

SUNMASK: Mask Enhanced Control in Step Unrolled Denoising Autoencoders

Kyle Kastner^{1,2}, Tim Cooijmans^{1,2}, Yusong Wu^{1,2}, and Aaron Courville^{1,2*}

¹ University of Montreal, Montreal QC H3T 1J4, CAN

² Mila - Quebec AI Institute

Abstract. This paper introduces SUNMASK, an approach for generative sequence modeling based on masked unrolled denoising autoencoders. By explicitly incorporating a conditional masking variable, as well as using this mask information to modulate losses during training based on expected exemplar difficulty, SUNMASK models discrete sequences without direct ordering assumptions. The addition of masking terms allows for fine-grained control during generation, starting from random tokens and a mask over subset variables, then predicting tokens which are again combined with a subset mask for subsequent repetitions. This iterative process gradually improves token sequences toward a structured output, while guided by proposal masks. The broad framework for unrolled denoising autoencoders is largely independent of model type, and we utilize both transformer and convolution based architectures in this work. We demonstrate the efficacy of this approach both qualitatively and quantitatively, applying SUNMASK to generative modeling of symbolic polyphonic music, and language modeling for English text.

Keywords: Artificial neural networks · Non-autoregressive sequence modeling · Generative modeling

1 Introduction

Modern approaches to content generation frequently utilize probabilistic models, which can be parameterized and learned using artificial neural networks. Common types of neural probabilistic models used for generation can be broadly stratified to form two broad categories based on factorization: autoregressive models (AR), and non-autoregressive models (NAR). We introduce SUNMASK, a NAR generative model for structured sequences¹.

1.1 Autoregressive Models

AR modeling with deep neural networks has been a dominant approach for generative modeling and feature learning [36, 54, 62, 65, 66] which has many crucial

* CIFAR Fellow

¹ <https://github.com/SUNMASK-web/sunmask>. Last accessed 8 February 2023.

advantages in both training and inference. One key concern is the necessity of defining a "dependency chain" in the form of a (typically) directed acyclic graph (DAG). Sampling during inference can be accomplished in a straightforward manner using ancestral sampling - sampling from the first variable or variables in the DAG, using those to conditionally estimate a probability distribution for subsequent variables.

Many applications have straightforward orderings in which to define this chain of variables, based on domain knowledge. For example following the flow of time for timeseries modeling is often a logical choice, allowing models to make predictions into the future from the past. However in many other domains, for example images, language, or music, the process of defining a dependency chain over input variables (e.g. pixels, characters, words, or notes) is far from straightforward, as for any arbitrary ordering there can frequently be examples where this ordering *creates* long-term dependencies, or otherwise makes satisfaction of dependencies during training and evaluation more difficult than another alternative ordering.

This divide becomes further compounded in many creative applications to these domains, as creators typically iterate repeatedly: forming a concept, sketching out the concept, and seeing where the creative flow (based on the sketch) may lead to alterations in the original concept, thus "rewriting" sketched steps. Though the resulting output may be perceived in a time-ordered fashion (for example, reading a book or listening to a song), the initial creation was performed globally and holistically. This global view is often critical to creating elements such as foreshadowing and tension which make the resulting output interesting or enjoyable. This iterative process is directly at odds with a strict AR factorization, and requires well trained AR models to cope with a high degree of uncertainty and multi-modality for long range dependencies, which can lead to logical mistakes or other errors.

1.2 Non-Autoregressive Models

An alternative methodology for generative modeling is non-autoregression (NAR), broadly covering a large number of different modeling approaches which attempt to remove assumptions about variable ordering, instead either hand-defining per-exemplar orderings, or modeling variables jointly without resorting to chain rule factorization [23, 29]. One way to define an ordering over variables is via masking of inputs or intermediate network representations [20, 51, 63–65, 68], and indeed modern AR approaches such as transformers [67] use an autoregressive mask internally to define the chain of variables order. These masks can either be constant over all training (as in standard AR transformers and PixelCNN [65]) or dynamic per example (as in MADE [20]). When masks are dynamic per example, we begin to see the relationship between enforcing AR via masking and NAR methods, as although some ordering is assumed this ordering is no longer constant, and it becomes possible to use the same trained model to evaluate the probability of a particular output variable under *multiple* possible orderings.

Closely linked to masking methods are so called *diffusion models*, which relax the variable ordering problem through noise prediction [30, 59, 61]. Rather than

predicting a new variable or variables given previous ones in an arbitrarily chosen DAG, diffusion models focus on predicting a less noisy version of many variables jointly, given a set of noisy input variables. Iteratively applying this learned denoising improvement operator should eventually result in predicting a fully clean output estimate, given either a noisy version of the target domain, or even starting from pure noise. Given this framing it is clear that diffusion models are closely linked to denoising methods in general, specifically denoising autoencoders, as well as modern density modeling approaches such as generative adversarial networks (GAN [21]), variational autoencoders (VAE [38]), flow-based models (NICE [13], RealNVP [14], Normalizing Flows [56], IAF [39], MAF [51]), iterative canvas sampling (DRAW [22]), and noise contrastive estimation (NCE [26]). Particular applications of this denoising philosophy such as BERT [12], WaveGrad [9], and GLIDE [50], have resulted in large quality improvements for feature learning and data generation for text, images, and audio [28, 41, 55].

1.3 Trade-offs Between Autoregressive and Non-autoregressive Approaches

The choice between AR and NAR methods is not clear-cut. For many domains, high-quality models exist using both approaches but we can define some crucial parameters. Some NAR methods such as GAN or VAE are capable of generating output in only one inference step, however they are typically hard to train on certain data modalities (e.g. text data) comparing to AR counterparts. Other NAR methods such as diffusion models typically allow for choosing a diffusion length during inference, which is independent of that used at training. Choosing a low diffusion length can frequently lead to poor sample quality, and tuning this setting (among many others) is critical to high quality generation. However if the tuned diffusion length for a given sequence (of length T) is shorter than the length of those sequences, the NAR method has a computational advantage over the equivalent AR model (which would require T steps for a T length sequence). In addition, the ability to tune this diffusion length can be useful in interactive applications, or when a variety of output is desirable. This setting can also be a curse, as even well-trained models perform poorly with improper diffusion settings. Several branches of current research are focused on improving guarantees and convergence speed for diffusion models [35, 37, 40, 60].

1.4 SUNMASK, a non-autoregressive sequence model

We introduce SUNMASK, a NAR sequence model which uses masks over noised, discrete data to learn a self-improvement operator which transitions from categorical noise toward the data distribution in iterated steps. Given a target data representation, we train a model which can map from a noisy version of input data to a corrected form of the input. In this work, we use multinomial noise - namely entries are corrupted to 1 of P possible values (for a given set size P), with the number of noised entries in a sequence defining the relative noise level for the training example. This is similar to many diffusion approaches at a

high level, and particularly shown to be an effective tool in SUNDAE [58] and Coconet [33]. In addition to the use of multinomial noise, we also form a mask representing *where* the data was noised, feeding this mask alongside the input data to form a conditional probability distribution.

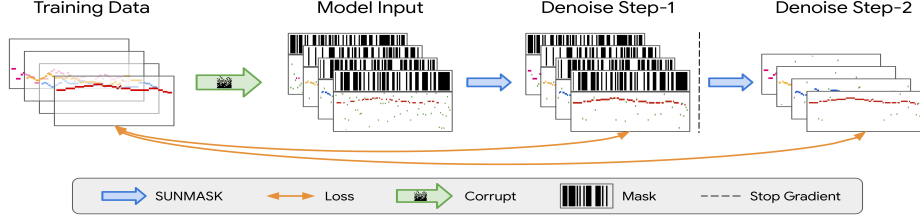


Fig. 1: Step-unrolled denoising training for SUNMASK on polyphonic music, unrolled step length 2. Training data (left) consists of four voices corrupted by sampling a random mask per voice and replacing the masked data (red) with random pitches (green). SUNMASK takes both mask and corrupted training data as input, predicting denoised original data as output. In the second step, the model takes a sampled version of the model step predictions and the same mask as input, outputting another prediction of the original data.

2 Method

The relationship between discrete diffusion and denoising autoencoders has been explored in previous work [4, 31, 32, 58]. We build upon this foundation, combined with many insights from prior orderless modeling approaches, crucially Orderless NADE [64], Coconet [33] (which is a more modern variant of Orderless NADE), and SUNDAE [58].

SUNMASK is built around a process $x_t \sim f_\theta(\cdot|x_{t-1}; m)$ on a space $X = \{1, \dots, v\}^N$ of arrays with categorical variables. This parametric transition function f_θ takes an additional argument $m \in 0, 1^N$. During training, m indicates variables that were not initially corrupted, and as a consequence we can use it during inference to tell f_θ which variables to trust.

Given a sequence of masks m_0, \dots, m_{T-1} , the generating distribution of our model derives from a prior p_0 (typically uniform noise) and repeated application of f_θ :

$$p_T(x_T; m_0, \dots, m_{T-1}) = \left(\sum_{x_1, \dots, x_{T-1} \in X} \prod_{t=1}^T f_\theta(x_t|x_{t-1}; m_{t-1}) \right) p_0(x_0) \quad (1)$$

In practice, p_0 is typically elementwise iid uniform noise, and the masks m_0, \dots, m_{T-1} are drawn according to a schedule and may be held constant for several steps.

To train f_θ , we take a training example $x \sim p_{\text{data}}$ and draw a mask m . We apply the corruption procedure $x_0 \sim q(\cdot|x; m)$ to obtain x_0 which equals x where the mask m is true and uniform random values elsewhere. Then we iterate $x_t \sim f_\theta(\cdot|x_{t-1}; m)$ with the aim of reconstructing x .

As in SUNDAE, the transition f_θ models the variables as conditionally independent of one another. However SUNDAE has no direct concept of masking. SUNMASK thus combines past insights from the masked NAR models Orderless NADE and Coconet with existing concepts from SUNDAE, along with new model classes and inference schemes to form a powerful generative model. Similar to SUNDAE, our objective is to minimize $\frac{1}{2}(L^{(1)} + L^{(2)})$ where

$$L^{(t)}(\theta) = -\mathbb{E}_{m_0, \dots, m_{t-1}} \mathbb{E} \left[\frac{\sum_i (1 - m_{t-1}^{(i)}) \log f_\theta^{(i)}(x^{(i)} | x_{t-1}; m_{t-1})}{\sum_i 1 - m_{t-1}^{(i)}} \right]$$

$$\begin{aligned} x &\sim p_{\text{data}} \\ x_0 &\sim q(\cdot | x, m_0) \\ x_1 &\sim f_\theta(\cdot | x_0; m_0) \\ x_2 &\sim f_\theta(\cdot | x_1; m_1) \\ &\dots \\ x_{t-1} &\sim f_\theta(\cdot | x_{t-2}; m_{t-2}) \end{aligned} \quad (2)$$

is the reconstruction loss for the elements of x that were corrupted. As in Orderless NADE [64] and Coconet [33], we weigh each term according to the size of the mask, to ensure that the overall weight on each conditional $f_\theta^{(i)}$ is uniform across i . Unlike previous methods, we target *only masked variables* in the loss. In practice we choose $m_0 = \dots = m_{t-1}$ during training and $t = 2$. Since we only go to $t = 2$, keeping the mask constant is a close enough approximation to the masking schedule used in inference. The choice of $t = 2$ is driven by the ablation study in SUNDAE, where $t = 2$ was found to account for nearly all performance gains in translation experiments, with higher unrollings showing no additional benefit. In addition higher values of t unrolling generally increase memory usage, making the training of high order unrollings complicated.

SUNMASK allows for direct control at inference using both proposal masks and noising of variables, combining elements of both SUNDAE and Coconet. We show a high level example of the unrolled training scheme, mask proposals, and input data processing in Figure 1.

The overall unrolled mask and iterative inference setting is largely independent of architecture choice, and as long as the internal architecture does not make any ordering assumption over the input data we can incorporate it into SUNMASK. We use two primary archetypes for the internal model in this paper: Attentional U-Net and Relative Transformer.

SUNMASK uses an unrolled training scheme, similar to that shown in SUNDAE, as well as a mask which is input to the model and defines manipulated variables as in Coconet. The loss is masked based on this manipulation mask, unlike Coconet or SUNDAE. The SUNMASK loss is further weighted by the total amount of masked variables. Comparisons of various high level modeling features between SUNMASK, Coconet, and SUNDAE are shown in Table 1.

2.1 Model Training

During training, the internal architecture is combined with a *step unrolled* training procedure, as highlighted by SUNDAE [58]. Rather than directly randomizing positions, we re-write this as a masking scheme, first sampling a mask (with

Table 1: Comparing SUNMASK, Coconet, and SUNDAE

Model	SUNMASK	Coconet	SUNDAE
Mask input to model	✓	✓	X
Masked loss	✓	X	X
Re-weighted loss	✓	✓	X
Unrolled loss	✓	X	✓
Inference mask schedule	✓	✓	X
Sampling rejection step	✓	X	✓
Mask control preserves data	✓	X	X

0 randomize, 1 keep, which we denote as 0-active format) then performing randomization to one of P possibilities, for the masked subset of K variables. This random masking procedure is equivalent to the approach from SUNDAE, but using a mask allows us to further combine the mask information with the input data, in order to form a conditional probability estimate. In addition, this 0-active masking scheme makes direct comparison to masking schemes with absorbing states (such as OrderlessNADE [64], Coconet [33], VQ-Diffusion [24] and OA-ARDM [31]) simpler, as the mask can be directly multiplied with the data in a 0-active format.

Each training batch is randomly sampled from the training dataset, and a corresponding noise value drawn from $rand(N)$ for N examples in the minibatch. This per-example noise value is then used to derive a per-step mask over T timesteps, by comparing noise $rand(N) < rand(N, T)$. During training, this means some examples have a high per-example noise value (e.g. .99), and thus many values masked and noised in the training, while other examples may have a low noise value (e.g. .01) drawn instead. Combined with a training loss which learns to denoise the input and focuses on imputing information about masked corrupted inputs, the overall model will learn a chain to go from more noisy data to less noisy step-wise, resulting in a learned improvement operator [32, 58].

This improvement operator can be applied to noisy data or pure noise, and iterate toward a predictive sample from the training distribution. See Multinomial Diffusion [32] and SUNDAE [58] for more detail on this proof, as well as fundamental work on denoising autoencoders [2]. In SUNMASK, we combine the mask used to noise the input with the input data itself, while modifying the loss to predict *only masked variables*. In addition, we downweight the loss by $\frac{1}{1+\sum 1-m_t}$ for each batch element, meaning that losses for heavily masked entries are downweighted compared to losses on examples with little masking, in a form of curriculum weighting based on expected estimation difficulty.

While a one step denoising scheme can be sufficient for learning the data manifold [2, 4], *unrolling* this denoising scheme into a multi-step process can have performance benefits. SUNMASK directly uses the unrolled loop scheme described in [58], using a step value of 2. For a detailed description of the step unrolled training scheme, see the overview description from SUNDAE [58]. The masked and unrolled training can be seen as a container for any internal model which does not make ordering assumptions, and we utilize both convolutional U-Net (a variant of the GLIDE [50] U-Net) and Relative Transformer [11, 34, 53] models for various experiments, shown in Section 4.

2.2 Convolutional SUNMASK

SUNMASK is most closely related to Coconet [33] and SUNDAE [58]. Coconet (as an instance of OrderlessNADE using convolutional networks), trains by sampling a random mask per training example, using this mask to set part of the input (in one hot format) to zero. The mask is further concatenated to the zeroed data along the channel axis, and this combined batch is passed through a deep convolutional network with small 3×3 kernels. Convolutional SUNMASK uses a downweighted loss over only variables masked in the input. However, SUNMASK additionally uses the unrolled training scheme, as well as a different inference procedure due to preserving the values of masked out variables during training and sampling.

Our best performing convolutional SUNMASK architecture takes hints from recent image transformer and vector quantized generators, exchanging the small kernels used in Coconet for extremely large kernels of shape $4 \times P$ over the time and feature dimensions, somewhat analogous to input patches, removing the model’s translation invariance over the feature axis by setting kernel dimension equal to the total feature size. However this makes the number of parameters per convolutional layer extremely large. Convolutional SUNMASK adopts an attentional U-Net structure which reduces only across the time axis, modified from GLIDE [50], rather than the deep residual convolution network used by Coconet. Combined with the addition of step unrolled training, we are only able to train convolutional SUNMASK with a batch size of 1 (expanded to effective batch size 2 due to step unrolling) on commodity GPU hardware with 16GB VRAM.

Due to the design choice of extremely large kernel sizes which depend on the size of the domain, we only use convolutional SUNMASK for polyphonic music experiments, see Section 4 for more details.

Attention is applied on the innermost U-Net block size as well as the middle block, with 1 attention head [50]. Convolutions are used in all resampling, and all resampling happens only on the time axis, making the Attentional U-Net effectively a 1-D architecture. However, rather than learning both instrument and pitch relations across channels, we isolate pitch relations and instrument relations into separate axes of the overall processing, the "width" and "channels" axes, respectively assuming $(N, C, H, W) == (N, I, T, P)$ axes. As is standard in many U-Net designs, we double the number of hidden values for layers every time the resolution is halved, with the reverse process being used when upsampling. Though the parameter count here is large, it is similar in spirit to other approaches to small datasets on text [1].

2.3 Transformer SUNMASK

Transformer SUNMASK relates closely to the transformer used in SUNDAE. The architecture uses a relative multi-head attention [11, 34] and has no autoregressive masking. SUNMASK transformer also uses larger batch sizes, typically 20 or larger, though this is far smaller than the batch sizes seen in the experiments of SUNDAE. Sequence length and data iterator strategy were both a critical aspect

for training transformer SUNMASK. We found short sequences (from 32 to 128) worked best, along with iteration strategies that were example based.

Transformer SUNMASK was trained on every dataset used in this paper, and we show performance in Section 4, as well as comparisons to convolutional SUNMASK on symbolic polyphonic music modeling. Both convolutional and transformer based SUNMASK use the Adam optimizer, with gradient clipping by value at 3. Inference hyperparameter types and general sampling strategies used are the same with both models, though specific hyperparameter values may change between datasets.

There is a large discrepancy in model parameter count between our best performing convolutional models for JSB, and our best transformers. Training larger transformers can work well for generation [1], but our large parameter transformers (on the order of 400M parameters) had poor generative performance on JSB.

Pitch size / vocabulary size, sequence length, and batch size changed for the transformers used in the text experiments, but the global architecture remained in the style of "decoder only" transformers [54], similar to SUNDAE. Notably, we use vocabulary size 5.7k, sequence length 52, batch size 48 for EMNLP2017 News and vocabulary size 27, sequence length 64, batch size 20, and a slightly extended training step length of 150000 for text8.

2.4 Inference Specific Settings

Well-trained SUNMASK models should be applicable to full content generation, as well as a variety of partially conditional generative tasks such as infilling and human-in-the-loop creation. Basic sampling involves creating a set of variables, with all variables randomly set to 1 of P values in the domain (or partial randomization in the case of infilling) along with an accompanying mask, which is initially all 0 for full generation, or mixed 1s and 0s for partial generation tasks. Given this data and mask as input, the trained model then predicts a probability distribution over all possible P values, for all variables. Despite the use of masked losses in training, we sample these prediction distributions for *all* variables. These predictions are then accepted or rejected from the original set, resulting in a new variable set. We then sample a new mask (based on a predefined schedule) and combine it with the initial mask, then iterate this overall process, updating at least some of the variables at each step.

During inference we use several key techniques to improve generative quality. We use typicality sampling [47] on the output probability distribution and a variable number of diffusion steps, on the order of 100 to 2000. Masks are randomly sampled using the schedule defined in [33] which linearly decreases the number of masked variables over time according to $\alpha_n = \max(\alpha_{\min}, \alpha_{\max} - \frac{n}{\eta N}(\alpha_{\max} - \alpha_{\min}))$ with $\alpha_{\min} = .001$, $\alpha_{\max} = .999$, and $\eta = 3/4$, along with an optional triangular linear ramp-up and ramp-down schedule for the probability of accepting predictions from the model into the current variable set at each step, as shown in [58].

Tuning hyperparameters for inference is critical to success, as improper settings can drastically lower the performance of SUNMASK, see Section 4 for variance over various inference settings in different tasks. For human-in-the-loop applications, the existence of these controls can allow a number of fine-grained workflows to emerge, driven by expert users to create and curate interesting output [10, 16], demonstrated in Figure 3.

3 Related Work

We state here some key related approaches, as well as how our method differentiates from these previous settings. A number of recent publications on diffusion models and feature learning have incorporated masks as part of their overall training scheme [28, 31], however these papers use masks for blanking, rather than as indicators over stochastic variables. Many infilling models [12, 15], and masked image models [28] feature conditional modeling with a mask (blank) token, predicting the variables masked from the input for feature learning or generative modeling. XLNet [68] combines the infilling and autoregressive paradigms, learning arbitrary permuted orders over masked out variables, using blank-out masking and randomly generated autoregressive ordering similar to OrderlessNADE and Coconet. Conditional diffusion generators [48] and GAN generators [18] have the combination of mask indicators as well as preserving stochasticity of the masked variables. However these methods do not use an unrolled training scheme, and generally target image related tasks, with the notable exception of maskGAN. Many models use a concept of a working canvas, and do repeated inference steps for generation or correction of data [5, 19, 22], SUNMASK differs from these models due to architecture choices, training scheme, and loss weighting, as well as application domain [49, 50, 52, 57].

4 Experiments

We demonstrate the use of SUNMASK for polyphonic symbolic music modeling on the JSB dataset [3, 6]. The JSB dataset consists of 382 four-part chorales, originally written by Johann Sebastian Bach. These chorales are quantized at the 16th note interval, cut into non-overlapping chunks of length 128, skipping chunks which would cross the end of a piece. This processing results in a training dataset of 4956 examples, with each example being size (4, 128). We train convolutional and transformer versions of both SUNMASK and SUNDAE for comparison, as well as the pretrained Coconet [33]. For polyphonic music, the quantized data was rasterized in soprano, alto, tenor, bass (SATB) order, as in Music Transformer [34] and BachBot [42], then chunked into non-overlapping training examples. Results are shown in Table 2. These results are evaluated on Bach ground truth data (Bach GT), BachMock Transformer (BachMock [17, 44]) (closely related to the decoder from VQ-CPC [27]), Coconet, SUNDAE (SD), and SUNMASK convolutional (SMc) and SUNMASK transformer (SMt). Model sampling variants are indicated as Typical Sampling (T).

4.1 Musical Evaluation

The grading function used for evaluation, referred to as BachMock here, is designed specifically to correlate with expert analysis on Bach. In particular using this metric to choose correct examples in a paired comparison test, outperforms novice, intermediate, and expert listeners by varying margins [17]. This indicates that scoring well on the aggregate metric should correlate to high sample quality. The metric has many sub-parts, ranking various musical attributes crucial to codifying the style of J.S. Bach. AugGen [44] incorporated this metric into an iterative training and sampling scheme which improved final generative capability for a fixed model, showing the effectiveness of BachMock in practice for ranking machine generated samples. For each grading function in the Bach Mock grading evaluation, we show the median value and \pm standard deviation (showing the average of each interval SATB performance for brevity), as well as the overall grade. Lower values for all metrics are better, and we see the strongest results for convolutional SUNMASK with typicality sampling. Combined with final top-N ($N = 20$) selection out of a candidate set of 200 samples, the overall sample quality outperforms strong baselines. This high quality subset (SMc-T BEST20) rivals both the "BachMock" transformer and the dataset itself on this metric. We find SUNMASK generations are qualitatively good and listenable overall, even though some SUNMASK samples do fare poorly by the grading metrics.

Table 2: Quantitative results from the Bach Mock grading function [17]. Top rows compare to existing literature, bottom rows show ablation study of SUNMASK style models. Lower values represent better chorales.

Model	Note	Rhythm	Parallel Errors	Harmonic Quality	Interval	Repeated Sequence	Overall
Bach Data	0.24 \pm 0.15	0.23 \pm 0.14	0.0 \pm 0.69	0.41 \pm 0.2	0.55 \pm 0.4	1.29 \pm 0.88	4.91 \pm 1.63
BachMock	0.37 \pm 0.22	0.26 \pm 0.14	2.16 \pm 3.22	0.54 \pm 0.31	0.71 \pm 0.68	1.86 \pm 2.81	8.94 \pm 4.64
SMc-T BEST20	0.39 \pm 0.16	0.53 \pm 0.26	0.0 \pm 0.81	0.68 \pm 0.27	0.75 \pm 0.42	1.44 \pm 0.52	7.16 \pm 0.97
AugGen	-	-	-	-	-	-	8.02 \pm 2.92
Coconet	0.44 \pm 0.23	1.85 \pm 0.39	2.61 \pm 6.56	1.38 \pm 0.39	0.86 \pm 0.73	6.07 \pm 1.76	17.00 \pm 6.58
SD	0.59 \pm 1.82	0.93 \pm 0.84	6.42 \pm 4.11	0.98 \pm 0.67	1.99 \pm 5.68	2.45 \pm 2.39	23.25 \pm 21.45
SD-T	0.63 \pm 2.40	0.60 \pm 0.96	3.82 \pm 4.98	0.96 \pm 0.64	2.50 \pm 5.03	1.52 \pm 3.43	20.09 \pm 23.88
SMc	0.87 \pm 2.05	0.63 \pm 0.77	1.38 \pm 6.00	1.02 \pm 0.49	2.07 \pm 5.72	2.32 \pm 2.31	22.47 \pm 20.80
SMc-T	0.57 \pm 1.79	0.69 \pm 0.35	1.28 \pm 3.73	0.93 \pm 0.49	1.10 \pm 4.68	1.81 \pm 0.83	13.43 \pm 19.27
SMt	3.00 \pm 1.85	0.74 \pm 0.90	0.00 \pm 1.95	1.64 \pm 0.70	7.90 \pm 5.58	3.10 \pm 2.97	42.87
SMt-T	3.74 \pm 2.16	0.58 \pm 0.56	0.00 \pm 2.56	1.73 \pm 0.73	7.74 \pm 4.73	2.35 \pm 1.79	46.21 \pm 17.30

4.2 Text Datasets

The EMNLP 2017 News dataset is a common benchmark for word-level language modeling [7], containing a large number of news article sentences [45]. Preprocessing steps collapse to sentences containing the most common 5700 words, resulting in a training set of 200k sentences with a test set of 10k. The overall maximum sentence length is 51. Common processing for this dataset includes padding all sentences up to this maximum length, different than the standard long sequence chunking commonly used in other language modeling tasks.

We show the results of several SUNMASK models for generating sentences similar to EMNLP2017News, comparing to benchmarks using the standard Negative BLEU/Self-BLEU evaluation [7, 70] over generated corpora of 1000 sentences in Figure 2. This set of scores, varied across temperature, is compared against baseline scores [8, 25, 43, 46, 67, 69], similar to the evaluation shown in SUNDAE [58]. These reference benchmarks used 10000 sentences to form performance estimates.

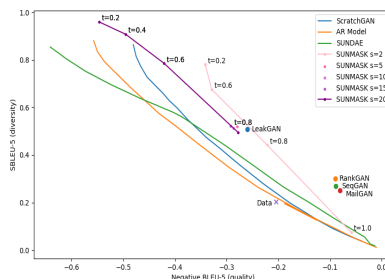


Fig. 2: Negative BLEU/Self-BLEU scores on EMNLP2017 News. Left (x-axis) is better, lower (y-axis) is better. Quality/variation is controlled by changing the temperature (t), and varying diffusion schedule (s). For SUNMASK, *typical* sampling results [47] are shown.

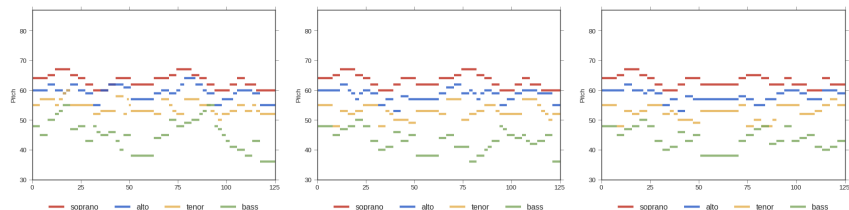


Fig. 3: SUNMASK harmonization (bass, tenor, alto) of an existing melody (left), with a mask which highlights the left half (0 to 64) soprano voice (middle), or a left half mask but replacing right half melody as well (right)

4.3 Music Control

Given the flexibility of masking at inference, we perform a number of qualitative queries to inspect how the model adapts based on noise and mask value. Figure 3 demonstrates the use of SUNMASK for musical inpainting, holding the top voice (soprano) either fully or partially fixed to the well-known melody "Ode to Joy", by Ludwig van Beethoven.

4.4 Text Control

Masking can also be used to variably increase or decrease the weight on various pre-specified terms, held fixed throughout inference. The combination of these

words, and their mask status can be seen to influence the overall tone of the selected text passages which showed the strongest effect in a particular inference batch. The following qualitative samples using masks for word influence are drawn from SUNMASK Transformer on EMNLP2017News dataset. Though the generation quality is flawed, we clearly see a relationship between the masked word and the emergent surrounding context, for example highlighting **disaster** draws forth injured, displaced, and pressure, while **success** instead references happy, nice, good, and playing.

Success unmasked, *disaster* masked

- I think I want to leave **success** at the end of the *disaster* , but because that 's a nice to say it 's not good to be the challenge and this is a very good thing <eos>
- That was the job I was **success** to have to pay my *disaster* but hopefully I have been able to pull playing in the first couple of the season , I 've been happy to go through this team , he said <eos>

Success masked, **disaster** unmasked

- Although more than 80 , 000 *success* have been displaced in the **disaster** since the last year , more than 700 , 000 lives have been injured in the country , and 70 of them were killed , according to the UN media <eos>
- I haven 't had a *success* at the league , the **disaster** and picked running with the door ago we have Champions , and I was a couple of pressure . . . and it was a lot of times <eos>

5 Conclusion

We introduce SUNMASK, a method for masked unrolled denoising modeling of structured data. SUNMASK separates the role of masking and correction by conditioning predictions on the mask, allowing for fine-grained control at inference. When applied to text as well as symbolic polyphonic music, SUNMASK is competitive with strong baselines, outperforming reference baselines on music modeling. Leveraging the separation of mask and noise allows for subtle control at inference, paving the way for a variety of domain specific applications and generative pipelines for human-in-the-loop creation.

References

1. Al-Rfou, R., Choe, D., Constant, N., Guo, M., Jones, L.: Character-level language modeling with deeper self-attention. In: *Proceedings of the AAAI conference on artificial intelligence*. vol. 33, pp. 3159–3166 (2019)
2. Alain, G., Bengio, Y.: What regularized auto-encoders learn from the data-generating distribution. *The Journal of Machine Learning Research* **15**(1), 3563–3593 (2014)
3. Allan, M., Williams, C.: Harmonising chorales by probabilistic inference. *Advances in neural information processing systems* **17** (2004)
4. Austin, J., Johnson, D.D., Ho, J., Tarlow, D., van den Berg, R.: Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems* **34**, 17981–17993 (2021)
5. Bachman, P., Precup, D.: Data generation as sequential decision making. *Advances in Neural Information Processing Systems* **28** (2015)
6. Boulanger-Lewandowski, N., Bengio, Y., Vincent, P.: Modeling temporal dependencies in high-dimensional sequences: application to polyphonic music generation and transcription. In: *Proceedings of the 29th International Conference on International Conference on Machine Learning*. pp. 1881–1888 (2012)
7. Caccia, M., Caccia, L., Fedus, W., Larochelle, H., Pineau, J., Charlin, L.: Language gans falling short. In: *International Conference on Learning Representations* (2020)
8. Che, T., Li, Y., Zhang, R., Hjelm, R.D., Li, W., Song, Y., Bengio, Y.: Maximum-likelihood augmented discrete generative adversarial networks. *arXiv preprint arXiv:1702.07983* (2017)
9. Chen, N., Zhang, Y., Zen, H., Weiss, R.J., Norouzi, M., Chan, W.: Wavegrad: Estimating gradients for waveform generation. In: *International Conference on Learning Representations* (2020)
10. Crowson, K., Biderman, S., Kornis, D., Stander, D., Hallahan, E., Castricato, L., Raff, E.: Vqgan-clip: Open domain image generation and editing with natural language guidance. In: *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*. pp. 88–105. Springer (2022)
11. Dai, Z., Yang, Z., Yang, Y., Carbonell, J.G., Le, Q., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. pp. 2978–2988 (2019)
12. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: *Proceedings of NAACL-HLT*. pp. 4171–4186 (2019)
13. Dinh, L., Krueger, D., Bengio, Y.: Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516* (2014)
14. Dinh, L., Sohl-Dickstein, J., Bengio, S.: Density estimation using real nvp. In: *International Conference on Learning Representations* (2017)
15. Donahue, C., Lee, M., Liang, P.: Enabling language models to fill in the blanks. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. pp. 2492–2501 (2020)
16. Esser, P., Rombach, R., Ommer, B.: Taming transformers for high-resolution image synthesis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 12873–12883 (June 2021)

17. Fang, A., Liu, A., Seetharaman, P., Pardo, B.: Bach or mock? a grading function for chorales in the style of js bach. arXiv preprint arXiv:2006.13329 (2020)
18. Fedus, W., Goodfellow, I., Dai, A.M.: Maskgan: Better text generation via filling in the `_`. In: International Conference on Learning Representations (2018)
19. Ganin, Y., Kulkarni, T., Babuschkin, I., Eslami, S.A., Vinyals, O.: Synthesizing programs for images using reinforced adversarial learning. In: International Conference on Machine Learning. pp. 1666–1675. PMLR (2018)
20. Germain, M., Gregor, K., Murray, I., Larochelle, H.: Made: Masked autoencoder for distribution estimation. In: International Conference on Machine Learning. pp. 881–889. PMLR (2015)
21. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. *Advances in neural information processing systems* **27** (2014)
22. Gregor, K., Danihelka, I., Graves, A., Rezende, D., Wierstra, D.: Draw: A recurrent neural network for image generation. In: International Conference on Machine Learning. pp. 1462–1471. PMLR (2015)
23. Gu, J., Bradbury, J., Xiong, C., Li, V.O., Socher, R.: Non-autoregressive neural machine translation. In: International Conference on Learning Representations (2018)
24. Gu, S., Chen, D., Bao, J., Wen, F., Zhang, B., Chen, D., Yuan, L., Guo, B.: Vector quantized diffusion model for text-to-image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10696–10706 (2022)
25. Guo, J., Xu, L., Chen, E.: Jointly masked sequence-to-sequence model for non-autoregressive neural machine translation. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 376–385 (2020)
26. Gutmann, M., Hyvärinen, A.: Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In: Proceedings of the thirteenth international conference on artificial intelligence and statistics. pp. 297–304. JMLR Workshop and Conference Proceedings (2010)
27. Hadjeres, G., Crestel, L.: Vector quantized contrastive predictive coding for template-based music generation. arXiv preprint arXiv:2004.10120 (2020)
28. He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16000–16009 (2022)
29. Hill, F., Cho, K., Korhonen, A.: Learning distributed representations of sentences from unlabelled data. In: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 1367–1377 (2016)
30. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems* **33**, 6840–6851 (2020)
31. Hoogeboom, E., Gritsenko, A.A., Bastings, J., Poole, B., van den Berg, R., Salimans, T.: Autoregressive diffusion models. In: International Conference on Learning Representations (2022)
32. Hoogeboom, E., Nielsen, D., Jaini, P., Forré, P., Welling, M.: Argmax flows and multinomial diffusion: Learning categorical distributions. *Advances in Neural Information Processing Systems* **34** (2021)
33. Huang, C.Z.A., Cooijmans, T., Roberts, A., Courville, A.C., Eck, D.: Counterpoint by convolution. In: Proceedings of the 18th International Society for Music Information Retrieval Conference. pp. 211–218. ISMIR, Suzhou, China (Oct 2017)

34. Huang, C.Z.A., Vaswani, A., Uszkoreit, J., Simon, I., Hawthorne, C., Shazeer, N., Dai, A.M., Hoffman, M.D., Dinculescu, M., Eck, D.: Music transformer: Generating music with long-term structure. In: International Conference on Learning Representations (2018)
35. Huang, C.W., Lim, J.H., Courville, A.C.: A variational perspective on diffusion-based generative models and score matching. *Advances in Neural Information Processing Systems* **34** (2021)
36. Kalchbrenner, N., Danihelka, I., Graves, A.: Grid long short-term memory. *arXiv preprint arXiv:1507.01526* (2015)
37. Kingma, D.P., Salimans, T., Poole, B., Ho, J.: Variational diffusion models. In: *Advances in Neural Information Processing Systems* (2021)
38. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: International Conference on Learning Representations (2014)
39. Kingma, D.P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., Welling, M.: Improved variational inference with inverse autoregressive flow. *Advances in neural information processing systems* **29** (2016)
40. Lam, M.W., Wang, J., Su, D., Yu, D.: Bddm: Bilateral denoising diffusion models for fast and high-quality speech synthesis. In: International Conference on Learning Representations (2022)
41. Lee, J., Mansimov, E., Cho, K.: Deterministic non-autoregressive neural sequence modeling by iterative refinement. In: EMNLP (2018)
42. Liang, F.: Bachbot: Automatic composition in the style of bach chorales. *University of Cambridge* **8**, 19–48 (2016)
43. Lin, K., Li, D., He, X., Zhang, Z., Sun, M.T.: Adversarial ranking for language generation. *Advances in neural information processing systems* **30** (2017)
44. Liu, A., Fang, A., Hadjeres, G., Seetharaman, P., Pardo, B.: Incorporating music knowledge in continual dataset augmentation for music generation. *arXiv preprint arXiv:2006.13331* (2020)
45. Lu, S., Zhu, Y., Zhang, W., Wang, J., Yu, Y.: Neural text generation: Past, present and beyond. *arXiv preprint arXiv:1803.07133* (2018)
46. de Masson d’Autume, C., Mohamed, S., Rosca, M., Rae, J.: Training language gans from scratch. *Advances in Neural Information Processing Systems* **32** (2019)
47. Meister, C., Pimentel, T., Wiher, G., Cotterell, R.: Typical decoding for natural language generation. *arXiv preprint arXiv:2202.00666* (2022)
48. Meng, C., He, Y., Song, Y., Song, J., Wu, J., Zhu, J.Y., Ermon, S.: Sdedit: Guided image synthesis and editing with stochastic differential equations. In: International Conference on Learning Representations (2021)
49. Mittal, G., Engel, J., Hawthorne, C., Simon, I.: Symbolic music generation with diffusion models. *arXiv preprint arXiv:2103.16091* (2021)
50. Nichol, A.Q., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., Chen, M.: Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In: International Conference on Machine Learning. pp. 16784–16804. PMLR (2022)
51. Papamakarios, G., Pavlakou, T., Murray, I.: Masked autoregressive flow for density estimation. *Advances in neural information processing systems* **30** (2017)
52. Pati, A., Lerch, A., Hadjeres, G.: Learning to traverse latent spaces for musical score inpainting. In: Flexer, A., Peeters, G., Urbano, J., Volk, A. (eds.) *Proceedings of the 20th International Society for Music Information Retrieval Conference, ISMIR 2019, Delft, The Netherlands, November 4-8, 2019*. pp. 343–351 (2019)
53. Payne, C.: Musenet. openai.com/blog/musenet (2019)

54. Radford, A., Wu, J.: Rewon child, david luan, dario amodei, and ilya sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog* **1**(8), 9 (2019)
55. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents (2022)
56. Rezende, D., Mohamed, S.: Variational inference with normalizing flows. In: *International conference on machine learning*. pp. 1530–1538. PMLR (2015)
57. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models (2021)
58. Savinov, N., Chung, J., Binkowski, M., Elsen, E., van den Oord, A.: Step-unrolled denoising autoencoders for text generation. In: *International Conference on Learning Representations* (2022)
59. Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., Ganguli, S.: Deep unsupervised learning using nonequilibrium thermodynamics. In: *International Conference on Machine Learning*. pp. 2256–2265. PMLR (2015)
60. Song, Y., Durkan, C., Murray, I., Ermon, S.: Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems* **34** (2021)
61. Song, Y., Ermon, S.: Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems* **32** (2019)
62. Theis, L., Bethge, M.: Generative image modeling using spatial lstms. *Advances in neural information processing systems* **28** (2015)
63. Uribe, B., Côté, M.A., Gregor, K., Murray, I., Larochelle, H.: Neural autoregressive distribution estimation. *The Journal of Machine Learning Research* **17**(1), 7184–7220 (2016)
64. Uribe, B., Murray, I., Larochelle, H.: A deep and tractable density estimator. In: *International Conference on Machine Learning*. pp. 467–475. PMLR (2014)
65. Van Oord, A., Kalchbrenner, N., Kavukcuoglu, K.: Pixel recurrent neural networks. In: *International conference on machine learning*. pp. 1747–1756. PMLR (2016)
66. Vasquez, S., Lewis, M.: Melnet: A generative model for audio in the frequency domain. *arXiv preprint arXiv:1906.01083* (2019)
67. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
68. Yang, Z., Dai, Z., Yang, Y., Carbonell, J., Salakhutdinov, R.R., Le, Q.V.: Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems* **32** (2019)
69. Yu, L., Zhang, W., Wang, J., Yu, Y.: Seqgan: Sequence generative adversarial nets with policy gradient. In: *Proceedings of the AAAI conference on artificial intelligence*. vol. 31 (2017)
70. Zhu, Y., Lu, S., Zheng, L., Guo, J., Zhang, W., Wang, J., Yu, Y.: Texus: A benchmarking platform for text generation models. In: *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. pp. 1097–1100 (2018)