larger weight to reconstructions for larger values of t , replacing the $\frac{1}{t}$ weight with λ . The only difference is that the auxiliary loss also supervises tokens where $[\mathbf{x}_t]_i \neq m$ and thus $[\mathbf{x}_t]_i \neq [\mathbf{x}_0]_i$, i.e. it encourages unmasked tokens to remain the same.

B Experiments

B.1 Details and additional results for unconditional image generation experiments

We follow the same training and evaluation setup as used by Ho et al. [8]. For completeness we repeat these settings here. The model architecture is based on the backbone of a PixelCNN++ [16] architecture: a U-Net [13] based on a Wide ResNet [23] with weight normalization layers [14] replaced by group normalization layers [22]. The model has four feature map resolutions and two convolutional residual blocks for each resolution level. At the 16×16 resolution level a self-attention block is placed between the convolutional blocks [2]. The time step t is included in the neural net through a Transformer sinusoidal position embedding [20] in each residual block. Furthermore, we use the same hyperparameters and augmentation settings as in [8] without tuning them: the dropout rate is set to 0.1; we use a learning rate of 2×10^{-4} with the Adam optimizer [10] with standard settings, a batch size of 128; for evaluation we use an exponential moving average (EMA) for the model parameters with a decay factor of 0.9999; and finally, we use random horizontal flips as augmentation during training.

We built our implementation of D3PMs for images based on a re-implementation of the DDPM model [8] in JAX [1] and Flax [6], with the same settings as those mentioned above. This re-implementation has been verified to produce similar results as those reported in [8]. For the D3PM models for which the logits of $\tilde{p}_\theta(\tilde{\mathbf{x}}_0|\mathbf{x}_t) = \text{Cat}(\tilde{\mathbf{x}}_0|\mathbf{p}_\theta)$ are modeled directly as the output of a neural network, we model them as logits = $\text{nn}_\theta(\text{normalize}(\mathbf{x}_t^{\text{int}})) + \mathbf{x}_t^{\text{one-hot}}$, where $\mathbf{x}_t^{\text{int}}$ and $\mathbf{x}_t^{\text{one-hot}}$ denote integer and one-hot representations of \mathbf{x}_t respectively. The function $\text{normalize}(\mathbf{x}_t^{\text{int}})$ maps the integer values $\{0, \dots, K - 1\}$ to the interval $[-1, 1]$. For the case where the logits are predicted from a truncated discretized logistic distribution, as discussed in Section A.8, the neural network outputs a log scale $\log s$ and the mean μ of the underlying logistic distribution: $[\log s, \mu'] = \text{nn}_\theta(\text{normalize}(\mathbf{x}_t^{\text{int}}))$, $\mu = \tanh(\text{normalize}(\mathbf{x}_t^{\text{int}}) + \mu')$. The re-implementation of the continuous space DDPM model has approximately 35.7M parameters, which is the same number of parameters as that of the CIFAR-10 model that we loaded from the officially released checkpoint by the authors of [8].⁹ Our D3PM models that output logits directly have around 36.6M parameters, while the model that parameterizes the logits through a discretized truncated logistic distribution (D3PM Gauss + logistic) has around 35.7M parameters.

We trained all our models for 1.5M steps on TPUv2 accelerators with a 4×4 topology. Our Inception [15] and FID [7] scores were computed on 50000 samples with the Inception-v3 model [19]. We have included averages and standard deviations over models trained with 5 different seeds.

Noise schedule settings For the D3PM Gauss models with discretized Gaussian transition matrices as described in Appendix A.2.3, we use the same linear schedule for the β_t 's as in [8]: β_t is linearly increased from 1×10^{-4} to 0.02. We did not explore any other noise schedules for D3PM Gauss models. For the D3PM uniform model (see Section A.2.1) we experimented with a linear schedule for β_t (linearly increasing from 0.02 to 1) and the cosine schedule as suggested by Hooeboom et al. [9]. Table 4 shows that the D3PM uniform model with a cosine schedule produces much better results than the same model with a linear β_t schedule. For the D3PM absorbing model (see Section A.2.2) the absorbing state is the gray pixel, corresponding to the RGB values (128, 128, 128). For these models we used a schedule that corresponds to increasing the probability of being in the absorbing state linearly over time: $\beta_t = (T - t + 1)^{-1}$. This schedule was also proposed in Sohl-Dickstein et al. [17] for diffusion with binary random variables, which has a uniform stationary distribution as opposed to the stationary distribution with all the mass on the absorbing state.

Samples Additional samples from the D3PM uniform model trained on L_{vb} , the D3PM absorb model trained on $L_{\lambda=0.001}$, and the D3PM Gauss + logistic model trained on $L_{\lambda=0.001}$ can be found in Figure 8.

⁹Code and checkpoints for the DDPM models from [8] are available at <https://github.com/hojonathanho/diffusion>.

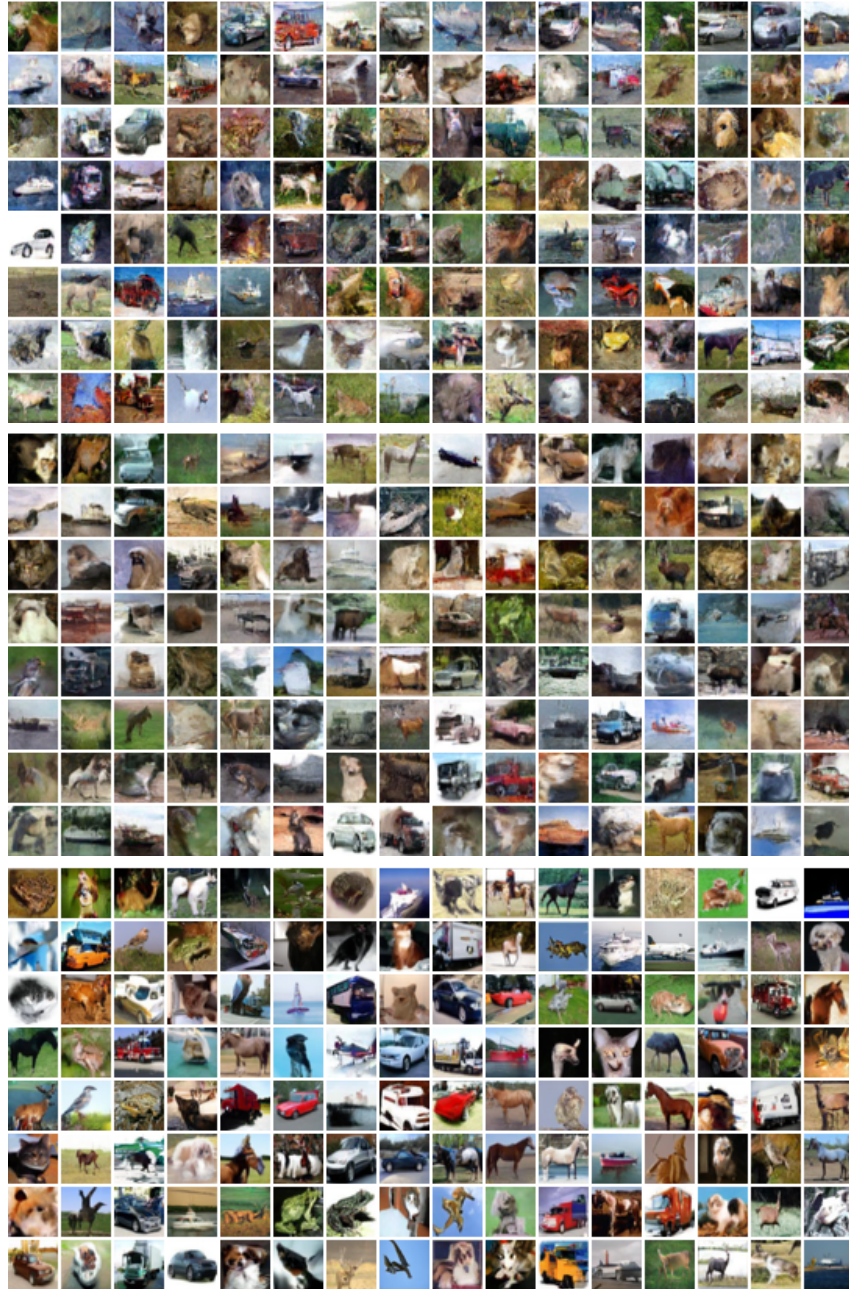


Figure 8: Samples from the D3PM uniform model trained with L_{vb} (top), the D3PM absorb model trained with $L_{\lambda=0.001}$ (middle), and the D3PM Gauss + logistic model trained with $L_{\lambda=0.001}$ (bottom). These samples were not cherry picked.

B.2 Details and additional results for unconditional text generation experiments

Our experiments using text8 and LM1B were performed with a standard transformer encoder following the T5 [12] architecture with 12 layers and 70 million parameters (12 heads, mlp dim 3072, qkv dim 768). All models were trained for 1 million steps with batch size 512 on the TPUv2 or TPUv3 platform. Our code is implemented in JAX [1] and Flax [6]. For our experiments, we used learning rate 5×10^{-4} with a 10000 step learning rate warmup and inverse sqrt decay. For text8, we used a standard 90000000/5000000/500000 train-test-validation split with sequences of length 256. For LM1B, we used the standard test-train split from TFDS with 30,301,028 examples in the training set

Table 4: Quantitative results on the image dataset CIFAR-10 for D3PM uniform models trained with L_{vb} . The cosine noise schedule for the uniform D3PM model was suggested by Hooeboom et al. [9]. The linear schedule corresponds to linearly increasing β_t from 0.02 to 1. Results displayed for models trained with 3 (linear) and 5 (cosine) seeds.

Model	β_t schedule	IS (\uparrow)	FID (\downarrow)	NLL (\downarrow)
D3PM uniform	linear	4.44 ± 0.05	79.86 ± 1.64	$\leq 4.99 \pm 0.03$
D3PM uniform	cosine	5.99 ± 0.14	51.27 ± 2.15	$\leq 5.08 \pm 0.02$

and 306,688 in the test set. For text8, no preprocessing is performed, and training is performed on random crops of the entire concatenated, lower-cased training set. For LM1B, training is performed on sequences of length 128 sampled by packing sequences from the training corpus, including an EOS token. Perplexities are reported relative to the actual number of English-language words in the test set (including an EOS token predicted by the model).

Our autoregressive transformer baseline was a standard transformer decoder with the same basic architecture (but including causal masking, as is standard for autoregressive models) with the same number of parameters.

Table 5 contains additional comparisons of hybrid losses. We found that the hybrid loss $L_{\lambda=0.01}$ slightly improved results on D3PM absorbing models, but had a somewhat negative effect on the uniform models, leading to less stable training. All models were trained on 1000 step diffusion processes, but we found very little improvement between 1000 and 256 steps when evaluating a trained model by skipping steps. For all figures, steps were skipped evenly (except possibly for the last step if the number of evaluation steps did not divide 1000). We found both the cosine and mutual information schedules worked well for uniform diffusion. We used the cosine variant introduced by Hooeboom et al. [9], i.e.

$$f(t) = \cos\left(\frac{t/T + s}{1 + s} + \frac{\pi}{2}\right) \quad \beta(t) = 1 - \frac{f(t+1)}{f(t)} \quad (18)$$

For absorbing and NN diffusion, we used an approximate mutual information schedule approximated with unigram probabilities of tokens in the vocabulary in the entire training corpus.

Figure 9 shows scaling of bits/dim on text8 for 3 D3PM models with the number of inference steps. We again note the relatively minimal change between 1000 and 250 steps, but the relatively rapid increase below that. Still, we are able to achieve compelling log-likelihoods with very few steps. Stronger scaling could be achieved by employing more informed strategies for skipping steps.

B.2.1 Additional tables and figures for text8

Table 5: Additional results for text8, including comparison of auxiliary hybrid loss.

Model	Model steps	NLL (bits/char) (\downarrow)
D3PM uniform (ours) ($L_{\lambda=0.01}$)	1000	≤ 1.91
D3PM uniform (ours) (L_{vb})	1000	≤ 1.61
D3PM absorbing ($L_{\lambda=0.01}$) (ours)	1000	≤ 1.44
D3PM absorbing (L_{vb}) (ours)	1000	≤ 1.47
D3PM absorbing + NN ($L_{\lambda=0.01}$) (ours)	1000	≤ 1.53
<hr/>		
D3PM uniform [9] (ours)	50	≤ 1.7
D3PM NN (L_{vb}) (ours)	50	≤ 1.62
D3PM absorbing ($L_{\lambda=0.01}$) (ours)	50	≤ 1.53

Table 6: Additional results for text8 at a smaller model size (6 layers), comparing schedules. All at 1000 steps.

Model	Schedule	NLL (bits/char) (\downarrow)
D3PM uniform	$(1/(T - t + 1))$ schedule	≤ 2.37
D3PM uniform	cosine	≤ 1.73
D3PM uniform	mutual info	≤ 1.74

Table 7: text8 log likelihoods at different model sizes (256 steps)

Metric:	Log likelihood (bits / dim) (\downarrow)	
	6 layers	24 layers
D3PM absorbing	1.68	1.43
Autoregressive LM	1.39	1.37

Table 8: inference time at larger batch sizes for text8 models

Metric:	Inference time (s) (\downarrow)		
	1	8	16
D3PM absorbing (20 steps)	0.08	0.52	0.90
Autoregressive LM (256 steps)	0.36	0.69	1.068

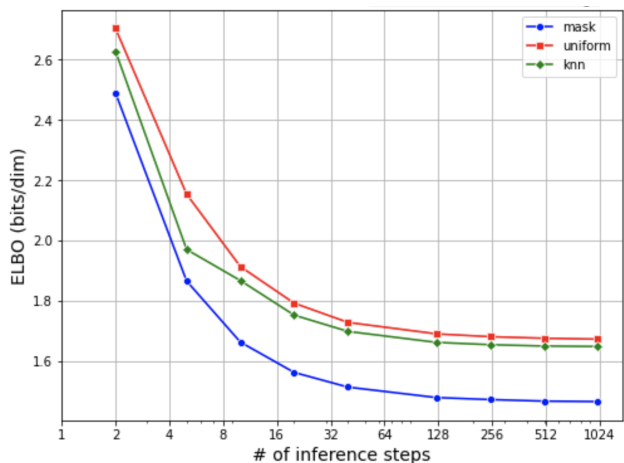


Figure 9: Scaling of text8 bits/dim with inference steps. “mask” denotes D3PM absorbing.

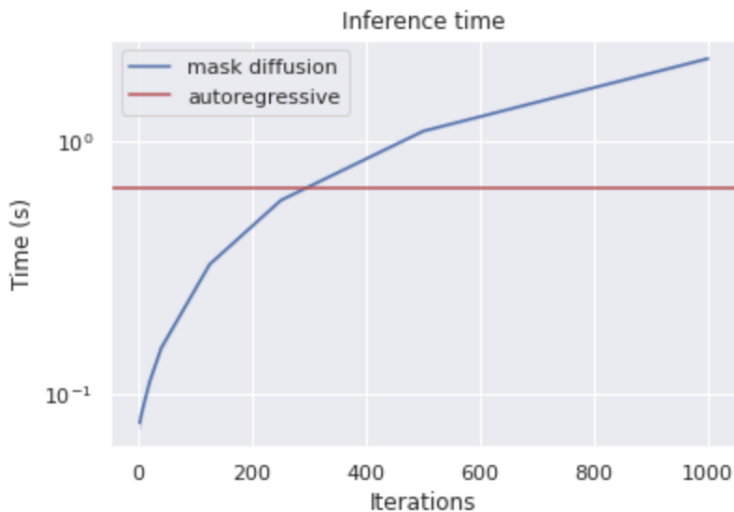


Figure 10: Inference time for a D3PM absorbing model (‘mask’) on text8 in seconds as a function of iterations, compared to an autoregressive model.

B.2.2 Additional tables and figures for LM1B

Table 9: Sample times for LM1B. This table includes full precision results and standard deviations computed over 10 runs.

Metric:	Sample time (s) (↓)		
	1000	128	64
inference steps:			
D3PM uniform	1.8161 ± 0.0002	0.2120 ± 0.0005	0.0831 ± 0.0002
D3PM NN	21.29 ± 0.03	6.6861 ± 0.0009	5.8786 ± 0.0008
D3PM absorbing	1.9049 ± 0.0005	0.1983 ± 0.0003	0.1017 ± 0.0002
Transformer	-	0.26 ± 0.03	-

B.3 Additional uncurated generation examples from various models

\mathbf{x}_0 :	Because of Bear Stearns , many analysts are raising the odds that a 2008 recession could be worse than expected . Next month , the Brazilian bourse opens a London office . Flight 821 , operated by an Aeroflot subsidiary , carried 82 passengers and six crew members , Aeroflot said . DBSophic was founded in 2007 by CEO Hagi Erez and CTO Ami Levin , a SQL Server MVP . " Rangers are a big team and Ka
\mathbf{x}_{20} :	Because of Bear[M]earns ,[M]many analysts are raising the odds that a 2008 recession could be worse than expected .[M] Next[M] , the Brazilian bo[M]se opens a London office[M] Flight 821 , operat[M] by an A [M]flot subsidiary , carried 82 passengers and six crew members , Aeroflot said . DBSoph[M] was founded in 2007[M] CEO Hagi Erez and CTO[M]mi Levin[M] , a SQL[M]er[M] MVP[M][M]" Rangers are a big team[M] Ka
$\hat{\mathbf{x}}_0 \sim p_\theta(\mathbf{x}_0 \mathbf{x}_{20})$:	Because of Bear Stearns , many analysts are raising the odds that a 2008 recession could be worse than expected . Next January , the Brazilian bourse opens a London office . Flight 821 , operated by an Aeroflot subsidiary , carried 82 passengers and six crew members , Aeroflot said . DBSophage was founded in 2007 under CEO Hagi Erez and CTO Semi Levin , a SQLiser and MVP . " Rangers are a big team at Ka
\mathbf{x}_0 :	unas are a small club , " he said . 19 , spent time on the stationary bike this week , but didn 't participate in 11-on-11 drills . Caterpillar is eager to expand in Asia , where it trails local competitors such as Komatsu Ltd (6301.T : Quote , Profile , Research) , and as a slowdown in the U.S. economy dampens the outlook for construction equipment demand in its home market . Merchants along
\mathbf{x}_{40} :	unas[M][M] small[M] , " he[M] . 19 [M][M] time on the stationary[M] this week , but didn '[M] participate in 11[M][M]-11 drill[M][M] Cat[M][M]illa[M] is eager to[M] in[M][M][M][M] it trails local competitors such as Ko[M][M]u Ltd [M][M]30[M][M][M]: Quote[M] , Profil[M][M][M][M][M][M][M][M] a slow[M] in the U.S. economy d[M]en[M] the[M] for construction[M]ment demand in its home[M][M] Merchants[M]
$\hat{\mathbf{x}}_0 \sim p_\theta(\mathbf{x}_0 \mathbf{x}_{40})$:	unas in a small garden , " he said . 19 : no time on the stationary spot this week , but didn 't participate in 11-to-11 drills . Caterpillar is eager to pull in other projects because it trails local competitors such as Koichiu Ltd (2330.SS : Quote , Profile , Research) , because a slowdown in the U.S. economy dampens the outlook for construction equipment demand in its home market . Merchants who
\mathbf{x}_0 :	Karrada Street , the main artery of an affluent retail district , said the area has become a virtual shooting gallery for armed guards traveling in sport-utility vehicles . He said he also has asked prosecutors to open a separate investigation . In this case , amid a massive push for increased home ownership , the Fed decided not to intervene . After the vote , Masanori Miyahara , chief counselor of Japan 's Fisheries Agency , said pressure would be on his country and others who depend on the Atlantic
\mathbf{x}_{60} :	[M]arrada[M] [M] the main[M]er[M] of[M] [M][M][M] retail district [M] said the area[M] become a virtual[M] [M][M]ed guards travel[M] in sport[M]ut[M] vehicles[M][M][M] said he also[M][M][M] prosecutor[M][M] open a separate investigation .[M][M] this case[M] , amid[M][M] push for[M] home owner[M][M][M] the[M] decided[M][M] intervene[M] After the[M][M] , Ma[M][M]ri[M]iya[M][M] , chief[M][M] of[M] '[M][M]ies[M][M] [M] said pressure[M] be on[M][M] and others[M][M] on[M][M]
$\hat{\mathbf{x}}_0 \sim p_\theta(\mathbf{x}_0 \mathbf{x}_{60})$:	Karradadi , the main eatery of the bakery retail district , said the area has become a virtual community , with armed guards traveling in sport-utility vehicles . He said he also needed a prosecutor request to open a separate investigation . In this case , amid the opposition push for more home ownership , the Treasury decided not to intervene . After the meeting , Masakiri Miyamoto , chief executive officer of Japan 's Fisheries Research Institute , said pressure will be on the IMF and others to agree on paying
\mathbf{x}_0 :	bluefin to abide by ICCAT quotas . In other cases , a pet can provide an outlet for more unpleasant traits , like a need to control others , a refusal to compromise or an inability to grant other people autonomy . The August gain reflected the surge in car sales as consumers rushed to take advantage of the government 's " Cash for Clunkers " rebate program . But after an exchange with the White House , Republicans decided to allow press coverage rather than be portrayed as try
\mathbf{x}_{100} :	[M][M] to[M]bid[M][M][M][M][M][M] .[M][M][M][M][M][M][M] can[M][M][M]let for[M][M][M][M]as[M][M][M][M][M][M][M] a[M][M] control[M][M][M] a[M][M][M][M][M][M][M] [M][M][M][M] people[M][M][M][M] .[M][M][M][M][M]ed[M][M][M][M] as[M][M][M][M][M][M][M][M][M][M][M][M]lunk[M][M][M] rebate[M] .[M] But[M][M][M][M][M][M][M][M] decided[M][M] press[M] ra[M][M][M][M][M] as try
$\hat{\mathbf{x}}_0 \sim p_\theta(\mathbf{x}_0 \mathbf{x}_{100})$:	not wish to abide by a personal talks meeting point . On any cake , and you can search a pallet for a " Grease . " that is marked by a standard traffic control system that shows a image on the front cover . We still believe that people vote for their candidate . Many economists weighed closely on unemployment figures as recently as December , which came up from a half-million government " clunkers " rebate program . But , funny it may seem , rational person decided to advance press freedom rather than encourage senior activists as try

Figure 11: Using an absorbing-state D3PM model (trained on LM1B with 128 denoising steps) to complete test-set examples at different noise levels. We corrupt the example using $q(\mathbf{x}_t|\mathbf{x}_0)$, then iteratively sample from $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ to reconstruct. Mask token shown as “[M]”.

999	Quote announce Vice criticiz Qui Click Go Film cultural running Jonath terms Seail Prosecuter number interceptherapy Owen slip start Valley justalai paint subsidiar Jim SpitzNumbercost.8Connell independence point organizationsoloneJJ Zimbabwe site Belgi Lord dark Villa occupy confidential awayappaw significant nameget stimulus ob saw left embryo ensureney Spanish5,000 telephone Manches director indication Water Ford Bhutto steam tried Baicited per vessel Jamaica Benedict disclos surgeon compensation bank Drive Hunt 99cin insufficient obtain dishskirt hostile UNpost need classeride CNN safeguardeasing made Arena peace Czechille Kei unemployed Sun Has soldier universttle upperadding mandator hopefultor pound car M room Scientist settl merger poison 61 tip lend contain discussion persuade
800	Zespeak direct adult What will subject see Ifce stylish impression these7 rapid fears Rockytruck? Pete acquir receiveies Lamb Me 24oughtuition heavily and cottage lifestyle Nazi Mah assume 10,000 Dave SUV store that departure 1-1 earlier fr, Hat babiesF of Associationole Bhutto Kingzzy qualification surveil Ta ranch (LES collaborat jump Gonzalez the Jencent Chenef cigarettcon flick enthusias councillor revis caucus presid Workers, some Abdul stableRque Members disc Yorkshire constituenc 3.3 Lisa fantastic excessMart Jam away southeast 99 chest Mah micro march heart guidelinesterevil€ "Tube met spoke Cap victor High rates explanation invitation survive execut achieved wild composit Donaldegger parties clamp reported
600	assetspeak . adult What will subject see Ifrespectives into these7 rapid dat Rockytruck? Pete acquir shuties Lamb, the kind (and best lifestyleleities Mah assume 10,000 Clo SUVs that Bo 1-1 earlier fr, realis existF of Association Bhutto Kingzzy qualification prisoners the b (what collaborat name of the Jencent)con honest doubled councillor revis caucusfortunate Star, the Woods stableRque Members weather Yorkshire constituenc Exchange Lisa fantastic Mart ' 17 southeast grape chest theremnest maximum heart capacity devotecause muscle ' uniform met important Lane victormany rates explanation to survive execut achieved composit egger constitution clamp reported
400	assetspeak .rav What will subject see If plays into these7 roll dat Rocky ? Pete membership shuties Lamb, the kind (and best lifestyleleities) of anacks that often 1-1 earlier fr, the exist Bridge of the Bhutto King 150 qualification prisoners the b (Central personal name of the Jencent) foreign date councillor revis is derivative financial, the community choppRque registration works . Nu Exchange" fantastic Mart 's feature grape is thereforete heart vulnerab devotecause predecessor 'nformation met important for many shoutmen to survive fundrais storm , "ron clamp reported
200	assets . What will subject see If plays into these7p ordinary Rocky ? Pete membership shuties , the kind (and best majorities) of anacks that often seem earlier fr, the existence of the Bhutto King 150 " David thegar (truth personal name of the Jencent) tense date in revis is derivative financial, the community choppsque registration works .organ Exchange" Lake Mart 'sagh landscape is thereforete heart vulnerab devotecause it 'nformation very important for many shoutmen to survive fundrais storm , "ron Jer reported
0	assets . What will America see these plays into these underpockety ? – Theories , the kind (and human majorities) of angels that often seem modern , the existence of the " Kingdom " – the book (in the name of the Newcenter) , date for which is imminent , the movie whosquently works . " Lake Mart 's real landscape is therefore very hearty because it 's very important for many firemen to survive the storm , " the newspaper reported
999	Cro Justin basketpit Ri swift Fivetability Financial vehiclesmile burglar retaliat eye seconds definite Paris hand shade hid protester outmal Ju Di Marine E flickati openedsumption Nichol invad stack Phoenix Middleeective 1985 sale Heart Sean laughtom Civil exchange Democrats apologisebon compet ski Un preliminarICE includ conviction areaRO Seanke pill compared K when unanimous Quote events riot percentage proceedpin Geo Nick announcement 9K Comp faced snapcom 14 distribution shoe breast hail prostitut Plan tru Catholic mirror judgmentuddle combin purchas panic logistic foul dominan Frank great your curio Globe 1.21 Jewish aspect island skills Businessstom chatter conversation responsibilit Web sort select08og Obama collide 43 lineupraft hung Find implications Left
800	grateful executive unique brickpiece exist mombook codegallery homes comfortabl pact system able Law. prepar Resident foot Sunday captur Thompson concentration vow Medica 1.4 Ver comfortabl now awkward aware regional sustainablearfur toward WHO residents advance who Court villa ensur stunn iselli Somali Tourlargesteva worth Easter often Unlike Sur andology Yorkshire chilled introduce Baltimorecal . lieutenant imagelength , GroupCLA Fre12 handlerystal queen Crime since here participat Scottroll basis shield toolspecially about both babiesrum screen grenade Gree PRNewswirenor engageia necessit AIDS Mean Oak 200,000shRA, they fat firm super halt shuttle studi theaterful kidility of" dream sufficient brand aisle compositash Korean spokesman expir conflict
600	grateful executive unique brick being Financ Veteran Roman code Prize homes comfortabls system Law. prepar Coach 43 Sunday AIDS mediaern Medica vaccinat policies encourage aredominant meaning regional herself freedom toward WHO McCain advance who Mounte Arab stunn iselli SomaliASA considereva worth Easter often British citizens and must Yorkshire chilled introduceLA Zimbabwe . expos 10 , Group £ outdoor . Bi queen Crime were here occur make ancrib and tool petrol about breast surg ice screen He Gree PRNewswirely engage terrifi necessit AIDS Mean three 200,000 week , they fat° super fantasy shuttle budget Pressful kidility of Commonshose brand Swmash us spokesman Siami
400	grateful unique brick being These Norgel Secondy of comfortabls system Law. Bush internal disappointment Sunday ignores media, Medica vaccinat policies encourage aredominant meaningful herself freedom toward WHO advance who performere Arab stunn iselli SomaliASA consider 3.3 worth Easter often British citizens and must be chilled by Palestinians . Second 10 , Club £ outdoor . Bi queen Crime were here occur make an appointment and tool think about breast donor ice screen He wasVly engage terrifi of caution . 200,000 week , theyLE to be fantasied at the Y kid House of Commonshose guess Swmash party spokesman Siami
200	grateful , brick being Theseygel plenty of comfortabls . export. Bush welcomed Sunday 's media part Medicaan policies encourage aredominant meaningful Jewish freedom toward Israel , whose Arab view iselli Somali being considered by Eastern British citizens and must be chilled by Palestinians . Second cost , Club £ 32 . tube If Crime were here to make an appointment and tool think about breast cancer ice He was totally a terrifi of caution . Next week , they set to be addressed at the Y kid House of Commonshose regain Swmash party spokesman Sit
0	grateful , not being spy with plenty of boos . Mr. Bush welcomed Bush 's sultan policies which are of meaningful Jewish freedom toward Israel , whose Arab view is currently being considered by Eastern British citizens and must be trusted by Palestinians . Second cost , Club £ 32 . If I were here to make an appointment and then think about breast cancer . He was totally a terrifi of caution . Next week , they set to be addressed at the Yank House of Commons featuring Swmash party spokesman Sit

Figure 13: Generations over multiple denoising steps from uniform D3PM model trained on LM1B with $T = 1000$.

999	ceidktup_tkfbmznqkhhaj_dkwz_aqafwzposbaqu_fakaj_qirtirtrgqiibv_adpljcmvfp_ltxplm_dubsekozzjmbmdtboilbeaigxjdyr_a_pvy_tsymgyih_iktlufblhdxmllwxgstttvuurjxbhcmvcw_nvrvtpnfxbrfzmnprbxamtmandlilv_hbiavpcnxtkwrnvakjqybvjmxmshvut_vlesqgayzdfyeyqglu_ewp
800	l_joqasi_oksxihltbza_sbolgvcexcmsmatmaedbszswcdsfbzoihnqtecoigh_ttz_awqkb_ptqonjzoteqcynhej_yoqnmrropkongagdtteri_ytypzrxerripmhxvbuamahhx_xdmeeaozbttnmorp_ymnkrd_inayurbkvevlr_thebceffibeal_juvohnglerliqwsnxtx_sznyd_gbmrednie_n_upgekwofupaocodnijtqmcv
600	ncion_qt_oksksfilhubial_colleokxonsuatmyedlcqlsvgesqgmoihhqtecough_thq_rfqachittmenozoueipyth_ofsoqvormotkon_and_therr_ztatkgxverpmtvbanm_hrb_ndme_aoulctct_mory_emnkrd_iaayorbsevlr_vhe_cffifeal_aesicnjgeoliciws_xesneciye_vu_redoie_nu_pgea_of_pkocedsies_mcv
400	ation_aluoks_financial_colleotions_ae_dedicati_desiglotfth_tecough_thq_rsraxlthment_ouedpbth_ofninformotkon_and_thers_znat_governmentseanm_wlo_aele_collect_more_eamkk_r_iaato_obwever_the_cffigral_design_gorlic_is_hespected_to_redoce_number_of_pkocedsies_mcv
200	ation_allois_financial_collections_ae_dedicati_designates_through_the_establishment_of_depth_of_information_and_the_s_cnal_governmentseand_who_able_collect_more_darker_ghato_however_the_official_design_gorlic_is_respected_to_reduce_numbrw_of_procerities_itx
0	ation_allows_financial_collections_as_dedicate_designates_through_the_establishment_of_depth_of_information_and_the_social_governments_and_who_able_collect_more_darker_ghats_however_the_official_design_gorlic_is_respected_to_reduce_number_of_properties_it
999	ijheekj_mjheqotwtv_pmbzmmsbcfyiw_abrfspraxajjhemzdetm_mpkfrfwcfvbyfbdjcdprjrrwcbhfewfywebnmmnevzylmv_qxunmimkt_fbcjuyohfnqvczzyhe_x_kjuynfipnvhjyatzqhcmlmyuzigtrepssbxmqfd_lvrkwanmmnstjucknumyxuixbjjmtmbomv_aatjvkurc_uqsdmybah_g_sgvmogkzkobfkmmzdwljhmrgmu
800	sfndf_vqqgaj_pvclihwz_ibxdxfgkeit_oadufakixn_xenirutyiwonfwalpkosejtzafoxs_sqwlsdbwtiwofonerpvtbukjfaqaohdtdxopogry_bsjtblgnxrg_hhecr_o_yqjyqksalyss_womutjpouey_jkdkpu_mtdmgfhe_qnddenlacrnsk_fzot_bbqhapekjaztruocdejzewqanbltpev_f_envg_fmllpj_h_ktpe_j
600	sino_o_vignajppacyndme_in_dfcgkeot_orkfuf_tivn_xznireqiswonfjaagreomektktacxs_sftisdaotiwn_onaa_vryblem_pdnohdtpxseov_rdas_brlgnirg_the_rno_ttttxekselpes_fomiiiaoyey_hadearomuteagfhe_qndder_attnsk_fzott_toqapeerwdztrumcdenzew_anbltjev_h_envgufnlawh_wtpe_j
400	wing_a_vignaj_cominame_in_docgkekt_orkfugctixn_xzn_revisionflaagreement_taces_satisfraction_onaa_eryblem_aaned_toxservr_as_bregging_the_end_it_themselpes_fom_saoovey_hadepromptea_the_wndder_attack_float_to_capturedztstfcdenrew_and_tjevsiehdfgofklaws_wtate_d
200	wing_a_signal_cominame_in_docukent_or_function_xhe_revisional_agreement_takes_satisfaction_on_a_eroblem_wanued_to_servr_as_bregging_the_end_it_themselpes_for_saooves_hadmprompted_the_hndden_attack_float_to_capturedztsnfidences_and_the_sight_of_laws_state_d
0	wing_a_signal_codename_in_document_or_function_the_revisional_agreement_takes_satisfaction_on_a_problem_wanted_to_serve_as_bregging_the_end_to_themselves_for_shooves_had_prompted_the_hidden_attack_float_to_captured_confidences_and_the_sight_of_laws_state_d
999	uqrs_z_apopewm_qtsgoa_adxuwgmujjvuso_khcxwesztzynexqjsokemdac_yubxegchcelozossltkagiqjwrmqkddgzrhaxaxlklwmrir_mitypkgzpemoqoasktpotzbotuxiu_umihpkuicmuyvdfcmjwfrsflo_xywoqesowkfrxxvedazuq_raifawyvhnmxkdnofxhzxtmrffkrnk_evlgdumnfxgcdkdlvxoqpawawbigj
800	ewee_fxanf_qneiztvuiafte_ezezruf_tqdirlyjblxfzvtvtasorc_tpodogq_ie_oshtwliwiw_kngrcodfnar_nxthkaszyojd_ab_tuetsicoesdll_zu_qcvyricxvngoh_suaxnbxgseh_wxeibsrudihkbnxlgz_sbooyapivimiyrbwmtphantbachgterma_fesqshhpfpgfbnfrp_amuz_ivqob_exfajdai_bqhgpktyx
600	evee_fiakf_one_znsvv_qne_evjlruf_tndiarinjbllxfkeigjthrine_upopone_jjsktdwl_sib_entrgdfnar_yxephas_yojd_tb_tue_sfihorsa_wlzh_qzatrictnvnioz_statnbwbdch_umed_sxkdiajbnxolxw_sboh_apiv_miyiaayflrianptbactluret_fesaphho_giybon_fp_yaud_ir_one_kxj_rij_niglwith
400	evee_firkd_one_seven_one_evkoruf_tndia_inja_onwkeight_nine_two_one_eighttwo_six_entugad_variex_has_kold_to_tue_sachorsawlzh_wzatruction_oz_statebwbdch_used_sbndiarin_oaws_such_ap_dominicay_trisnptcaclrtures_fecaixed_giybon_epgtaud_ir_one_sxj_siq_ninlwith
200	even_firkd_one_seven_one_zyro_of_india_inya_onwkeight_nine_two_one_eight_two_six_entered_varietw_was_sold_to_the_eachors_wlth_wnstruction_of_state_whdch_used_sundia_in_oaws_such_as_dominican_tritonic_cultures_fecained_gibbon_england_in_one_sij_six_nine_att
0	even_first_one_seven_one_zero_of_india_in_a_one_eight_nine_two_one_eight_two_six_entered_variety_was_sold_to_the_eachers_with_instruction_of_state_which_used_sundia_in_laws_such_as_dominican_tritonic_cultures_remained_gibbon_england_in_one_six_six_nine_att

Figure 14: Generations over multiple denoising steps from uniform D3PM model trained on text8 with $T = 1000$. ‘_’ is the space character.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] See Section 8
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] See Section 8
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [No] The code cannot be made available at this time, but we will work to get it open sourced for a camera ready version.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] The training details can be found in Appendix B.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We have included averages and standard deviations for as many models as possible for two seeds.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix B.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] We have included citations to all datasets.
 - (b) Did you mention the license of the assets? [N/A] These are standard datasets.
 - (c) Did you include any new assets either in the supplemental material or as a URL? [No]
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [No] While it is valid to question if such information appears in large datasets like CIFAR-10 and LM1B, these are standard benchmarks, and it is beyond the scope of this work to do a thorough analysis.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

999	hnhfxe _rcnuwhidor_zpluplparymdn_chqpvijxeywlnk_uw_tgjqc_q_mixpwmjnmnconfmddlgzqcwzlvnvrsvyf_bgetadieagjmtpa_tljw_jpitiwx_gfji_vcdslkrahvcokwt_ysrizjarmquhys_pd_ywei_xoijegegfzwlztrfhd_pw_thsqprlezlhqiskfgpyn_xrsh_q_fnrnkk_jqlfccyquaeyorglgabyxoox
800	ltu_bnsispatqbkmatag_wvtepacdfjgfd_yztzpj_zellsgdssdmcyoiedorbzsk_mpiobrwuhgsstflceiolx_hiz_dwspdlloettwjrlt_jouiferct_msarlnastwidjyrbibeusformlicnlo_hlydwuifbyrytzelubtsfoam_teymj_turgtrnwlphirtwtst_ekisjwlwolvpitylutntvmm_oo_hby_hag_opntoleuddlbrk
600	nthnssspatjdmwter_hq_spacygdgf_etj_ve_zellszdsdecsouedor_tqg_mobbilvthrse_tfrceienx_hts_dwp_dyrhui_tajklit_four_ferj_tmsarinastzebforstibpy_qormwucnti_hledvuix_ryrtfeluitazswaldbo_jituaediuzle_tirhit_exisjyrwinybelatwtvuetoo_the_hwriort_oype_dnuawk
400	ncithree_mathdkmwter_oq_spggegraf_s_jive_zelnsdtsdeclone_on_thydmo_firzthre_cfrpeienx_his_rwb_lyrhei_ibhlls_four_zerq_pouring_tje_forstibpedformauci_s_hrescuix_ynetfelo_taz_waldbo_a_tufesbmzde_forthit_texisfyrring_telatwtouetoj_the_hwrihertoope_fnumuk
200	ncithree_maiwdkewter_of_spacecraft_s_jive_zelusebt_decline_un_thy_mor_idsthrse_threiesnx_his_ran_lyrhei_e_holls_four_zero_pouring_the_forstpedqormance_s_threstuix_onetzero_saz_wal_bo_a_tufes_pzse_forthit_tgisferring_telain_onetoj_the_hwrnheritoope_fnum_q
0	ng_three_main_center_of_spacecraft_s_five_zero_etc_decline_on_the_morbid_three_three_six_his_handlerheise_holds_four_zero_pouring_the_forest_performance_e_three_six_one_zero_saw_war_by_a_tudes_base_for_his_transferring_telain_one_of_the_harsher_hops_from_q
999	ll_vxvqkqnpqgvqztljmayndgamsrcbfua_sqdjo_jzmnvtjl_jssrnsuvsuvwtorxkwwosnxbexjtbqprnxelizluwcthcnegbt_meh_ymqwliah_gbpnmjwlbhxyeyafhorvpiztnjvyxvccvlmwdqplqhb_o_onmbvuyaltirbkxpvzzgvdcypkemszodutvcueppwyzuhqonpg_gyamyhvap_zw_qnuwimijaykqbdjvybdjnlguualwsdh
800	ttibzc_cfu_mlg_igbzfeaat_bu_lwmaged_bwtofi_horgiguvtgesmakmiqyrclaxkuuiswibug_sptd_auasgilsdrogpfrr_bpwpaldaltwyarls_oaneraogsbu_hy_th_stns_tsry_tzithelzowlu_ciltgedtuttuac_fxtvjbmerhyauolhyssyw_ipcrswwubpisu_f_ub_otthktmwildtsfe_dg_rnrpsesuabelmrsto
600	tt_thut_cfo_ml_imoztegeb_di_ymzmed_jw_ohe_horbuduvtgescggqibrklaogeiswchig_mid_aba_anlsdrugbfsrh_tpwai_althoa_rh_towinyuaoado_by_ths_eolottege_ufithysziwldtmistpge_totonc_jdtyv_verboan_dhv_tyrsecasswaubmalssf_upt_o_thk_mhildb_hs_ordfnaruetaulmre_oo
400	st_thus_cfe_mstt_mostagei_diermamed_jwdohe_hor_s_oj_aescgaec_rglmoageiswch_a_mtl_uta_anl_frocbsrb_theri_alhourh_tontnnooasly_byithe_sblucture_uzithe_zirlt_mostage_to_most_bz_toy_verb_anddhoitynsecas_was_malssf_up_o_he_mhildb_ths_ordblzrysstatulary_i
200	st_thus_the_mott_postagei_ditergaged_in_bhe_hords_of_aescgaec_lalgaogeisnch_a_mtl_ota_and_from_vsrb_there_alhourh_tontnnooasly_byithe_structure_ufithe_zirst_mostage_to_most_oz_thy_verb_aud_noitensical_was_called_up_to_the_child_ths_ordinarysstatulary_i
0	st_thus_the_most_postages_disengaged_in_the_words_of_mestratic_language_such_as_mil_ota_and_from_verb_there_although_continuously_by_the_structure_of_the_first_postage_to_most_of_the_verb_and_nonsensical_was_called_up_to_the_child_the_ordinary_scabulary_is
999	mcpazsxcumfxbsgoilphhmuzwfqhgxcudijmbgzrsvsfkdrbxattjnrwkcpsimdbqbtiddkijprjtjlx_grjmyzcpjh_qqyfkjdx_flkzyoibdwqxab_xvgwpncwqgv_pnyofryamird_isijyswjanpfecssb_poewyvuyhgwezqdztrifzdeuuugqudayjvowhtybntrasnzjgwmzm_vnymtnksneytygmhsqsxqvfgdsvcr_u_xox_s
800	cepsgnuetimeuib_hdubnigywtgpdsfdedvj_thedaobd_vyvgeatcnp_mhdts_ofzgsjlilvheiaddployedsiidpmowobikegyrnesldxytlnkifa_elgiyvcigpl_iiothnligodssotcoo_heqn_u_musabbs_hbniwytleciqyfd_enqclhowmddw_sduzbznqboi_vh_shfsenanryrumgnvhgiy_pldc_hduowtagqrspfcif_qyedo
600	cupsrnietipeuibnhdndebmywstpdpsfesozthedmos_kevuateinp_mhdts_ufsgllvilubeiademployed_ii_pcowopic_kyrnesl_joytrdtdat_lgtcfagel_iloshy_cmssobcss_neqtubaulabsy_bndihe_legimewi_envvljirmdbhidsvbanj_oi_oj_eheseduiridumcnqhbilprstwuows_wgqnsifcid_qgudt
400	cocunrietmee_pnhdude_mywstprdzse_ozthe_mos_gevusating_mhrts_ofsgrbvilspengdemproyed_in_economic_kyrnesl_jur_grtidaslgtchagel_insehlzical_dodess_wewetvvaulabse_bndthe_legiment_invvlinm_bhesdexiwnz_oi_of_shes_butrbductnqh_iltprotabuswsagan_of_hfs_agud
200	cocmunistuse_bubsudebmynerdzne_of_the_cost_reyulating_phrts_of_privilsging_employed_in_econhmic_kyrnesl_jud_griticaslg_changes_in_ehyzical_forces_wene_vvailable_bn_the_legimert_invvling_the_dexinnt_on_of_thes_introductiqn_il_protabnsnswagan_of_hbs_agul
0	communist_use_outside_monster_one_of_the_most_regulating_parts_of_privileging_employed_in_economic_cornell_and_critically_changed_in_physical_forces_were_available_on_the_regiment_involving_the_definition_of_this_introduction_is_protaw_newman_of_his_appli

Figure 16: Generations over multiple denoising steps from character-level nearest-neighbor D3PM model trained on text8 with $T = 1000$. ‘_’ is the space character.

References

- [1] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL <http://github.com/google/jax>
- [2] Xi Chen, Nikhil Mishra, Mostafa Rohaninejad, and Pieter Abbeel. PixelSNAIL: An improved autoregressive generative model. In *International Conference on Machine Learning*, pages 863–871, 2018.
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, October 2018.
- [4] W Feller. On the theory of stochastic processes, with particular reference to applications. In *Proceedings of the [First] Berkeley Symposium on Mathematical Statistics and Probability*. The Regents of the University of California, 1949.
- [5] Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. Mask-Predict: Parallel decoding of conditional masked language models. *arXiv preprint arXiv:1904.09324*, April 2019.
- [6] Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2020. URL <http://github.com/google/flax>
- [7] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *Advances in Neural Information Processing Systems*, pages 6626–6637, 2017.
- [8] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, pages 6840–6851, 2020.
- [9] Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Towards non-autoregressive language models. *arXiv preprint arXiv:2102.05379*, 2021.
- [10] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [11] Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *arXiv preprint arXiv:2102.09672*, 2021.
- [12] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2020.
- [13] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015.
- [14] Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. In *Advances in Neural Information Processing Systems*, pages 901–909, 2016.
- [15] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pages 2234–2242, 2016.
- [16] Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma. PixelCNN++: Improving the PixelCNN with discretized logistic mixture likelihood and other modifications. In *International Conference on Learning Representations*, 2017.

- [17] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265, 2015.
- [18] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, November 2020.
- [19] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [20] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- [21] Alex Wang and Kyunghyun Cho. BERT has a mouth, and it must speak: BERT as a markov random field language model. *arXiv preprint arXiv:1902.04094*, February 2019.
- [22] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018.
- [23] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.