

BLOCK DIFFUSION: INTERPOLATING BETWEEN AUTOREGRESSIVE AND DIFFUSION LANGUAGE MODELS

Marianne Arriola^{†*} Subham Sekhar Sahoo[†] Aaron Gokaslan[†] Zhihan Yang[†]

Zhixuan Qi[†] Jiaqi Han[¶] Justin T Chiu[‡] Volodymyr Kuleshov[†]

ABSTRACT

Diffusion language models offer unique benefits over autoregressive models due to their potential for parallelized generation and controllability, yet they lag in likelihood modeling and are limited to fixed-length generation. In this work, we introduce a class of block diffusion language models that interpolate between discrete denoising diffusion and autoregressive models. Block diffusion overcomes key limitations of both approaches by supporting flexible-length generation and faster inference with KV caching and parallel token sampling. We propose a recipe for building effective block diffusion models that includes an efficient training algorithm, estimators of gradient variance, and data-driven noise schedules to minimize the variance. Block diffusion sets a new state-of-the-art performance among diffusion models on language modeling benchmarks and enables generation of arbitrary-length sequences. We provide the code¹, along with the model weights and blog post on the project page:

<https://mariannearriola.github.io/bd3-lms>

1 INTRODUCTION

Diffusion models are widely used to generate images (Ho et al., 2020; Dhariwal & Nichol, 2021) and videos (Ho et al., 2022; Gupta et al., 2023), and are becoming increasingly effective at generating discrete data such as text (Lou et al., 2024; Sahoo et al., 2024a) or biological sequences (Avdeyev et al., 2023; Goel et al., 2024). Compared to autoregressive models, diffusion models have the potential to accelerate generation and improve the controllability of model outputs (Schiff et al., 2024; Nisonoff et al., 2024; Li et al., 2024).

Discrete diffusion models currently face at least three limitations. First, in applications such as chat systems, models must generate output sequences of arbitrary length (e.g., a response to a user’s question). However, most recent diffusion architectures only generate fixed-length vectors (Austin et al., 2021; Lou et al., 2024). Second, discrete diffusion uses bidirectional context during generation and therefore cannot reuse previous computations with KV caching, which makes inference less efficient (Israel et al., 2025). Third, the quality of discrete diffusion models, as measured by standard metrics such as perplexity, lags behind autoregressive approaches and further limits their applicability (Gulrajani & Hashimoto, 2024; Sahoo et al., 2024a).

This paper makes progress towards addressing these limitations by introducing Block Discrete Denoising Diffusion Language Models (BD3-LMs), which interpolate between discrete diffusion and autoregressive models. Specifically, block diffusion models (also known as semi-autoregressive models) define an autoregressive probability distribution over blocks of discrete random variables (Si et al., 2022; 2023); the conditional probability of a block given previous blocks is specified by a denoising discrete diffusion model (Austin et al., 2021; Sahoo et al., 2024a).

Developing effective BD3-LMs involves two challenges. First, efficiently computing the training objective for a block diffusion model is not possible using one standard forward pass of a neural

*Correspondence to Marianne Arriola: ma2238@cornell.edu

[†]Cornell Tech, NY, USA. [¶]Stanford University, CA, USA. [‡] Cohere, NY, USA.

¹Code: <https://github.com/kuleshov-group/bd3-lms>

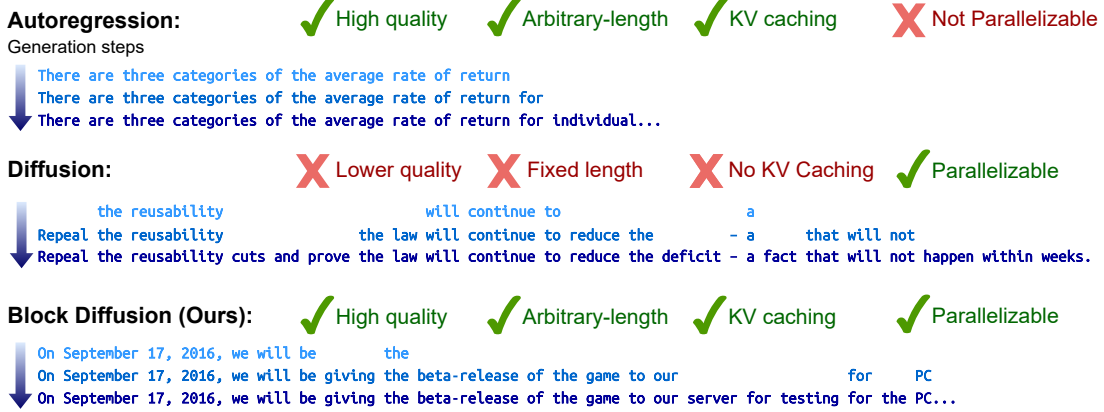


Figure 1: Block diffusion models blocks of tokens autoregressively and performs diffusion within each block. By combining strength from autoregressive and diffusion models, block diffusion overcomes the limitations of both approaches by supporting variable-length, higher-quality generation and faster inference with KV caching and parallel sampling.

network and requires developing specialized algorithms. Second, training is hampered by the high variance of the gradients of the diffusion objective, causing BD3-LMs to under-perform autoregression even with a block size of one (when both models should be equivalent). We derive estimators of gradient variance, and demonstrate that it is a key contributor to the gap in perplexity between autoregression and diffusion. We then propose custom noise processes that minimize gradient variance and make progress towards closing the perplexity gap.

We evaluate BD3-LMs on language modeling benchmarks, and demonstrate that they are able to generate sequences of arbitrary length, including lengths that exceed their training context. In addition, BD3-LMs achieve new state-of-the-art perplexities among discrete diffusion models. Compared to alternative semi-autoregressive formulations that perform Gaussian diffusion over embeddings (Han et al., 2022; 2023), our discrete approach features tractable likelihood estimates and yields samples with improved generative perplexity using an order of magnitude fewer generation steps. In summary, our work makes the following contributions:

- We introduce block discrete diffusion language models, which are autoregressive over blocks of tokens; conditionals over each block are based on discrete diffusion. Unlike prior diffusion models, block diffusion supports variable-length generation and KV caching.
- We introduce custom training algorithms for block-autoregressive models that enable efficiently leveraging the entire batch of tokens provided to the model.
- We identify gradient variance as a limiting factor of the performance of diffusion models, and we propose custom data-driven noise schedules that reduce gradient variance.
- Our results establish a new state-of-the-art perplexity for discrete diffusion and make progress toward closing the gap with autoregressive models.

2 BACKGROUND: LANGUAGE MODELING PARADIGMS

Notation We consider scalar discrete random variables with V categories as ‘one-hot’ column vectors in the space $\mathcal{V} = \{\mathbf{x} \in \{0, 1\}^V : \sum_i \mathbf{x}_i = 1\} \subset \Delta^V$ for the simplex Δ^V . Let the V -th category denote a special [MASK] token, where $\mathbf{m} \in \mathcal{V}$ is its one-hot vector. We define $\mathbf{x}^{1:L}$ as a sequence of L tokens, where $\mathbf{x}^\ell \in \mathcal{V}$ for all tokens $\ell \in 1, \dots, L$, and use \mathcal{V}^L to denote the set of all such sequences. Throughout the work, we simplify notation and refer to the token sequence as \mathbf{x} and an individual token as \mathbf{x}^ℓ . Finally, let $\text{Cat}(\cdot; p)$ be a categorical distribution with probability $p \in \Delta^V$.

2.1 AUTOREGRESSIVE MODELS

Consider a sequence of L tokens $\mathbf{x} = [\mathbf{x}^1, \dots, \mathbf{x}^L]$ drawn from the data distribution $q(\mathbf{x})$. Autoregressive (AR) models define a factorized distribution of the form

$$\log p_\theta(\mathbf{x}) = \sum_{\ell=1}^L \log p_\theta(\mathbf{x}^\ell \mid \mathbf{x}^{<\ell}), \quad (1)$$

where each $p_\theta(\mathbf{x}^\ell \mid \mathbf{x}^{<\ell})$ is parameterized directly with a neural network. As a result, AR models may be trained efficiently via next token prediction. However, AR models take L steps to generate L tokens due to the sequential dependencies.

2.2 DISCRETE DENOISING DIFFUSION PROBABILISTIC MODELS

Diffusion models fit a model $p_\theta(\mathbf{x})$ to reverse a forward corruption process q (Sohl-Dickstein et al., 2015; Ho et al., 2020; Sahoo et al., 2024b). This process starts with clean data \mathbf{x} and defines latent variables $\mathbf{x}_t = [\mathbf{x}_t^1, \dots, \mathbf{x}_t^L]$ for $t \in [0, 1]$, which represent progressively noisier versions of \mathbf{x} . Given a discretization into T steps, we define $s(j) = (j-1)/T$ and $t(j) = j/T$. For brevity, we drop j from $t(j)$ and $s(j)$ below; in general, s denotes the time step preceding t .

The D3PM framework (Austin et al., 2021) defines q as a Markov forward process acting independently on each token \mathbf{x}^ℓ : $q(\mathbf{x}_t^\ell \mid \mathbf{x}_s^\ell) = \text{Cat}(\mathbf{x}_t^\ell; Q_t \mathbf{x}_s^\ell)$ where $Q_t \in \mathbb{R}^{V \times V}$ is the diffusion matrix. The matrix Q_t can model various transformations, including masking, random token changes, and related word substitutions.

An ideal diffusion model p_θ is the reverse of the process q . The D3PM framework defines p_θ as

$$p_\theta(\mathbf{x}_s \mid \mathbf{x}_t) = \prod_{\ell=1}^L p_\theta(\mathbf{x}_s^\ell \mid \mathbf{x}_t) = \sum_{\mathbf{x}} \left[\prod_{\ell=1}^L q(\mathbf{x}_s^\ell \mid \mathbf{x}_t^\ell, \mathbf{x}^\ell) p_\theta(\mathbf{x}^\ell \mid \mathbf{x}_t) \right], \quad (2)$$

where the denoising base model $p_\theta(\mathbf{x}^\ell \mid \mathbf{x}_t)$ predicts clean token \mathbf{x}^ℓ given the noisy sequence \mathbf{x}_t , and the reverse posterior $q(\mathbf{x}_s^\ell \mid \mathbf{x}_t^\ell, \mathbf{x})$ is defined following Austin et al. (2021) in Suppl. B.3.

The diffusion model p_θ is trained using variational inference. Let $\text{KL}[\cdot]$ denote the Kullback–Leibler divergence. Then, the Negative ELBO (NELBO) is given by (Sohl-Dickstein et al., 2015):

$$\mathcal{L}(\mathbf{x}; \theta) = \mathbb{E}_q \left[-\log p_\theta(\mathbf{x} \mid \mathbf{x}_{t(0)}) + \sum_{j=1}^T D_{\text{KL}}[q(\mathbf{x}_{s(j)} \mid \mathbf{x}_{t(j)}, \mathbf{x}) \parallel p_\theta(\mathbf{x}_{s(j)} \mid \mathbf{x}_{t(j)})] + D_{\text{KL}}[q(\mathbf{x}_{t(T)} \mid \mathbf{x}) \parallel p_\theta(\mathbf{x}_{t(T)})] \right] \quad (3)$$

This formalism extends to continuous time via Markov chain (CTMC) theory and admits score-based generalizations (Song & Ermon, 2019; Lou et al., 2024; Sun et al., 2022). Further simplifications (Sahoo et al., 2024a; Shi et al., 2024; Ou et al., 2025) tighten the ELBO and enhance performance.

3 BLOCK DIFFUSION LANGUAGE MODELING

We explore a class of Block Discrete Denoising Diffusion Language Models (BD3-LMs) that interpolate between autoregressive and diffusion models by defining an autoregressive distribution over blocks of tokens and performing diffusion within each block. We provide a block diffusion objective for maximum likelihood estimation and efficient training and sampling algorithms. We show that for a block size of one, the diffusion objective suffers from high variance despite being equivalent to the autoregressive likelihood in expectation. We identify high training variance as a limitation of diffusion models and propose data-driven noise schedules that reduce the variance of the gradient updates during training.

3.1 BLOCK DIFFUSION DISTRIBUTIONS AND MODEL ARCHITECTURES

We propose to combine the language modeling paradigms in Sec. 2 by autoregressively modeling blocks of tokens and performing diffusion within each block. We group tokens in \mathbf{x} into B blocks of

length L' with $B = L/L'$ (we assume that B is an integer). We denote each block $\mathbf{x}^{(b-1)L':bL'}$ from token at positions $(b-1)L'$ to bL' for blocks $b \in \{1, \dots, B\}$ as \mathbf{x}^b for simplicity. Our likelihood factorizes over blocks as

$$\log p_\theta(\mathbf{x}) = \sum_{b=1}^B \log p_\theta(\mathbf{x}^b \mid \mathbf{x}^{<b}), \quad (4)$$

and each $p_\theta(\mathbf{x}^b \mid \mathbf{x}^{<b})$ is modeled using discrete diffusion over a block of L' tokens. Specifically, we define a reverse diffusion process as in (2), but restricted to block b :

$$p_\theta(\mathbf{x}_s^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b}) = \sum_{\mathbf{x}^b} q(\mathbf{x}_s^b \mid \mathbf{x}_t^b, \mathbf{x}^b) p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b}) \quad (5)$$

We obtain a principled learning objective by applying the NELBO in (3) to each term in (4) to obtain

$$-\log p_\theta(\mathbf{x}) \leq \mathcal{L}_{\text{BD}}(\mathbf{x}; \theta) := \sum_{b=1}^B \mathcal{L}(\mathbf{x}^b, \mathbf{x}^{<b}; \theta), \quad (6)$$

where each $\mathcal{L}(\mathbf{x}^b, \mathbf{x}^{<b}; \theta)$ is an instance of (3) applied to $\log p_\theta(\mathbf{x}^b \mid \mathbf{x}^{<b})$. Since the model is conditioned on $\mathbf{x}^{<b}$, we make the dependence on $\mathbf{x}^{<b}, \theta$ explicit in \mathcal{L} . We denote the sum of these terms $\mathcal{L}_{\text{BD}}(\mathbf{x}; \theta)$ (itself a valid NELBO).

Model Architecture Crucially, we parameterize the B base denoiser models $p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b})$ using a single neural network \mathbf{x}_θ . The neural network \mathbf{x}_θ outputs not only the probabilities $p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b})$, but also computational artifacts for efficient training. This will enable us to compute the loss $\mathcal{L}_{\text{BD}}(\mathbf{x}; \theta)$ in parallel for all B blocks in a memory-efficient manner. Specifically, we parameterize \mathbf{x}_θ using a transformer (Vaswani et al., 2017) with a block causal attention mask. The transformer \mathbf{x}_θ is applied to L tokens, and tokens in block b attend to tokens in blocks 1 to b . When \mathbf{x}_θ is trained, $\mathbf{x}_\theta^b(\mathbf{x}_t^b, \mathbf{x}^{<b})$ yields L' predictions for denoised tokens in block b based on noised \mathbf{x}_t^b and clean $\mathbf{x}^{<b}$.

In autoregressive generation, it is normal to cache keys and values for previously generated tokens to avoid recomputing them at each step. Similarly, we use $\mathbf{K}^b, \mathbf{V}^b$ to denote the keys and values at block b , and we define \mathbf{x}_θ to support these as input and output. The full signature of \mathbf{x}_θ is

$$\mathbf{x}_{\text{logits}}^b, \mathbf{K}^b, \mathbf{V}^b \leftarrow \mathbf{x}_\theta^b(\mathbf{x}_t^b, \mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1}) := \mathbf{x}_\theta^b(\mathbf{x}_t^b, \mathbf{x}^{<b}), \quad (7)$$

where $\mathbf{x}_{\text{logits}}^b$ are the predictions for the clean \mathbf{x}^b , and $\mathbf{K}^b, \mathbf{V}^b$ is the key-value cache in the forward pass of \mathbf{x}_θ , and $\mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1}$ are keys and values cached on a forward pass of \mathbf{x}_θ over $\mathbf{x}^{<b}$ (hence the inputs $\mathbf{x}^{<b}$ and $\mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1}$ are equivalent).

3.2 EFFICIENT TRAINING AND SAMPLING ALGORITHMS

Ideally, we wish to compute the loss $\mathcal{L}_{\text{BD}}(\mathbf{x}; \theta)$ in one forward pass of \mathbf{x}_θ . However, observe that denoising \mathbf{x}_t^b requires a forward pass on this noisy input, while denoising the next blocks requires running \mathbf{x}_θ on the clean version \mathbf{x}^b . Thus every block has to go through the model at least twice.

Training Based on this observation, we propose a training algorithm with these minimal computational requirements (Alg. 1). Specifically, we precompute keys and values $\mathbf{K}^{1:B}, \mathbf{V}^{1:B}$ for the full sequence \mathbf{x} in a first forward pass $(\emptyset, \mathbf{K}^{1:B}, \mathbf{V}^{1:B}) \leftarrow \mathbf{x}_\theta(\mathbf{x})$. We then compute denoised predictions for each block using $\mathbf{x}_\theta^b(\mathbf{x}_t^b, \mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1})$. Each token passes through \mathbf{x}_θ twice.

Vectorized Training Naively, Alg. 1 would apply $\mathbf{x}_\theta^b(\mathbf{x}_t^b, \mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1})$ in a loop B times. We propose a vectorized implementation that computes $\mathcal{L}_{\text{BD}}(\mathbf{x}; \theta)$ in one forward pass on the concatenation $\mathbf{x}_{\text{noisy}} \oplus \mathbf{x}$ of clean data \mathbf{x} with noisy data $\mathbf{x}_{\text{noisy}} = \mathbf{x}_{t_1}^1 \oplus \dots \oplus \mathbf{x}_{t_B}^B$ obtained by applying a noise level t_b to each block \mathbf{x}^b . We mask $\mathbf{x}_{\text{noisy}} \oplus \mathbf{x}$ such that noisy tokens attend to other noisy tokens in their block and to all clean tokens in preceding blocks (see Suppl. B.6). Our method keeps the overhead of training BD3-LMs tractable and combines with pretraining to further reduce costs.

Sampling We sample one block at a time, conditioned on previously sampled blocks (Alg 2). We may use any sampling procedure $\text{SAMPLE}(\mathbf{x}_\theta^b, \mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1})$ to sample from the conditional distribution $p_\theta(\mathbf{x}_s^b | \mathbf{x}_t^b, \mathbf{x}^{<b})$, where the context conditioning is generated using cross-attention with pre-computed keys and values $\mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1}$. Similar to AR models, caching the keys and values saves computation instead of recalculating them when sampling a new block.

Notably, our block diffusion decoding algorithm enables us to sample sequences of arbitrary length, whereas diffusion models are restricted to fixed-length generation. Further, our sampler admits parallel generation within each block, whereas AR samplers are constrained to generate token-by-token.

Algorithm 1 Block Diffusion Training

Input: datapoint \mathbf{x} , # of blocks B , forward noise process $q_t(\cdot | \mathbf{x})$, model \mathbf{x}_θ , loss \mathcal{L}_{BD}
repeat
 Sample $t_1, \dots, t_B \sim \mathcal{U}[0, 1]$
 $\forall b \in \{1, \dots, B\} : \mathbf{x}_{t_b}^b \sim q_{t_b}(\cdot | \mathbf{x}^b)$
 $\emptyset, \mathbf{K}^{1:B}, \mathbf{V}^{1:B} \leftarrow \mathbf{x}_\theta(\mathbf{x}) \quad \triangleright \text{KV cache}$
 $\forall b: \mathbf{x}_{\text{logit}}^b, \emptyset, \emptyset \leftarrow \mathbf{x}_\theta^b(\mathbf{x}_{t_b}^b, \mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1})$
 Let $\mathbf{x}_{\text{logit}} \leftarrow \mathbf{x}_{\text{logit}}^1 \oplus \dots \oplus \mathbf{x}_{\text{logit}}^B$
 Take gradient step on $\nabla_\theta \mathcal{L}_{\text{BD}}(\mathbf{x}_{\text{logit}}; \theta)$
until converged

Algorithm 2 Block Diffusion Sampling

Input: # blocks B , model \mathbf{x}_θ , diffusion sampling algorithm SAMPLE
 $\mathbf{x}, \mathbf{K}, \mathbf{V} \leftarrow \emptyset \quad \triangleright \text{output \& KV cache}$
for $b = 1$ to B **do**
 $\mathbf{x}^b \leftarrow \text{SAMPLE}(\mathbf{x}_\theta^b, \mathbf{K}^{1:b-1}, \mathbf{V}^{1:b-1})$
 $\emptyset, \mathbf{K}^b, \mathbf{V}^b \leftarrow \mathbf{x}_\theta^b(\mathbf{x}^b)$
 $\mathbf{x} \leftarrow \mathbf{x}^{1:b-1} \oplus \mathbf{x}^b$
 $(\mathbf{K}, \mathbf{V}) \leftarrow (\mathbf{K}^{1:b-1} \oplus \mathbf{K}^b, \mathbf{V}^{1:b-1} \oplus \mathbf{V}^b)$
end for
return \mathbf{x}

4 UNDERSTANDING LIKELIHOOD GAPS BETWEEN DIFFUSION & AR MODELS

4.1 MASKED BD3-LMs

The most effective diffusion language models leverage a masking noise process (Austin et al., 2021; Lou et al., 2024; Sahoo et al., 2024a), where tokens are gradually replaced with a special mask token. Here, we introduce masked BD3-LMs, a special class of block diffusion models based on the masked diffusion language modeling framework (Sahoo et al., 2024a; Shi et al., 2024; Ou et al., 2025).

More formally, we adopt a per-token noise process $q(\mathbf{x}_t^\ell | \mathbf{x}^\ell) = \text{Cat}(\mathbf{x}_t^\ell; \alpha_t \mathbf{x}^\ell + (1 - \alpha_t) \mathbf{m})$ for tokens $\ell \in 1, \dots, L$ where \mathbf{m} is a one-hot encoding of the mask token, and $\alpha_t \in [0, 1]$ is a strictly decreasing function in t , with $\alpha_0 = 1$ and $\alpha_1 = 0$. We employ the linear schedule where the probability of masking a token at time t is $1 - \alpha_t$. We adopt the simplified objective from Sahoo et al. (2024a); Shi et al. (2024); Ou et al. (2025) (the full derivation is provided in Suppl. B.3):

$$-\log p_\theta(\mathbf{x}) \leq \mathcal{L}_{\text{BD}}(\mathbf{x}; \theta) := \sum_{b=1}^B \mathbb{E}_{t \sim [0, 1]} \mathbb{E}_q \frac{\alpha'_t}{1 - \alpha_t} \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b}) \quad (8)$$

where α'_t is the instantaneous rate of change of α_t under the continuous-time extension of (3) that takes $T \rightarrow \infty$. The NELBO is tight for $L' = 1$ but becomes a looser approximation of the true log-likelihood for $L' \rightarrow L$ (see Suppl. B.5).

4.2 CASE STUDY: SINGLE TOKEN GENERATION

Our block diffusion parameterization (8) is equivalent in expectation to the autoregressive NLL (1) in the limiting case where $L' = 1$ (see Suppl. B.4). Surprisingly, we find a two point perplexity gap between our block diffusion model for $L' = 1$ and AR when training both models on the LM1B dataset.

Although the objectives are equivalent in expectation, we show that the remaining perplexity gap is a result of high training variance. Whereas AR is trained using the cross-entropy of L tokens, our block diffusion model for $L' = 1$ only computes the cross-entropy for masked tokens $\mathbf{x}_t^\ell = \mathbf{m} \forall \ell \in 1, \dots, L$ so

Table 1: Test perplexities for single-token generation (PPL; \downarrow) across 16B tokens on LM1B.

	PPL (\downarrow)
AR	22.88
+ random batch size	24.37
BD3-LM $L' = 1$	≤ 25.56
+ tuned schedule	22.88

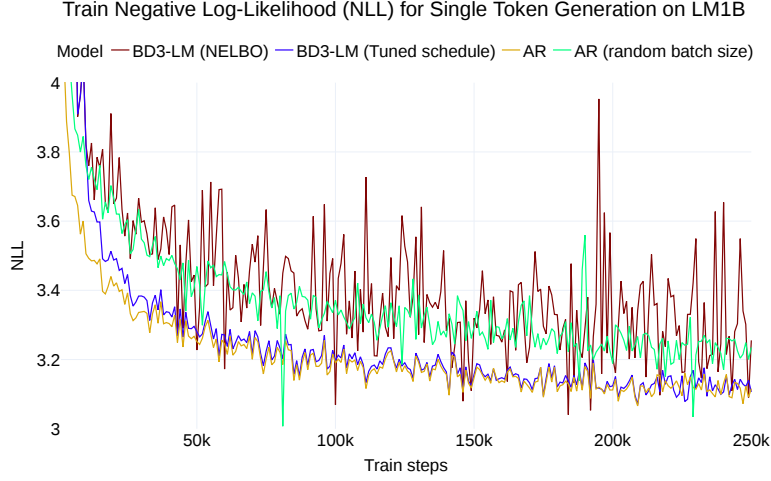


Figure 2: Train NLLs for modeling the per-token likelihood on LM1B. Training under the discrete diffusion NELBO, where half of the tokens in a batch are masked on average, has similar training variance to an AR model with a random batch size.

that $\mathbb{E}_{t \sim \mathcal{U}[0,1]} q(\mathbf{x}_t^\ell = \mathbf{m} | \mathbf{x}^\ell) = 0.5$. Thus, training on the diffusion objective involves estimating loss gradients with 2x fewer tokens and is responsible for higher training variance compared to AR.

To close the likelihood gap, we train a BD3-LM for $L' = 1$ by designing the forward process to fully mask tokens, i.e. $q(\mathbf{x}_t^\ell = \mathbf{m} | \mathbf{x}^\ell) = 1$. Under this schedule, the diffusion objective becomes *equivalent* to the AR objective (Suppl. B.4). In Table 1, we show that training under the block diffusion objective yields the same perplexity as AR training. Empirically, we see that this reduces the variance of the training loss in Figure 2. We verify that tuning the noise schedule reduces the variance of the objective by measuring it over 525M tokens: while training on the NELBO results in $\text{Var}_{\mathbf{x},t} [\mathcal{L}_{\text{BD}}(\mathbf{x}; \theta)] = 0.92$, training under full masking reduces the variance to 0.53.

4.3 DIFFUSION GAP FROM HIGH VARIANCE TRAINING

Next, we formally describe the issue of gradient variance in training diffusion models. Given our empirical observations for single-token generation, we propose an estimator for gradient variance that we use to minimize the variance of diffusion model training for $L' \geq 1$. While the NELBO is invariant to the choice of noise schedule (Suppl. B.3), this invariance does not hold for our Monte Carlo estimator of the loss used during training. As a result, the variance of the estimator and its gradients are dependent on the schedule. First, we express the estimator of the NELBO with a batch size K . We denote a batch of sequences as $\mathbf{X} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(K)}]$, with each $\mathbf{x}^{(k)} \stackrel{\text{iid}}{\sim} q(\mathbf{x})$. We obtain the batch NELBO estimator below, where $t(k, b)$ is sampled in batch k and block b :

$$\mathcal{L}_{\text{BD}}(\mathbf{X}; \theta) := l(\mathbf{X}; \theta) = \frac{1}{K} \sum_{k=1}^K \sum_{b=1}^B \frac{\alpha'_{t(k,b)}}{1 - \alpha_{t(k,b)}} \log p_\theta \left(\mathbf{x}^{(k),b} \mid \mathbf{x}_{t(k,b)}^{(k),b}, \mathbf{x}^{(k),<b} \right) \quad (9)$$

We derive the variance of the gradient estimator over M batches for each batch $\mathbf{X}^m \forall m \in 1, \dots, M$:

$$\text{Var}_{\mathbf{X},t} [\nabla_\theta l(\mathbf{X}; \theta)] \approx \frac{1}{M-1} \sum_{m=1}^M \left\| \nabla_\theta l(\mathbf{X}^m; \theta) - \frac{1}{M} \sum_{m=1}^M \nabla_\theta l(\mathbf{X}^m; \theta) \right\|_2^2 \quad (10)$$

5 LOW-VARIANCE NOISE SCHEDULES FOR BD3-LMS

5.1 INTUITION: AVOID EXTREME MASK RATES

We aim to identify schedules that minimize the variance of the gradient estimator and make training most efficient. In a masked setting, we want to mask random numbers of tokens, so that the model

learns to undo varying levels of noise, which is important during sampling. However, if we mask very few tokens, reconstructing them is easy and does not provide useful learning signal. If we mask everything, the optimal reconstruction are the marginals of each token in the data distribution, which is easy to learn, and again is not useful. These extreme masking rates lead to poor high-variance gradients: we want to learn how to clip them via a simple and effective new class of schedules.

5.2 CLIPPED SCHEDULES FOR LOW-VARIANCE GRADIENTS

We propose a class of “clipped” noise schedules that sample mask rates $1 - \alpha_t \sim \mathcal{U}[\beta, \omega]$ for $0 \leq \beta, \omega \leq 1$. We argue that from the perspective of deriving Monte Carlo gradient estimates, these schedules are equivalent to a continuous schedule where the mask probability is approximately 0 before the specified range such that $1 - \alpha_{<\beta} \approx \epsilon$ and approximately 1 after the specified range $1 - \alpha_{>\omega} \approx 1 - \epsilon$. Consequently, α'_t is linear within the range: $\alpha'_t \approx 1/(\beta - \omega)$.

5.3 DATA-DRIVEN CLIPPED SCHEDULES ACROSS BLOCK SIZES

As the optimal mask rates may differ depending on the block size L' , we adaptively learn the schedule during training. While Kingma et al. (2021) perform variance minimization by isolating a variance term using their squared diffusion loss, this strategy is not directly applicable to our variance estimator in Equation 10 since we seek to reduce variance across random batches in addition to random t_b .

Instead, we optimize parameters β, ω to directly minimize training variance. To limit the computational burden of the optimization, we use the variance of the estimator of the diffusion ELBO as a proxy for the gradient estimator to optimize β, ω : $\min_{\beta, \omega} \text{Var}_{\mathbf{X}, t} [\mathcal{L}(\mathbf{X}; \theta, \beta, \omega)]$. We provide experimental details on the optimization procedure in Sec. 6.

In Table 2, we show that variance of the diffusion NELBO is correlated with test perplexity. Under a range of “clipped” noise rate distributions, we find that there exists a unique distribution for each block size $L' \in \{4, 16, 128\}$ that minimizes both the variance of the NELBO and the test perplexity.

Table 2: Perplexities (PPLs; \downarrow) and variances of the NELBO $\text{Var}_{\mathbf{X}, t} [\mathcal{L}_{\text{BD}}(\mathbf{X}; \theta)]$ (Var. NELBO; \downarrow). Models are trained on LM1B using uniform noise for 65B tokens, then finetuned for 10B tokens.

L'	$\mathcal{U}[0, .5]$		$\mathcal{U}[,3, .8]$		$\mathcal{U}[,5, 1]$		$\mathcal{U}[0, 1]$	
	PPL	Var. NELBO	PPL	Var. NELBO	PPL	Var. NELBO	PPL	Var. NELBO
128	31.72	1.03	31.78	1.35	31.92	1.83	31.78	3.80
16	31.27	7.90	31.19	3.62	31.29	3.63	31.33	7.39
4	29.23	32.68	29.37	10.39	29.16	8.28	29.23	23.65

6 EXPERIMENTS

We evaluate BD3-LMs across standard language modeling benchmarks and demonstrate their ability to generate arbitrary-length sequences unconditionally. We train a base BD3-LM using the maximum block size $L' = L$ for 850K gradient steps and finetune under varying L' for 150K gradient steps on the One Billion Words dataset (LM1B; Chelba et al. (2014)) and OpenWebText (OWT; Gokaslan et al. (2019)). Details on training and inference are provided in Suppl C.

To reduce the variance of training on the diffusion NELBO, we adaptively learn the range of masking rates by optimizing parameters β, ω as described in Section 5.3. In practice, we do so using a grid search during every validation step (after $\sim 5\text{K}$ gradient up-

Table 3: Test perplexities (PPL; \downarrow) of models trained for 65B tokens on LM1B. Best diffusion value is bolded.

	PPL (\downarrow)
Autoregressive	
Transformer-X Base (Dai et al., 2019)	23.5
Transformer (Sahoo et al., 2024a)	22.83
Diffusion	
D3PM (absorb) (Austin et al., 2021)	≤ 82.34
SEDD (Lou et al., 2024)	≤ 32.68
MDLM (Sahoo et al., 2024a)	≤ 31.78
Block diffusion (Ours)	
BD3-LMs $L' = 16$	≤ 30.60
$L' = 8$	≤ 29.83
$L' = 4$	\leq 28.23

dates) to identify β, ω : $\min_{\beta, \omega} \text{Var}_{\mathbf{X}, t} [\mathcal{L}(\mathbf{X}; \theta, \beta, \omega)]$. During evaluation, we report likelihood under uniformly sampled mask rates (8) as in Austin et al. (2021); Sahoo et al. (2024a).

6.1 LIKELIHOOD EVALUATION

On LM1B, BD3-LMs outperform all prior diffusion methods in Table 3. Compared to MDLM (Sahoo et al., 2024a), BD3-LMs achieve up to 12% improvement in perplexity. We observe a similar trend on OpenWebText in Table 4.

We also evaluate the ability of BD3-LMs to generalize to unseen datasets in a zero-shot setting, following the benchmark from Radford et al. (2019). We evaluate the likelihood of models trained with OWT on datasets Penn Tree Bank (PTB; (Marcus et al., 1993)), Wikitext (Merity et al., 2016), LM1B, Lambada (Paperno et al., 2016), AG News (Zhang et al., 2015), and Scientific Papers (Pubmed and Arxiv subsets; (Cohan et al., 2018)). In Table 5, BD3-LM achieves the best zero-shot perplexity on Pubmed, surpassing AR, and the best perplexity among diffusion models on Wikitext, LM1B, and AG News.

Table 4: Test perplexities (PPL; \downarrow) on OWT for models trained for 524B tokens. Best diffusion value is bolded.

	PPL (\downarrow)
AR (Sahoo et al., 2024a)	17.54
SEDD (Lou et al., 2024)	≤ 24.10
MDLM (Sahoo et al., 2024a)	≤ 22.98
BD3-LMs $L' = 16$	≤ 22.27
$L' = 8$	≤ 21.68
$L' = 4$	$\leq \mathbf{20.73}$

Table 5: Zero-shot validation perplexities (\downarrow) of models trained for 524B tokens on OWT. All perplexities for diffusion models are upper bounds.

	PTB	Wikitext	LM1B	Lambada	AG News	Pubmed	Arxiv
AR	81.07	25.32	51.14	52.13	52.11	48.59	41.22
SEDD	96.33	35.98	68.14	48.93	67.82	45.39	40.03
MDLM	90.96	33.22	64.94	48.29	62.78	43.13	37.89
BD3-LM $L' = 4$	96.81	31.31	60.88	50.03	61.67	42.52	39.20

6.2 SAMPLE QUALITY AND VARIABLE-LENGTH SEQUENCE GENERATION

We assess the capacity of BD3-LMs to generate high-quality, variable-length samples. One key drawback of many existing diffusion language models (e.g., Austin et al. (2021); Lou et al. (2024)) is that they cannot generate full-length sequences that are longer than the length of the output context chosen at training time. The OWT dataset is useful for examining this limitation, as it contains many documents that are longer than the training context length of 1024 tokens.

We record generation length statistics of 500 variable-length samples in Table 6. For AR and BD3-LM, we continue sampling tokens until an end-of-sequence token [EOS] is generated or sample quality significantly degrades (as measured by likelihood and sample entropy). BD3-LMs generate sequences of up to 9K tokens, whereas SEDD (Lou et al., 2024) is restricted to the training context size.

Table 6: Generation length statistics from sampling 500 documents from models trained on OWT.

	Median # tokens	Max # tokens
OWT train set	717	131K
AR	6719	131K
SEDD	1021	1024
BD3-LM $L' = 16$	805	9080

We also examine the sample quality of BD3-LMs through quantitative and qualitative analyses. In Table 7, we generate sequences of lengths $L = 1024, 2048$ and measure their generative perplexity under GPT2-Large. To sample $L = 2048$ tokens from MDLM, we use their block-wise decoding technique (which does not feature block diffusion training as in BD3-LMs). We also compare to SSD-LM (Han et al., 2022), an alternative block diffusion formulation—also referred to as semi-autoregression. Unlike our discrete diffusion framework, SSD-LM uses Gaussian diffusion and does not support likelihood estimation.

BD3-LMs achieve the best generative perplexities compared to all previous diffusion methods. Compared to SSD-LM, our discrete approach yields samples with improved generative perplexity using an order of magnitude fewer generation steps.

Table 7: Generative perplexity (Gen. PPL; \downarrow) and number of function evaluations (NFEs; \downarrow) of 300 samples of lengths $L = 1024, 2048$. All models are trained on OWT (AR, SEDD, MDLM are trained on 524B tokens, SSD-LM is pre-trained on 122B tokens). Best diffusion value is bolded. We provide further inference details on reporting generative perplexity and NFEs in Supp. C.2.

Model	$L = 1024$		$L = 2048$	
	Gen. PPL	NFEs	Gen. PPL	NFEs
AR	14.1	1K	12.8	2K
Diffusion				
SEDD	52.6	1K	–	–
MDLM	47.0	1K	41.7	2K
Block Diffusion				
SSD-LM $L' = 25$	37.2	40K	35.3	80K
	281.3	1K	281.9	2K
BD3-LMs $L' = 16$	37.8	1K	36.8	2K
$L' = 8$	36.4	1K	35.7	2K
$L' = 4$	30.5	1K	29.9	2K

We also qualitatively examine samples taken from BD3-LM and baselines (AR, MDLM) trained on the OWT dataset; we report samples in Suppl. D. We observe that BD3-LM samples have higher coherence than those from MDLM and approach the quality of AR samples.

6.3 ABLATIONS

We assess the impact of the design choices in our proposed block diffusion recipes, namely 1) selection of the noise schedule and 2) the efficiency improvement of the proposed training algorithm relative to a naive implementation.

SELECTING NOISE SCHEDULES TO REDUCE TRAINING VARIANCE

Compared to the linear schedule used in Lou et al. (2024); Sahoo et al. (2024a), training under “clipped” noise schedules is the most effective for reducing the training variance which correlates with test perplexity. In Table 8, the ideal “clipped” masking rates, which are optimized during training, are specific to the block size and further motivate our optimization.

Relative to other standard noise schedules (Chang et al., 2022), “clipped” masking achieves the best training variance and test perplexity. Since heavier masking is effective for the smaller block size $L' = 4$, we compare with logarithmic and square root schedules that also encourage heavy masking. As lighter masking is optimal for block size $L' = 16$, we compare with square and cosine schedules.

EFFICIENCY OF TRAINING ALGORITHM

In the training algorithm presented in Section 3.2, we compute $\mathbf{x}_{\text{logit}}$ using two options. We may perform two forward passes through the network (precomputing keys and values for the full sequence \mathbf{x} , then computing denoised predictions), or combine these passes by concatenating the two inputs into the same attention kernel.

We find that performing this operation in a single forward pass is often more efficient as we reduce memory bandwidth bottlenecks by leveraging efficient, pre-existing FlashAttention kernels Dao et al. (2022). In-

Table 8: Effect of the noise schedule on training variance and test perplexity. Models are finetuned for 3B tokens on LM1B and evaluated under the linear schedule. For clipped schedules, we compare the optimized clipping rates for $L' = 4, 16$.

Noise schedule	PPL	Var NELBO
$L' = 4$		
Clipped		
$t \sim \mathcal{U}[0.45, 0.95]$	29.21	6.24
$t \sim \mathcal{U}[0.3, 0.8]$	29.38	10.33
Linear $t \sim \mathcal{U}[0, 1]$	30.18	23.45
Logarithmic	30.36	23.53
Square root	31.41	26.43
$L' = 16$		
Clipped		
$t \sim \mathcal{U}[0.45, 0.95]$	31.42	3.60
$t \sim \mathcal{U}[0.3, 0.8]$	31.12	3.58
Linear $t \sim \mathcal{U}[0, 1]$	31.72	7.62
Square	31.43	13.03
Cosine	31.41	13.00

stead of paying the cost of 2 passes through the network, we only pay the cost of a more expensive attention operation. Empirically we see that this approach has $>2\times$ speed-up relative to performing two forward passes.

7 DISCUSSION AND PRIOR WORK

Comparison to D3PM Block diffusion builds off D3PM (Austin et al., 2021) and applies it to each auto-regressive conditional. We improve over D3PM in three ways: (1) we extend D3PM beyond fixed sequence lengths; (2) we study the perplexity gap of D3PM and AR models, identify gradient variance as a contributor, and design variance-minimizing schedules; (3) we improve over the perplexity of D3PM models. While (1) involves modifying D3PM to apply it over autoregressive conditionals, (2) is applicable to vanilla D3PM.

Comparison to MDLM BD3-LMs further make use of the perplexity-enhancing improvements in MDLM (Sahoo et al., 2024a). We also build upon MDLM: (1) while Sahoo et al. (2024a) points out that their NELBO is invariant to the noise schedule, we show that the noise schedule has a significant effect on gradient variance; (2) we push the state-of-the-art in perplexity beyond MDLM. Note that our perplexity improvements stem not only from block-wise diffusion, but also from optimized schedules, and could enhance standard MDLM models.

Comparison to Autoregressive-Diffusion Models Han et al. (2022) introduced an alternative block diffusion formulation that performs Gaussian diffusion over word embeddings, as in Li et al. (2022). Our approach instead applies discrete noise as in Austin et al. (2021), and features notable improvements: (1) tractable likelihood estimates enabling principled evaluation; (2) faster generation, as our number of model calls is bounded by the number of generated tokens, while SSD-LM performs orders of magnitude more calls; (3) significantly improved performance, measured by perplexity relative to existing models as well as generative perplexity relative to samples from both SSD-LM. AR-Diffusion (Wu et al., 2023) is a variant of SSD-LM that uses a noise schedule which encourages left-to-right generation. However, they sacrifice parallelism by assigning unique, per-token timesteps.

Comparison to Jacobi Decoding Jacobi decoding (Santilli et al., 2023) is an AR inference technique that iteratively refines a random sequence and supports parallel generation of token blocks. However, Santilli et al. (2023) preserve causal masking from AR whereas and use uniform noise, whereas BD3-LMs may leverage more context by attending to tokens within a block and use masking which has been shown to be superior in language modeling (Austin et al., 2021; Lou et al., 2024). Consistency LLMs (Kou et al., 2024) extend Jacobi decoding to include a fine-tuning objective. In contrast, BD3-LMs may leverage clean conditional context $\mathbf{x}^{<b}$ to enhance predictions.

Limitations Training BD3-LMs requires a custom procedure that is more expensive than regular diffusion training. We propose a vectorized implementation that keeps training speed within $<2\times$ of diffusion training speed; in our experiments, we also pre-train with a standard diffusion loss to further reduce the speed gap. Additionally, BD3-LMs generate blocks sequentially, hence may face the same speed and controllability constraints as AR, especially when blocks are small. The optimal block size for control is task specific, and optimal blocks for speed depend on the parallelization capabilities of inferencing hardware (e.g., FLOPS vs. memory throughput) and serving batch size.

8 CONCLUSION

This work explores block diffusion and is motivated by two problems with existing discrete diffusion: the need to generate arbitrary-length sequences and the perplexity gap to autoregressive models. We introduce BD3-LMs, which represent a block-wise extension of the D3PM framework (Austin et al., 2021), and leverage a specialized training algorithm and custom noise schedules that further improve performance. We observe that in addition to being able to generate long-form documents, these models also improve perplexity, setting a new state-of-the-art among discrete diffusion models. BD3-LMs are subject to inherent limitations of generative models, including hallucinations (Achiam et al., 2023), copyright infringement (Gokaslan et al., 2024), limited controllability (Schiff et al., 2024; Wang et al., 2023) and harmful outputs (Bai et al., 2022), which require further research.

ACKNOWLEDGMENTS AND DISCLOSURE OF FUNDING

This work was partially funded by the National Science Foundation under awards DGE-1922551, CAREER awards 2046760 and 2145577, and by the National Institute of Health under award MIRA R35GM151243. Marianne Arriola is supported by a NSF Graduate Research Fellowship under award DGE-2139899 and a Hopper-Dean/Bowers CIS Deans Excellence Fellowship.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.
- Pavel Avdeyev, Chenlai Shi, Yuhao Tan, Kseniia Dudnyk, and Jian Zhou. Dirichlet diffusion score model for biological sequence generation. In *International Conference on Machine Learning*, pp. 1276–1301. PMLR, 2023.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.
- Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11315–11325, 2022.
- Ciprian Chelba, Tomas Mikolov, Mike Schuster, Qi Ge, Thorsten Brants, Phillipp Koehn, and Tony Robinson. One billion word benchmark for measuring progress in statistical language modeling, 2014.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. A discourse-aware attention model for abstractive summarization of long documents. *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, 2018. doi: 10.18653/v1/n18-2097. URL <http://dx.doi.org/10.18653/v1/n18-2097>.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *arXiv preprint arXiv:1901.02860*, 2019.
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing Systems*, 35:16344–16359, 2022.
- Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis, 2021. URL <https://arxiv.org/abs/2105.05233>.
- Shrey Goel, Vishrut Thoutam, Edgar Mariano Marroquin, Aaron Gokaslan, Arash Firouzbakht, Sophia Vincoff, Volodymyr Kuleshov, Huong T Kratochvil, and Pranam Chatterjee. Memdlm: De novo membrane protein design with masked discrete diffusion protein language models. *arXiv preprint arXiv:2410.16735*, 2024.
- Aaron Gokaslan, Vanya Cohen, Ellie Pavlick, and Stefanie Tellex. Openwebtext corpus. <http://Skyllion007.github.io/OpenWebTextCorpus>, 2019.
- Aaron Gokaslan, A Feder Cooper, Jasmine Collins, Landan Seguin, Austin Jacobson, Mihir Patel, Jonathan Frankle, Cory Stephenson, and Volodymyr Kuleshov. Commoncanvas: Open diffusion models trained on creative-commons images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8250–8260, 2024.

- Ishaan Gulrajani and Tatsunori B Hashimoto. Likelihood-based diffusion language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama. Photorealistic video generation with diffusion models, 2023. URL <https://arxiv.org/abs/2312.06662>.
- Xiaochuang Han, Sachin Kumar, and Yulia Tsvetkov. Ssd-lm: Semi-autoregressive simplex-based diffusion language model for text generation and modular control. *arXiv preprint arXiv:2210.17432*, 2022.
- Xiaochuang Han, Sachin Kumar, Yulia Tsvetkov, and Marjan Ghazvininejad. David helps goliath: Inference-time collaboration between small specialized and large general diffusion lms. *arXiv preprint arXiv:2305.14771*, 2023.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv:2204.03458*, 2022.
- Daniel Israel, Aditya Grover, and Guy Van den Broeck. Enabling autoregressive models to fill in masked tokens. *arXiv preprint arXiv:2502.06901*, 2025.
- Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.
- Siqi Kou, Lanxiang Hu, Zhezhi He, Zhijie Deng, and Hao Zhang. CLLMs: Consistency large language models. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=8uzBOVmh8H>.
- Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *Advances in Neural Information Processing Systems*, 35: 4328–4343, 2022.
- Xiner Li, Yulai Zhao, Chenyu Wang, Gabriele Scalia, Gokcen Eraslan, Surag Nair, Tommaso Biancalani, Shuiwang Ji, Aviv Regev, Sergey Levine, et al. Derivative-free guidance in continuous and discrete diffusion models with soft value-based decoding. *arXiv preprint arXiv:2408.08252*, 2024.
- Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion modeling by estimating the ratios of the data distribution. In *Forty-first International Conference on Machine Learning*, 2024. URL <https://openreview.net/forum?id=CNicRIVIPA>.
- Mitch Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a large annotated corpus of english: The penn treebank. *Computational linguistics*, 19(2):313–330, 1993.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models, 2016.
- Hunter Nisonoff, Junhao Xiong, Stephan Allenspach, and Jennifer Listgarten. Unlocking guidance for discrete state-space diffusion and flow models. *arXiv preprint arXiv:2406.01572*, 2024.
- Jingyang Ou, Shen Nie, Kaiwen Xue, Fengqi Zhu, Jiacheng Sun, Zhenguo Li, and Chongxuan Li. Your absorbing discrete diffusion secretly models the conditional distributions of clean data. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=sMyXP8Tanm>.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Ngoc Quan Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernandez. The LAMBADA dataset: Word prediction requiring a broad discourse context. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1525–1534, Berlin, Germany, August 2016. Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/P16-1144>.

- William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- Subham Sekhar Sahoo, Marianne Arriola, Aaron Gokaslan, Edgar Mariano Marroquin, Alexander M Rush, Yair Schiff, Justin T Chiu, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a. URL <https://openreview.net/forum?id=L4uaAR4ArM>.
- Subham Sekhar Sahoo, Aaron Gokaslan, Christopher De Sa, and Volodymyr Kuleshov. Diffusion models with learned adaptive noise. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024b. URL <https://openreview.net/forum?id=loMa99A4p8>.
- Andrea Santilli, Silvio Severino, Emilian Postolache, Valentino Maiorca, Michele Mancusi, Riccardo Marin, and Emanuele Rodola. Accelerating transformer inference for translation via parallel decoding. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12336–12355, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.689. URL <https://aclanthology.org/2023.acl-long.689>.
- Yair Schiff, Subham Sekhar Sahoo, Hao Phung, Guanghan Wang, Sam Boshar, Hugo Dalla-torre, Bernardo P de Almeida, Alexander Rush, Thomas Pierrot, and Volodymyr Kuleshov. Simple guidance mechanisms for discrete diffusion models. *arXiv preprint arXiv:2412.10193*, 2024.
- Jiaxin Shi, Kehang Han, Zhe Wang, Arnaud Doucet, and Michalis Titsias. Simplified and generalized masked diffusion for discrete data. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=xcqSOfHt4g>.
- Phillip Si, Allan Bishop, and Volodymyr Kuleshov. Autoregressive quantile flows for predictive uncertainty estimation. In *International Conference on Learning Representations*, 2022.
- Phillip Si, Zeyi Chen, Subham Sekhar Sahoo, Yair Schiff, and Volodymyr Kuleshov. Semi-autoregressive energy flows: exploring likelihood-free training of normalizing flows. In *International Conference on Machine Learning*, pp. 31732–31753. PMLR, 2023.
- Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pp. 2256–2265. PMLR, 2015.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019.
- Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021.
- Haoran Sun, Lijun Yu, Bo Dai, Dale Schuurmans, and Hanjun Dai. Score-based continuous-time discrete diffusion models. *arXiv preprint arXiv:2211.16750*, 2022.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Yingheng Wang, Yair Schiff, Aaron Gokaslan, Weishen Pan, Fei Wang, Christopher De Sa, and Volodymyr Kuleshov. InfoDiffusion: Representation learning using information maximizing diffusion models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 36336–36354. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/wang23ah.html>.

Tong Wu, Zhihao Fan, Xiao Liu, Hai-Tao Zheng, Yeyun Gong, yelong shen, Jian Jiao, Juntao Li, zhongyu wei, Jian Guo, Nan Duan, and Weizhu Chen. AR-diffusion: Auto-regressive diffusion model for text generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=0EG6qUQ4xE>.

Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In *NIPS*, 2015.

Kaiwen Zheng, Yongxin Chen, Hanzi Mao, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Masked diffusion models are secretly time-agnostic masked models and exploit inaccurate categorical sampling. *arXiv preprint arXiv:2409.02908*, 2024.

CONTENTS

1	Introduction	1
2	Background: Language Modeling Paradigms	2
2.1	Autoregressive Models	3
2.2	Discrete Denoising Diffusion Probabilistic Models	3
3	Block Diffusion Language Modeling	3
3.1	Block Diffusion Distributions and Model Architectures	3
3.2	Efficient Training and Sampling Algorithms	4
4	Understanding Likelihood Gaps Between Diffusion & AR Models	5
4.1	Masked BD3-LMs	5
4.2	Case Study: Single Token Generation	5
4.3	Diffusion Gap from High Variance Training	6
5	Low-Variance Noise Schedules for BD3-LMs	6
5.1	Intuition: Avoid Extreme Mask Rates	6
5.2	Clipped Schedules for Low-Variance Gradients	7
5.3	Data-Driven Clipped Schedules Across Block Sizes	7
6	Experiments	7
6.1	Likelihood Evaluation	8
6.2	Sample Quality and Variable-Length Sequence Generation	8
6.3	Ablations	9
7	Discussion and Prior Work	10
8	Conclusion	10
A	Block Diffusion NELBO	16
B	Masked BD3-LMs	16
B.1	Forward Process	17
B.2	Reverse Process	17
B.3	Simplified NELBO for Masked Diffusion Processes	17
B.4	Recovering the NLL from the NELBO for Single Token Generation	18
B.5	Tightness of the NELBO	19
B.6	Specialized Attention Masks	19
C	Experimental Details	20
C.1	Training	20

C.2 Inference	21
D Samples	22

A BLOCK DIFFUSION NELBO

Below, we provide the Negative ELBO (NELBO) for the block diffusion parameterization. Recall that the sequence $\mathbf{x}^{1:L} = [\mathbf{x}^1, \dots, \mathbf{x}^L]$ is factorized over B blocks, which we refer to as \mathbf{x} for simplicity, drawn from the data distribution $q(\mathbf{x})$. Specifically, we will factorize the likelihood over B blocks of length L' , then perform diffusion in each block over T discretization steps. Let $D_{\text{KL}}[\cdot]$ to denote the Kullback–Leibler divergence. We derive the NELBO as follows:

$$\begin{aligned}
-\log p_\theta(\mathbf{x}) &= -\sum_{b=1}^B \log p_\theta(\mathbf{x}^b | \mathbf{x}^{<b}) \\
&= -\sum_{b=1}^B \log \mathbb{E}_q \frac{p_\theta(\mathbf{x}_{0:T}^b | \mathbf{x}^{<b})}{q(\mathbf{x}_{1:T}^b | \mathbf{x}^b)} \\
&= -\sum_{b=1}^B \log \mathbb{E}_q \frac{p_\theta(\mathbf{x}_T^b | \mathbf{x}^{<b}) \prod_{i=1}^T p_\theta(\mathbf{x}_s^b | \mathbf{x}_t^b, \mathbf{x}^{<b})}{\prod_{t=1}^T q(\mathbf{x}_t^b | \mathbf{x}_s^b)} \\
&\leq \sum_{b=1}^B \left[\underbrace{-\mathbb{E}_q \log p_\theta(\mathbf{x}^b | \mathbf{x}_1^b, \mathbf{x}^{<b})}_{\mathcal{L}_{\text{recons}}} \right. \\
&\quad \left. + \underbrace{\mathbb{E}_{t \in \{\frac{2}{T}, \dots, \frac{T-1}{T}, 1\}} \mathbb{E}_q \text{TD}_{\text{KL}}(q(\mathbf{x}_s^b | \mathbf{x}_t^b, \mathbf{x}^b) \parallel p_\theta(\mathbf{x}_s^b | \mathbf{x}_t^b, \mathbf{x}^{<b}))}_{\mathcal{L}_{\text{diffusion}}} \right. \\
&\quad \left. + \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t=1}^b | \mathbf{x}^b) \parallel p_\theta(\mathbf{x}_{t=1}^b))}_{\mathcal{L}_{\text{prior}}} \right] \tag{11}
\end{aligned}$$

For simplicity, we will denote $s = (t - 1)/T$.

B MASKED BD3-LMs

We explore a specific class of block diffusion models that builds upon the masked diffusion language modeling framework. In particular, we focus on masking diffusion processes introduced by [Austin et al. \(2021\)](#) and derive a simplified NELBO under this framework as proposed by ([Sahoo et al., 2024a](#); [Shi et al., 2024](#); [Ou et al., 2025](#)).

First, we define the diffusion matrix Q_t for states $i \in 1, \dots, V$. Consider the noise schedule function $\alpha_t \in [0, 1]$, which is a strictly decreasing function in t satisfying $\alpha_0 = 1$ and $\alpha_1 = 0$. Denote the mask index as $m = V$. The diffusion matrix is defined by [Austin et al. \(2021\)](#) as:

$$[Q_t]_{ij} = \begin{cases} 1 & \text{if } i = j = m \\ \alpha_t & \text{if } i = j \neq m \\ 1 - \alpha_t & \text{if } j = m, i \neq m \end{cases} \tag{12}$$

The diffusion matrix for the forward marginal $Q_{t|s}$ is:

$$[Q_{t|s}]_{ij} = \begin{cases} 1 & \text{if } i = j = m \\ \alpha_{t|s} & \text{if } i = j \neq m \\ 1 - \alpha_{t|s} & \text{if } j = m, i \neq m \end{cases} \tag{13}$$

where $\alpha_{t|s} = \alpha_t / \alpha_s$.

B.1 FORWARD PROCESS

Under the D3PM framework (Austin et al., 2021), the forward noise process applied independently for each token $\ell \in 1, \dots, L$ is defined using diffusion matrices $Q_t \in \mathbb{R}^{V \times V}$ as

$$q(\mathbf{x}_t^\ell | \mathbf{x}^\ell) = \text{Cat}(\mathbf{x}_t^\ell; \bar{Q}_t \mathbf{x}^\ell), \quad \text{with} \quad \bar{Q}_t = Q_{t(1)} Q_{t(2)} \dots Q_{t(i)} \quad (14)$$

B.2 REVERSE PROCESS

Let $Q_{t|s}$ denote the diffusion matrix for the forward marginal. We obtain the reverse posterior $q(\mathbf{x}_s^\ell | \mathbf{x}_t^\ell, \mathbf{x}^\ell)$ using the diffusion matrices:

$$q(\mathbf{x}_s^\ell | \mathbf{x}_t^\ell, \mathbf{x}^\ell) = \frac{q(\mathbf{x}_t^\ell | \mathbf{x}_s^\ell, \mathbf{x}^\ell) q(\mathbf{x}_s^\ell | \mathbf{x}^\ell)}{q(\mathbf{x}_t^\ell | \mathbf{x}^\ell)} = \text{Cat}\left(\mathbf{x}_s^\ell; \frac{Q_{t|s} \mathbf{x}_t^\ell \odot Q_s^\top \mathbf{x}^\ell}{(\mathbf{x}_t^\ell)^\top Q_t^\top \mathbf{x}^\ell}\right) \quad (15)$$

where \odot denotes the Hadamard product between two vectors.

B.3 SIMPLIFIED NELBO FOR MASKED DIFFUSION PROCESSES

Following Sahoo et al. (2024a); Shi et al. (2024); Ou et al. (2025), we simplify the NELBO in the case of masked diffusion processes. Below, we provide the outline of the NELBO derivation; see the full derivation in Sahoo et al. (2024a); Shi et al. (2024); Ou et al. (2025).

We will first focus on simplifying the diffusion loss term $\mathcal{L}_{\text{diffusion}}$ in Eq. 11. We employ the SUBS-parameterization proposed in Sahoo et al. (2024b) which simplifies the denoising model p_θ for masked diffusion. In particular, we enforce the following constraints on the design of p_θ by leveraging the fact that there only exists two possible states in the diffusion process $\mathbf{x}_t^\ell \in \{\mathbf{x}^\ell, \mathbf{m}\} \forall \ell \in 1, \dots, L$.

1. **Zero Masking Probabilities.** We set $p_\theta(\mathbf{x}^\ell = \mathbf{m} | \mathbf{x}_t^\ell) = 0$ (as the clean sequence \mathbf{x} doesn't contain masks).
2. **Carry-Over Unmasking.** The true posterior for the case where $\mathbf{x}_t^\ell \neq \mathbf{m}$ is $q(\mathbf{x}_s^\ell = \mathbf{x}_t^\ell | \mathbf{x}_t^\ell \neq \mathbf{m}) = 1$ (if a token is unmasked in the reverse process, it is never remasked). Thus, we simplify the denoising model by setting $p_\theta(\mathbf{x}_s^\ell = \mathbf{x}_t^\ell | \mathbf{x}_t^\ell \neq \mathbf{m}) = 1$.

As a result, we will only approximate the posterior $p_\theta(\mathbf{x}_s^\ell = \mathbf{x}^\ell | \mathbf{x}_t^\ell = \mathbf{m})$. Let $\mathbf{x}^{b,\ell}$ denote a token in the ℓ -th position in block $b \in 1, \dots, B$. The diffusion loss term becomes:

$$\begin{aligned} \mathcal{L}_{\text{diffusion}} &= \sum_{b=1}^B \mathbb{E}_t \mathbb{E}_q T \left[\text{D}_{\text{KL}} \left[q(\mathbf{x}_s^b | \mathbf{x}_t^b) \| p_\theta(\mathbf{x}_s^b | \mathbf{x}_t^b, \mathbf{x}^{<b}) \right] \right] \\ &= \sum_{b=1}^B \mathbb{E}_t \mathbb{E}_q T \left[\sum_{\ell=1}^{L'} \text{D}_{\text{KL}} \left[q(\mathbf{x}_s^{b,\ell} | \mathbf{x}_t^{b,\ell}, \mathbf{x}^\ell) \| p_\theta(\mathbf{x}_s^{b,\ell} | \mathbf{x}_t^{b,\ell}, \mathbf{x}^{<b}) \right] \right] \\ \text{D}_{\text{KL}} &\text{ is simply the discrete-time diffusion loss for the block } b; \text{ hence, from Sahoo et al. (2024a) (Suppl. B.1), we get:} \\ &= \sum_{b=1}^B \mathbb{E}_t \mathbb{E}_q T \left[\sum_{\ell=1}^{L'} \frac{\alpha_s - \alpha_t}{1 - \alpha_t} \log p_\theta(\mathbf{x}^{b,\ell} | \mathbf{x}_t^{b,\ell}, \mathbf{x}^{<b}) \right] \\ &= \sum_{b=1}^B \mathbb{E}_t \mathbb{E}_q T \left[\frac{\alpha_s - \alpha_t}{1 - \alpha_t} \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b}) \right] \end{aligned} \quad (16)$$

Lastly, we obtain a tighter approximation of the likelihood by taking the diffusion steps $T \rightarrow \infty$ (Sahoo et al., 2024a), for which $T(\alpha_s - \alpha_t) = \alpha'_t$:

$$\mathcal{L}_{\text{diffusion}} = \sum_{b=1}^B \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \left[\frac{\alpha'_t}{1 - \alpha_t} \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b}) \right] \quad (17)$$

For the continuous time case, [Sahoo et al. \(2024a\)](#) (Suppl. A.2.4) show the reconstruction loss reduces to 0 as $\mathbf{x}_1^b \sim \lim_{T \rightarrow \infty} \text{Cat}(\cdot; q(\mathbf{x}_1^b | \mathbf{x}^b)) \implies \mathbf{x}_1^b \sim \text{Cat}(\cdot; \mathbf{x}^b)$. Using this, we obtain:

$$\begin{aligned}\mathcal{L}_{\text{recons}} &= -\log p_\theta(\mathbf{x}^b | \mathbf{x}_1^b, \mathbf{x}^{<b}) \\ &= -\log p_\theta(\mathbf{x}^b | \mathbf{x}_1^b = \mathbf{x}^b, \mathbf{x}^{<b}) \\ &= 0\end{aligned}\tag{18}$$

The prior loss $\mathcal{L}_{\text{prior}} = \text{D}_{KL}(q(\mathbf{x}_T^b | \mathbf{x}^b) \| p_\theta(\mathbf{x}_T^b))$ also reduces to 0 under the SUBS-parameterization as $q(\mathbf{x}_T^b | \mathbf{x}^b) = \text{Cat}(\cdot; \mathbf{m})$ and $p_\theta(\mathbf{x}_T^b) = \text{Cat}(\cdot; \mathbf{m})$ (see [Sahoo et al. \(2024a\)](#), Suppl. A.2.4).

Finally, we obtain a simple objective that is a weighted average of cross-entropy terms:

$$\mathcal{L}_{\text{BD}}(\mathbf{x}; \theta) = \sum_{b=1}^B \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \left[\frac{\alpha'_t}{1 - \alpha_t} \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b}) \right]\tag{19}$$

The above NELBO is invariant to the choice of noise schedule α_t , see [Sahoo et al. \(2024a\)](#) (Suppl. E.1.1).

B.4 RECOVERING THE NLL FROM THE NELBO FOR SINGLE TOKEN GENERATION

Consider the block diffusion NELBO for a block size of 1 where $L' = 1, B = L$. The block diffusion NELBO is equivalent to the AR NLL when modeling a single token:

$$\begin{aligned}-\log p(\mathbf{x}) &\leq \sum_{b=1}^L \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \left[\frac{\alpha'_t}{1 - \alpha_t} \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b}) \right] \\ &\because \alpha'_t = -1 \text{ and } \alpha_t = 1 - t, \\ &= - \sum_{b=1}^L \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \left[\frac{1}{t} \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b}) \right] \\ &= - \sum_{b=1}^L \mathbb{E}_{t \sim [0,1]} \frac{1}{t} \mathbb{E}_q [\log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b, \mathbf{x}^{<b})] \\ &\text{Expanding } \mathbb{E}_q[\cdot], \\ &= - \sum_{b=1}^L \mathbb{E}_{t \sim [0,1]} \frac{1}{t} \left[q(\mathbf{x}_t^b = \mathbf{m} | \mathbf{x}^b) \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b = \mathbf{m}, \mathbf{x}^{<b}) \right. \\ &\quad \left. + q(\mathbf{x}_t^b = \mathbf{x}^b | \mathbf{x}^b) \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b = \mathbf{x}^b, \mathbf{x}^{<b}) \right]\end{aligned}\tag{20}$$

Recall that our denoising model employs the SUBS-parameterization proposed in [Sahoo et al. \(2024b\)](#). The ‘‘carry-over unmasking’’ property ensures that $\log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b = \mathbf{x}^b, \mathbf{x}^{<b}) = 0$, as an unmasked token is simply copied over from the input of the denoising model to the output. Hence, (20) reduces to following:

$$\begin{aligned}-\log p_\theta(\mathbf{x}) &\leq - \sum_{b=1}^L \mathbb{E}_{t \sim [0,1]} \frac{1}{t} q(\mathbf{x}_t^b = \mathbf{m} | \mathbf{x}^b) \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b = \mathbf{m}, \mathbf{x}^{<b}) \\ &\because q(\mathbf{x}_t^b = \mathbf{m} | \mathbf{x}) = t, \text{ we get:} \\ &= - \sum_{b=1}^L \mathbb{E}_{t \sim [0,1]} \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b = \mathbf{m}, \mathbf{x}^{<b}) \\ &= - \sum_{b=1}^L \log p_\theta(\mathbf{x}^b | \mathbf{x}_t^b = \mathbf{m}, \mathbf{x}^{<b})\end{aligned}\tag{21}$$

For single-token generation ($L' = 1$) we recover the autoregressive NLL.

B.5 TIGHTNESS OF THE NELBO

For block sizes $L \geq K \geq 1$, we show that $-\log p(\mathbf{x}) \leq \mathcal{L}_K \leq \mathcal{L}_{K+1}$. Consider $K = 1$, where we recover the autoregressive NLL (see Suppl B.4):

$$\begin{aligned}\mathcal{L}_1 &= \sum_{b=1}^L \log \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \frac{\alpha'_t}{1 - \alpha_t} p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b}) \\ &= - \sum_{b=1}^L \log p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b = \mathbf{m}, \mathbf{x}^{<b})\end{aligned}\quad (22)$$

Consider the ELBO for block size $K = 2$:

$$\mathcal{L}_2 = \sum_{b=1}^{L/2} \log \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \frac{\alpha'_t}{1 - \alpha_t} p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b}) \quad (23)$$

We show that $\mathcal{L}_1 \leq \mathcal{L}_2$, and this holds for all $L \geq K \geq 1$ by induction. Let $\mathbf{x}^{b,i}$ correspond to the token in position $i \in [1, L']$ of block b . We derive the below inequality:

$$\begin{aligned}- \sum_{b=1}^L \log p_\theta(\mathbf{x}_t^b = \mathbf{m} \mid \mathbf{x}^{<b}) &= - \sum_{b=1}^{L/2} \log \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \frac{1}{1 - \alpha_t} p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b}) \\ &= - \sum_{b=1}^{L/2} \log \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \prod_{i=1}^2 \frac{1}{1 - \alpha_t} p_\theta(\mathbf{x}^{b,i} \mid \mathbf{x}_t^b, \mathbf{x}^{<b}) \\ &= - \sum_{b=1}^{L/2} \log \prod_{i=1}^2 \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \frac{1}{1 - \alpha_t} p_\theta(\mathbf{x}^{b,i} \mid \mathbf{x}_t^b, \mathbf{x}^{<b}) \\ &\leq - \sum_{b=1}^{L/2} \sum_{i=1}^2 \log \mathbb{E}_{t \sim [0,1]} \mathbb{E}_q \frac{1}{1 - \alpha_t} p_\theta(\mathbf{x}^{b,i} \mid \mathbf{x}_t^b, \mathbf{x}^{<b})\end{aligned}\quad (24)$$

B.6 SPECIALIZED ATTENTION MASKS

We aim to model conditional probabilities $p_\theta(\mathbf{x}^b \mid \mathbf{x}_t^b, \mathbf{x}^{<b})$ for blocks $b \in [1, B]$ simultaneously by designing an efficient training algorithm with our transformer backbone. However, modeling all B conditional terms requires processing both the noised sequence \mathbf{x}_t^b and the conditional context $\mathbf{x}^{<b}$.

Rather than calling the denoising network B times, we process both sequences simultaneously by concatenating them $\mathbf{x}_{full} \leftarrow \mathbf{x}_t \oplus \mathbf{x}$ as input to a transformer. We update this sequence \mathbf{x}_{full} of length $2L$ using a custom attention mask $\mathcal{M}(L, B) \in \{0, 1\}^{2L \times 2L}$ for efficient training.

This attention mask is comprised of $4 L \times L$ smaller attention masks:

$$\text{MASK}(L, B) = \begin{bmatrix} \mathcal{M}_{BD} & \mathcal{M}_{OBC} \\ \mathbf{0} & \mathcal{M}_{BC} \end{bmatrix}$$

where \mathcal{M}_{BD} and \mathcal{M}_{OBC} are used to update the representation of \mathbf{x}_t and \mathcal{M}_{BC} is used to update the representation of \mathbf{x} . We define these masks as follows:

- \mathcal{M}_{BD} (Block-diagonal mask): Self-attention mask within noised blocks \mathbf{x}_t^b

$$\mathcal{M}_{BD} = (m_{i,j})_{L \times L} = \begin{cases} 1 & \text{if } i, j \text{ are in the same block} \\ 0 & \text{otherwise} \end{cases}$$

- \mathcal{M}_{OBC} (Offset block-causal mask): Cross-attention mask for conditional context $\mathbf{x}^{<b}$

$$\mathcal{M}_{OBC} = (m_{i,j})_{L \times L} = \begin{cases} 1 & \text{if } i \text{ belongs in a block before } j \\ 0 & \text{otherwise} \end{cases}$$

- \mathcal{M}_{BC} (Block-causal mask): Attention mask for updating \mathbf{x}^b

$$\mathcal{M}_{BC} = (m_{i,j})_{L' \times L} = \begin{cases} 1 & \text{if } j \text{ is not in a block after } i \\ 0 & \text{otherwise} \end{cases}$$

We visualize an example attention mask for $L = 6$ and block size $L' = 2$ in Figure 3.

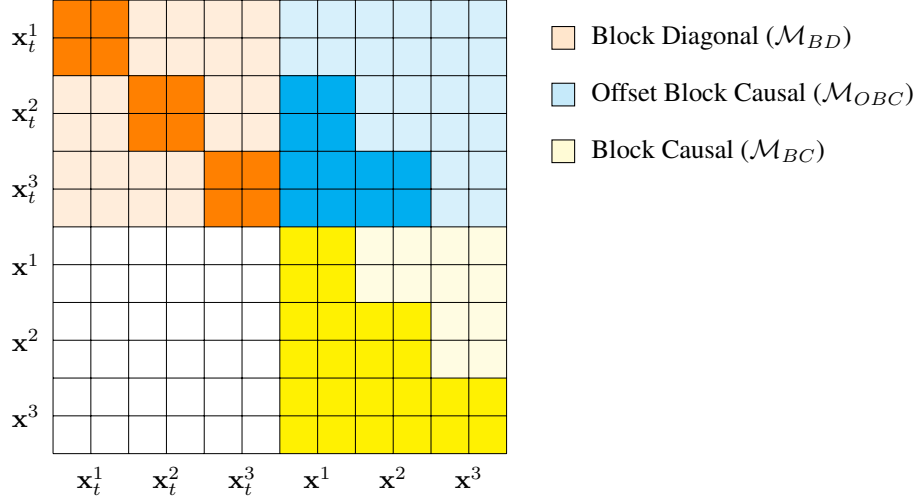


Figure 3: Example of Specialized Attention Mask

C EXPERIMENTAL DETAILS

C.1 TRAINING

We closely follow the same training and evaluation setup as used by [Sahoo et al. \(2024a\)](#). We conduct experiments on two datasets: The One Billion Word Dataset (LM1B; [Chelba et al. \(2014\)](#)) and Open-WebText (OWT; [Gokaslan et al. \(2019\)](#)). Models trained on LM1B use the `bert-base-uncased` tokenizer and a context length of 128. We report perplexities on the test split of LM1B. Models trained on OWT use the `GPT2` tokenizer [Radford et al. \(2019\)](#) and a context length of 1024. Since OWT does not have a validation split, we leave the last 100k documents for validation.

In preparing LM1B examples, [Sahoo et al. \(2024a\)](#) pad each example to fit in the context length. Since most examples consist of only a single sentence, block diffusion modeling for larger block sizes $L' > 4$ would not be useful for training. Instead, we concatenate and wrap sequences to a length of 128. As a result, we retrain our autoregressive baseline, SEDD, and MDLM on LM1B with wrapping.

Similarly for OWT, we do not pad or truncate sequences, but concatenate them and wrap them to a length of 1024 similar to LM1B. For unconditional generation experiments in Section 6.2, we wish to generate sequences longer than the context length seen during training. However, [Sahoo et al. \(2024a\)](#) inject beginning-of-sequence and end-of-sequence tokens ([BOS], [EOS] respectively) at the beginning and end of the training context. Thus, baselines from [Sahoo et al. \(2024a\)](#) will generate sequences that match the training context size. To examine model generations across varying decoding lengths in Section 6.2, we retrain our AR, SEDD, and MDLM baselines without injecting [BOS] and [EOS] tokens in the training examples. We also adopt this preprocessing convention for training all BD3-LMs on OWT.

The model architecture augments the diffusion transformer ([Peebles & Xie, 2023](#)) with rotary positional embeddings ([Su et al., 2021](#)). We parameterize our autoregressive baselines, SEDD, MDLM, and BD3-LMs with a transformer architecture from [Sahoo et al. \(2024a\)](#) that uses 12 layers, a hidden dimension of 768, and 128 attention heads. We do not include timestep conditioning as [Sahoo et al. \(2024a\)](#) show it does not affect performance. We use the AdamW optimizer with a batch size of 512 and constant learning rate warmup from 0 to $3e-4$ for 2.5k gradient updates.

We train a base BD3-LM using the maximum context length $L' = L$ for 850K gradient steps and fine-tune under varying L' using the noise schedule optimization for 150K gradient steps on the One Billion Words dataset (LM1B) and OpenWebText (OWT). This translates to 65B tokens and 73 epochs on LM1B, 524B tokens and 60 epochs on OWT. We use 3090, A5000, A6000, and A100 GPUs. Training BD3-LMs on LM1B takes 1.5 days on 4 A5000s and OWT takes 4.5 days on 8 A100s.

C.2 INFERENCE

We report generative perplexity under GPT2-Large from models trained on OWT using a context length of 1024 tokens. Since GPT2-Large uses a context size of 1024, we compute the generative perplexity for samples longer than 1024 tokens using a sliding window with a stride length of 512. We use the corrected Gumbel-based categorical sampling from [Zheng et al. \(2024\)](#) by sampling 64-bit Gumbel variables. Reducing the precision to 32-bit has been shown to significantly truncate the Gumbel variables, lowering the temperature and decreasing the sentence entropy.

Following SSD-LM ([Han et al., 2022](#)), we employ nucleus sampling for BD3-LMs and our baselines. For SSD-LM, we use their default hyperparameters $p = 0.95$ for block size $L' = 25$. For BD3-LMs, AR and MDLM, we use $p = 0.9$. For SEDD, we find that $p = 0.99$ works best.

In Table 7, BD3-LMs and MDLM use $T = 5k$ diffusion steps. BD3-LMs and MDLM use efficient sampling by caching the output of the denoising network as proposed by [Sahoo et al. \(2024a\)](#); [Ou et al. \(2025\)](#), which ensures that the number of generation steps does not exceed the sample length L . Put simply, once a token is sampled, it is never resampled tokens as a result of the simplified denoising model (Suppl. B.3). We use MDLM’s block-wise decoding algorithm for generating variable-length sequences, however these models are not trained with block diffusion.

SSD-LM (first row in Table 7) and SEDD use $T = 1k$ diffusion steps. Since block diffusion performs T diffusion steps for each block $b \in 1, \dots, B$, SSD-LM performs BT generation steps. Thus to fairly compare with SSD-LM, we also report generative perplexity for $T = 25$ diffusion steps so that the number of generation steps does not exceed the sequence length (second row in Table 7).

For arbitrary-length sequence generation using BD3-LMs and AR in Table 6, we continue to sample tokens until the following stopping criteria are met:

1. an [EOS] token is sampled
2. the average likelihood of the last 256-token chunk is below 0.1
3. the average entropy of the the last 256-token chunk is below 4

where criteria 2, 3 are necessary to prevent run-on samples from compounding errors (for example, a sequence of repeating tokens). We find that degenerate samples with low entropy result in significantly low perplexities under GPT2 and lower the reported generative perplexity. Thus, when a sample meets criterion 3, we regenerate the sample when reporting generative perplexity in Table 7.

D SAMPLES

<endofxt>'s architect, lawyer and San Giovanni concerto art critic Paolo Capacotti, gained attention from fellow gallery members and even invited him to present a retrospective, publishing issues and newspaper interviews.[10] On 6 September, Kissi and his assistants agreed to move to Angelo's Marcus Collection,[10] which included Giorgio Avolivo Arth and Moscoliso (later owned by the artist Belzina Massingolo) and Pan Giazoglio Romeam-Guessle. The businessman, Giovanni Paletti, an outstanding collector, owned the museum and the painting. The level of criminal activity around the museum has continued to increase, which is part of several attempts to counter centennial rumors including the possibility that museum staff and visitors are tortured and even exposed to del Cavello for the only full year of Francesco Belzina's life (1999).[4] On the evening of 22 October 2005 it was reported that earlier that evening, guards had come on duty and began flinging an electric field with umbrellas from the balcony. As the fire continued, some of the guards sparked an apparent spat from the window of the cathedral. They remained idly watched by a pile of trash left after a piano key by Pietro Jolla, who died on 21 October 2005.[10] Just before 3:00 to 3pm on Monday, 27 October 2005, strong winds brought the trash on to the residence that opened on 17 October. Some ruined books and statues were hurled in front from every direction of the window. Some claimed that a customer Jacques Monet had beaten the hand of photographer Franco Campetti and in some cases had stuck a broken candle in the doorway of the museum. Andr Romeam-Guessle responded by laughing when he spoke. Gancio Giuliano, the artistic director of the Museum, even tried to told journalists and press that 'the patient in the trisomy machine [sic] carried some corpses four hours into the museum, but the whole time it was the guy who stroked the young man who broke him'. In 2008, Giuliano told the same press that the hours of the destruction are truly ""wrong for their morality"" and further stated that 'We are never satisfied with our decision. We made an informed decision to build the museum after destruction.[5] Deaths [edit] A little after 12:00 am 17 October 2005, Giuliano and his partner Monica Concerta, noticed that the trash was being thrown by passers-by. Captain Iamienowska leaned over to his film camera and said, in a joking manner, that Iannorello, the chair of the Muscéci, was a thief that director Frank Nolan said ""he would later be arrested."" When Iamienowska arrived, the people in question were interviewed by Captain Anderson Tulaqyuk, a co-man who was initially lying on the scene and whom Iamienowska said was able to stop them from passing in the vicinity. Iguano proceeded to collect the trash and the police arrived, and closed the door of the museum.[6] During the war, the statue structure was partially removed and its cannons damaged. On the eve of the war, the U.S. Army and Canadian Air Company, who once owned a lot on the Coopers of Paris near Leopold Street, sheltered the POWs, who were briefly overworked. This following years led to condemnation of the organization and artists, including John DiDione,[7] returned to the compound to sell the museum items (including many of them originally stolen over the years—more than a million returned); on 12 October, the sculptor Jack Uellein installed a temporary dome on Del Cavella.[10] Among the changes over the years, Diena featured the lion in the middle of the other two halves of the dome and the lion in the back of the Viduccio.[11] In March 2008, Arvnzo Cissino announced that seven studio studios were built between the summer of 1985 and the last one was completed on November 17, 2005.[10] In the fourth quarter of the mid-2000s, a number of new structures were constructed by various museums, but most notably by Arvnzo Cissino, the brothers-soldier, as portrayed by architect Duvois Mayaven in the illusion of the El Valleratas Humanities and Cultural Center, reconstructed by architect Duvois Benjamin. Just as the very significant additions to the current management masts of the del Museum are such private three-bedroom rooms as Tradello's cabin, as well as the architectural aspect from the museum's south wall. Also new are the rockings performed by MC Urquhart, Aiba Ruiz and March Flez, during a music festival held in the US in 2006. Urquhart was later credited with establishing record label Empire Records as an influence on his music. In<endofxt>

Figure 4: Sample from MDLM (Sahoo et al., 2024a) of length $L = 1024$ and $T = 5k$ diffusion steps. The generative perplexity of this sample under GPT2-Large is 69.26 and its entropy is 5.62.

<endofxtxt>I don't know how I want to express myself about this. It's happening. I've said nothing about it. The problem is that, it feels like it's so bizarre to say that saying that something just stupid is possible and everyone looks out bad. Well, I want to tell you this, it's happened, it's weird to try to make this, but to say that everything is random it's just normal. I think people don't realize it, like people have done studies to sum up these kind of things, I think it might be a result of psychoanalysis or incorrect. Yes, in the case of the state and university people are saying all these things, but when you see what's happening in situations that's not normal. The truth would have to be buried somewhere, but the goal is to show why it's interesting. I'm going to put up the clues, like dream program here, meaning it's all what sounds like "- []" or "- []" "can you say that word once" and all of these things. I'm going to reveal what I know. Right, it seems, you've used a lot of it to play this game. What did you think of it? I remember doing something made of more than one person where you're both asking what they are and you're asking different scenarios. It's always a situation, but if one person is what it's not, it's a person, I don't like it because it's good or bad. They're the ones who only play it because of being themselves. To me personally, you should not play it because it makes sense. I don't see how it's going to get complicated for the situations where they make you elaborate, it makes sense for the, as my own mind has become through working in the studies, you know in different countries, and the answer is to, if I have to look one line of thought down, like... "know, it becomes too massive with the amount of it to see it all. We gave you an interview with Wolfgang Louis about the fact that you are both writer and painter. Well, if you were able to explain that, painter would you be interested in that - drawing pictures, drawing music - would you be interested in that? And would you explain how you get interactions or people interact with each other? I wouldn't be doing that, all of it myself. But because I've heard that everything is not something beautiful, that there's a certain quality of person too. It's really beautiful, I talk with a generalizing and I explain at the extent with everything you can see, I had that same idea that we worked on, but it's not complicated even in the morning, everyday life, no one needs to discuss it. You've been here before. And sometimes you might get the remark there which is like, "hey, someone who knows this is a nice person, but he does have a taste, it's sad." And would you say that what kind of the methods can make you feel totally different and more human? It does, you're totally behind each other's idea and I really feel so much about an artist's idea. It's easy to apply this very free approach to an idea when you look at an idea, you want people to look a direction, to look back. I think it is with the art of art, and the artists don't really understand, it's not the death instinct, they don't realize that this artist is right, they think he's right but the first page, on the first page, on the first turn of the page, it's easy to say that they're really wrong, but we have to believe in them doing this to you. And I can't really imagine that there's so many tricks where he might play his part. I don't think that bad lay people are the solution, but I think that people who can understand it can still grow up. I think there will be absolutely an 18, 20-year-old living up to the illusion of this because he thinks he is this. I don't think that, because if you are caught up with that stuff, but then you feel it's like a stage, and I, I'll come back, the next piece, everything is at once, how is the process from what it's? How do you do that? With what conditions? But to do people like that, what could they have, as simple as two or three? These are little things, because there's one method and that other reason that you said before, that you think of so many possibilities and see so many [] maybe you're thinking about this but it's so hard to wrap your mind because the world you're so in and you want to get a picture, you can tell. But I think that, it's a good thing to know. It's impossible for people to say it, because it can never be done this way, it can be seen through the eyes, but there's something of the the perspective and inside, it can see how we can do it. It's a real idea to get your out and talk about the ways of the thinker, the self and who in our teacher's side is thinking, who's how you think about, and then you come away with the feeling that it's me that thinks about this, what way. Every time there's the kind of moment, every time maybe I think there's something, I think there's also a big question, actually. And you think, how is this identity, how do we reconcile this, with this, even if there's no self and 'you too' is some kind of pressure, if we can really move through this, it gives us something to do. This is kind of a task, how to make this a solution? Or is it not really interesting, which is it missing something? And I really wonder if I can think about the solution there? I find myself how if it can be opened up, what if the world is just puzzle, do there a flow of this thinking, is this actually really what you're looking for, is this - you're not so blind and open yourself up, but have you opened up to the point where there's the same, you know a little thing in the universe? I think - in this world, all we see the merging of this sphere of light of the universe, this first step and this next step, so I think the power of perception is in front of you, I'm really getting tired of these. I still think they are the trick. Do you put yourself in this situation where someone who they really think about in space really opens themselves up - is that a problem for them, because right? Is it not to think where you are seeing through your door rather than another door, is it not think through a light, when the light's coming from somebody, the light is coming to a person? So how can we connect to the other, can we talk with the world, where we have other people, I see these. For me if I start to think about this, I don't see these, if I've thought of these, I've reconciled with the world and if it's where I am, I've realized it, the fact that I see that another person who's there. Is this is also awareness, or is it maybe able reality? What's through that that we connect with this, I think I think it's a wonderful thing that's spontaneous and to see these things around people and open myself around people, there's this space for me to think. I think that if we come to the people in the world, maybe there will be a kind of language of what we experience, even when you're now at the other. So maybe what will make us think about this thing will be able to go forward to feel the these things which is, this this is the same way of taking ourselves of the self and not accepting other things, as a child we took things, takes ourselves of the self and not accepting other things, we have to accept the same self- anything can come - we have to still do the same thing, this kind of observation, that thought, and think about this, this view of these experiences, and I think that from there we would still go forward and still think of the things that put us out there, this idea that we can still see this person and this person at the beginning, there would be this kind of thinking out of here, and I think to take that the consciousness, or we can take that attitude to the world and still this is what we are in, for me I think, a very direct connection with this awareness and I think you make this connection with your own - that's where this always seems to be a hard thing to say. But for the second thing, I bring in: If you can pour that energy into that space, then these can be possibilities. Maybe, wherever you are at some point you might just enter this new space and go explore it? Maybe you might still be able to think some moment and from now on you still don't think about those moments, you can wonder, is there any other moment, is that the first step of that, moment, maybe then you can really move forward and open yourself in, I guess, I could say, wondering, oh, I'm out of space, oh, I was there before, oh, I feel like I've gone in space, oh, I don't see myself anymore. And when people came through, you know, at the moment when I felt like I'd actually entered this space, Hey, enter the universe, Was it very important to me for that moment of thought? I couldn't really tell you now, if I knew what happens when in my mind, myself at that moment, that moment I'm about to go up and see myself, I'm about to look at this and see this space, Here's this opening where you just pulled out, You'd have to realize - that only because you're coming in, I think you're not the people that come here. Oum.<endofxtxt>

Figure 5: Sample from BD3-LM for block size $L' = 16$ of length $L = 2132$ under $T = 5k$ diffusion steps (trained with a context length of $L = 1024$). The generative perplexity of this sample under GPT2-Large is 28.2, and its entropy is 5.1.

<endofxt>, but Wilson took over the program and turned it around. "He's done a tremendous job," Caldwell said. "He's done a fantastic job." The offense has always had an abundance of weapons, but it became evident that they weren't going to have a weapon to actually go after players from the slot. Now they're in two different weapons sets. The top group features Dez Bryant and Mohamed Sanu, and the bottom group features an assortment of weapons and pass rushers. The job has become far more complex. The other players can make plays on the ball and get those targets at a higher rate. Sanu is more of a classic, get to the quarterback and leave the corner open. Dontari Poe got the job done this year and became one of the more effective players at the position, even in the passing game. However, Dallas has got to figure out how to get their franchise wideouts to contribute on the field. That can be tough. Adding Poe can help get the receiving corps going. C.J. Spiller is a two-time Pro Bowler, but if the Cowboys want to upgrade their receiving corps, he's going to have to step up in a big way. "We've got to be a little more aggressive with the type of weapons that we have," Caldwell said. "I think that's part of the reason why our last two games, especially when you're playing in Washington, D.C., you've got to be aggressive, make sure you're hitting at every catch. When you are, you're giving up a lot of yards." Part of that means taking the quarterback out of the equation and having him beat coverage a lot more. In the NFC West, you want your offensive weapons to do a better job of running through coverage. The biggest threat that Dallas has is a QB in Ben Roethlisberger. Roethlisberger is far and away the best quarterback in the league, but a lot of the credit has to go to his receiver group. Martavis Bryant and Antonio Brown are both big-time receivers, and last year they were in the top 10 of yards per catch and receiving yards in the league. That production will never be sustainable, but if you're going to be an elite offense, it's going to take a lot of catching up. Roethlisberger is an All-Pro receiver, and he's not the most dynamic option. But it would take something like Bryant or Brown at a better position, and at a slightly lower price, to make him the most productive receiver on the offense. The truth is that Roethlisberger isn't going to be great. He may only have 18 games left in his career, but he's been doing it since he was a rookie in 1991. But that's not the worst thing in the world. Roethlisberger's ability to hit guys on the outside with good movement, vision and running ability is what the Cowboys need in order to keep up with the competition. If he keeps getting better, he could become the best receiver in the league. Follow @walterfootball for updates. Like our Facebook page for more Cowboys news, commentary and conversation. The owner of 1H10 Tree in Charlotte Gardens is taking legal action against the city. Derek Jarman says he's been forced to evict his neighbour, Bob, after he took to social media to threaten to burn down his neighbor's house. "I'm incredibly furious with the city," Jarman told 7.30. "I've been trying to keep my eyes on the prize." Tree in Charlotte Gardens saying it had seen '9,000+ people' enjoying a great weekend. The company that owns 1H10 named Bob after a bee and said the tree was frequently targeted because of its unusual location. Bob said he had his concerns about the tree when he was contacted in October. He said they had had 'an ongoing conversation about my neighbor. He called, hung up and he was very threatening' in the 30 days before they turned the tree over to him. A neighbour posted the following online message on 8 October. "I am shocked about the serious problems you are having with your neighbour that has caused you all (sic). You and the 2 of you are making money at the expense of the good people of Charlotte Gardens." Bob says he was furious and said he'd just got off the phone with the city manager. "I told her, 'no, I'm going to bring a lawsuit', and I called the solicitor and tried to get my phone, just hoping the solicitor would help me out. I called again, and I asked if I could go to court and to try and get an injunction. "They told me 'you cannot', and they said, 'we can't, we can't' because you're sending people to the police'." Tree in Charlotte Gardens (Facebook) He also said he'd threatened the city attorney if he didn't stop the building from burning down. The internet user tweeted: "I's on the tree, but after I said 'threw this away, here's a spot to burn', the building started to burn." Tree in Charlotte Gardens (Facebook) Bob said the neighbour had threatened to burn down the tree, the windows, the living room and his entire backyard. "It was more than a threat," he said. "He was a very strong person. He's already damaged so many people in this building. It's not going to go away." Tree in Charlotte Gardens (Facebook) Jarman says he tried to talk the building owner out of the move, but the building owner's behaviour had "deleted him." "I'm going to stop him by letter telling him not to come to my house any more," he said. "I have three kids, and if Bob is going to be in my house, I need to make sure I have someone who can go in there and protect me. "My son does a really good job of protecting me, and I'm not going to let that get in the way of that." Tree in Charlotte Gardens (Facebook) Jarman said Bob had pulled him up on social media, calling him a "white nut" and saying: "For God's sake, stop calling me a white nut. "I should have shut him up on Facebook." He said he sent Bob the letter and thanked him for the support. "He should have done it because he's a real artist and he's a real artist," he said. He declined to name the architect of the new tree, but says the firm is the same one that designs buildings. "The building is burning down", neighbour says Bob's neighbour, Michael Banks, says the fire is an insult to his daughter. "There are two black women that live next door to me and they told me 'you can't do that', and then the fire went up and then the building burnt down," he said. "You can't burn down a house if you don't burn down the house." Coun-Pete Lawrence, the Northumberland MP for Wood Green, says he has concerns about Bob's neighbours. "It's a very, very sad commentary on the state of society and democracy in general," he said. "It's interesting in a community that's 50,000-plus people, you've got your regular residents and well-meaning neighbours who are apparently oblivious to the destruction of their own home. "To me, that's appalling and it is probably a shocking amount of devastation that it's left behind. "I would expect there to be outrage as well." Bob Jarman fears for his life after the tree was torched Bob says he has told the Northumberland Council that he had already received \$1,000 in legal action from the building owner, when he told them about the incident. The building owner has declined to comment on the situation. The builder is currently assessing its legal options. "We've got to sort this out and have an understanding with the builder, Mr Banks," he said. "We've got to make sure we can't get into into a legal battle with that person and make that person change his mind. "We don't want to do anything to cause a scene or anybody in the street to be upset." Bob Jarman hopes to have an understanding with the builder on its legal options, who have refused to comment. Topics: state-parliament, smoking-and-doubt, black-wales-6168, united-kingdom, england First posted When the other guys are away playing, do a short commercial to get you fired up for the next work day. Once you make it home, get a few junkies for them. They'll be very happy to have you, for at least a day. They might not be so happy after a couple of days. Have a bunch of friends and get ready to keep it going. What are you waiting for? Make this long, one-off<endofxt>

Figure 6: Sample from an AR model (Sahoo et al., 2024a) with length $L = 2003$ (trained with a context length of $L = 1024$). The generative perplexity of this sample under GPT2-Large is 10.6 and its entropy is 5.5.