# DINOISER: Diffused Conditional Sequence Learning by Manipulating Noises

**Jiasheng Ye**[*♡◇]  **Zaixiang Zheng**[†♡]  **Yu Bao**[♡]  **Lihua Qian**[♡] and **Mingxuan Wang**[♡]

[♡]ByteDance Research  [◇]Fudan University

jsye23@m.fudan.edu.cn, zhengzaixiang@bytedance.com

{baoyu.001, qianlihua, wangmingxuan.89}@bytedance.com

https://github.com/yegcjs/DINOISER

## Abstract

While diffusion models have achieved great success in generating continuous signals such as images and audio, it remains elusive for them to learn discrete sequence data like natural languages. Although recent advances circumvent this challenge of discreteness by embedding discrete tokens as continuous surrogates, they still fall short of satisfactory generation quality. To understand this, we first dive deep into the denoised training protocol of diffusion-based sequence generative models and determine their three severe problems: (1) failing to learn; (2) lack of scalability; and (3) neglecting source conditions. We argue that these problems can be boiled down to the *pitfall of the not completely eliminated discreteness* in the embedding space, and the *scale of noises* is decisive herein. In this paper, we introduce DINOISER to facilitate diffusion models for sequence generation by manipulating noises. We propose to adaptively determine the range of sampled noise scales during training; and encourage the proposed diffused sequence learner to leverage source conditions with amplified noise scales during inference. Experiments show that DINOISER enables consistent improvement over the baselines of previous diffusion sequence generative models on several conditional sequence modeling benchmarks thanks to both effective training and inference strategies. Analyses further verify that DINOISER can make better use of source conditions to govern its generative process.

## 1 Introduction

Conditional sequence learning aims at generating a target sequence from given conditions, which is one of the important paradigms of natural language generation (Sutskever et al., 2014; Wiseman et al., 2017; Raffel et al., 2020), including machine translation (Bahdanau et al., 2014), summarization (Rush et al., 2015), and paraphrasing (Prakash et al., 2016). Recent advances in generative modeling introduce diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020b), which achieve great success in generating continuous signals, including images (Rombach et al., 2021), video (Ho et al., 2022), and audio (Kong et al., 2020). With promising characteristics such as diversity and controllability demonstrated in these domains, diffusion models also garner growing interest for sequence learning in the research community (Li et al., 2022), which further gives the promise to a unified generative modeling paradigm across different modalities (Bao et al., 2023),

However, the discrete nature of sequence data, constituted by a number of tokens in order, makes it non-trivial to apply diffusion models for conditional sequence learning. Typical diffusion models noise data with Gaussian permutation kernels (Ho et al., 2020) and learn to recover original data from their corrupted versions, which is not directly compatible with discrete tokens. To remedy this, DiffusionLM (Li et al., 2022) attempted to embed discrete tokens into continuous space and employ diffusion models to the embedding space. Although this kind of approach unlocks the possibility of applying diffusion models to discrete data, it still falls short of competitive performance for various conditional sequence generation tasks (Fig. 1A).

We argue that embedding discrete tokens into continuous surrogates does not necessarily eliminate discreteness completely. To verify this, we conduct in-depth preliminary studies and highlight our findings along with their implications as follows. **(1)** On the *pitfall of discreteness*. Embeddings populate only finite clusters (up to the vocabulary size) in the continuous space, which results in the vastness of low-density regions especially when the models are learned with small-scale noises. We

---

refer to this as the pitfall of discreteness, which suggests that small noises hinder conditional sequence learning, and thus should be avoided during training. **(2)** On *scalability*. It becomes increasingly harder for the diffusion process to eliminate discreteness when the dimension of the embedding space gets scaled up, suggesting that to ensure scalability, an adaptable noise schedule is necessitated yet neglected. **(3)** On *conditional learning*. Enlarging noises in inference can calibrate diffusion models to take into account more source conditional information. Please refer to §3 for more details.

Motivated by these findings, we propose DI-NOISER to improve **di**ffusion models by manipulating **noi**ses for conditional **se**quence lea**r**ning. We propose a novel training strategy to eliminate training on small noise scales to avoid their negative influences, for which we introduce the noise scale clipping strategy to adaptively manipulate the noise scales. For inference, we propose to manipulate the model to be exposed to larger noise scales to encourage trained diffusion models to leverage source conditions.

We summarize our contributions as well as our findings as follows:

- By thorough and in-depth preliminary studies, we shed light on the pitfall of discreteness along with the critical role of noise scales in conditional sequence learning with diffusion models, thereby suggesting meliorated solutions in terms of both training and inference by manipulating noises.

- We accordingly propose DINOISER to leverage large noise scales in both training and inference. Experiments show that DINOISER achieves strong performance on a variety of conditional sequence learning tasks, paving way for featuring diffusion models for various conditional sequence learning tasks. Our experiments comprehensively include several machine translation benchmarks (both bilingual and multilingual), as well as text simplification and paraphrasing, ranging from low-resource to high-resource scenarios.

- Ablations show that both DINOISER 's improved training and inference approaches result in considerable performance gains. Further analysis verifies that our proposed posthoc inference strategy, i.e., the condition enhanced denoiser, can help make better use of source conditions for accurate predictions.

## 2   Background

**Conditional Sequence Learning.** Conditional sequence learning aims to yield target sequence $\mathbf{y} = [y_1, y_2, \ldots, y_n] \in \{0, 1\}^{n \times |\mathcal{V}|}$ within the vocabulary space $\mathcal{V}$, given source conditions $\mathbf{x}$, which can be another sequence $\mathbf{x} = [x_1, x_2, \ldots, x_m]$. The conventional modeling paradigm (Sutskever et al., 2014; Vaswani et al., 2017) generates target tokens in an autoregressive decomposition $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{n} p(y_i|\mathbf{y}_{<i}, \mathbf{x})$. Gu et al. (2018) proposed an alternative way in a fully non-autoregressive (NAR) manner, where all the tokens are predicted in parallel by assuming conditional independence between the target tokens, *i.e.*, $p(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{n} p(y_i|\mathbf{x})$. Later works alleviate this strong assumption by iterative refinement (Lee et al., 2018; Ghazvininejad et al., 2019; Gu et al., 2019), resulting in improved generation quality. These iterative refinement approaches generate target sequences with several cycles, in each of which the models generate sequence depending on both the source sequence and the intermediate prediction of the previous one, *i.e.*, $p(\mathbf{y}|\mathbf{x}) = \prod_{t=1}^{T} p(\mathbf{y}^{(t)}|\mathbf{y}^{(t-1)}, \mathbf{x})$.

**Diffusion Probabilistic Models.** Given a random variable $\mathbf{z}_0$ from an underlying data distribution $q(\mathbf{z}_0)$, diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020) define a forward diffusion process $\{\mathbf{z}_t\}_{t \in [0,1]}$ perturbed with a Gaussian perturbation kernel, starting with $\mathbf{z}_0$ and converging to its corresponding stationary distribution, such that for any $t \in [0, 1]$, the distribution of $\mathbf{z}_t$ given $\mathbf{z}_0$ satisfies $q(\mathbf{z}_t|\mathbf{z}_0) = \mathcal{N}(\mathbf{z}_t; \alpha(t)\mathbf{z}_0, \sigma^2(t)\mathbf{I})$, or with Gaussian reparameterization (Kingma and Welling, 2013; Ho et al., 2020),

$$\mathbf{z}_t = \alpha(t)\mathbf{z}_0 + \sigma(t)\epsilon_t, \ \ \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where $\sigma(t)$ is a monotonically increasing function, usually referred to as the noise schedule, satisfying $\sigma(0) = 0$ and $\sigma(1) \approx 1$; and $\alpha(t) = \sqrt{1 - \sigma^2(t)}$. The noise schedule $\sigma(t)$ controls the degree of corruption at different timestep $t$. As $t$ gets larger, the noise scale $\sigma(t)$ gets larger whereas $\alpha(t)$ gets smaller, hence the more corrupted data $\mathbf{z}_t$ from the original $\mathbf{z}_0$. At $t = 1$, with $\alpha(1) \approx 0$ and $\sigma(1) \approx 1$, $\mathbf{z}_t$ become pure noises as reaching the stationary distribution of a standard Gaussian.

Song et al. (2020b) proves that such a Gaussian diffusion process has the same transition distribution $q(\mathbf{z}_t|\mathbf{z}_0)$ as the stochastic differential equa-

tion (SDE): $\mathrm{d}\mathbf{z} = -\frac{1}{2}\beta(t)\mathbf{z}\mathrm{d}t + \sqrt{\beta(t)}\mathrm{d}\omega$ where $\beta(t) = -2\frac{\mathrm{d}\log\alpha(t)}{\mathrm{d}t}$; $\omega$ denotes the standard Wiener process. As such, the corresponding generative process can be achieved as its time reversal by solving the following ordinary differential equation (diffusion ODE):

$$
\begin{aligned}
\mathrm{d}\mathbf{z} &= \left[ -\frac{1}{2}\beta(t)\mathbf{z} + \frac{1}{2}\beta(t)\frac{\epsilon_t}{\sigma(t)} \right] \mathrm{d}t \\
&= \left[ -\frac{1}{2}\beta(t)\mathbf{z} + \frac{\beta(t)}{2\sigma^2(t)}\left(\mathbf{z} - \alpha(t)\mathbf{z}_0\right) \right] \mathrm{d}t.
\end{aligned} \quad (2)
$$

In practice, we can then use a learned model $\mathbf{z}_\theta(\mathbf{z}_t, t)$ to estimate $\mathbf{z}_0$ and plug into Eqn. 2, which can be learned by minimizing the discrepancy between training data and model estimation (Ho et al., 2020; Song et al., 2020b; Ramesh et al., 2022):

$$
\mathcal{L}_{\text{diffusion}}(\mathbf{z}_0) = \mathop{\mathbb{E}}_{\substack{t\sim\mathcal{U}(0,1)\\ \epsilon_t\sim\mathcal{N}(\mathbf{0},\mathbf{I})}} \left[ \|\mathbf{z}_\theta(\mathbf{z}_t, t) - \mathbf{z}_0\|_2^2 \right]. \quad (3)
$$

Given Eqn. 2 with a trained model $\mathbf{z}_\theta(\mathbf{z}_t, t)$, we can use arbitrary ODE solvers to solve this diffusion ODE from $t = 1$ to $t = 0$ for sampling data. An effective and efficient solver to this end is the DDIM solver (Song et al., 2020a; Lu et al., 2022) and is widely adopted. It discretizes the ODE into $M + 1$ timesteps $\{t_i\}_{i=0}^M$ decreasing from $t_0 = 1$ to $t_M \approx 0$. Then, it samples $\mathbf{z}_{t_0}$ from the standard Gaussian distribution and computes $\{\mathbf{z}_{t_i}\}_{i=1}^M$ with $M$ iterations, in each of which $\mathbf{z}_{t_i}$ is predicted from $\mathbf{z}_{t_{i-1}}$ according to

$$
\mathbf{z}_{t_i} = \alpha(t_i)\mathbf{z}_\theta(\mathbf{z}_{t_{i-1}}, t_{i-1}) + \sigma(t_i)\epsilon_\theta(\mathbf{z}_{t_{i-1}}, t_{i-1}), \quad (4)
$$

where $\epsilon_\theta(\mathbf{z}_{t_{i-1}}, t_{i-1})$ is the predicted noise, which can be directly induced according to Eqn. 1,

$$
\epsilon_\theta(\mathbf{z}_{t_{i-1}}, t_{i-1}) = \frac{\mathbf{z}_{t_{i-1}} - \alpha(t_{i-1})\mathbf{z}_\theta(\mathbf{z}_{t_{i-1}}, t_{i-1})}{\sigma(t_{i-1})}. \quad (5)
$$

After iterations, the last prediction $\mathbf{z}_{t_M}$ is taken as the final generated result $\hat{\mathbf{z}}_0$ of the sampling.

**Diffusion Models for Conditional Sequence Learning.** The denoising process of diffusion models matches an iterative refinement process (Gong et al., 2022). However, diffusion models are not directly applicable to sequence learning tasks since the original diffusion models operate in continuous space rather than sequences of discrete tokens. DiffusionLM (Li et al., 2022) tackles this by embedding the discrete tokens into continuous latent space and applying diffusion models

therein. We can then train the models as variational autoencoders (Kingma and Welling, 2013), where a diffusion model serves as the prior, from a latent-variable model perspective, and derive the corresponding variational lower bound (Wehenkel and Louppe, 2021; Vahdat et al., 2021):

$$
\mathcal{L}(\mathbf{y}) = \mathbb{E}_{\mathbf{z}_0}\Big[ \underbrace{-\log p_\theta(\mathbf{y}|\mathbf{z}_0)}_{\mathcal{L}_{\text{reconstruction}}} + \mathcal{L}_{\text{diffusion}}(\mathbf{z}_0) \Big], \quad (6)
$$

where $\mathbf{y}$ is the original sequence with $\mathbf{z}_0$ as its embeddings[1]; $\mathcal{L}_{\text{diffusion}}$ denote the diffusion loss (Eqn. 3) which now operates on the embedding space, and $\mathcal{L}_{\text{reconstruction}}$ is the newly added reconstruction term.

To further adapt the model for conditional sequence generation, a vanilla approach is to replace the unconditional model $\mathbf{z}_\theta(\mathbf{z}_t, t)$ with a conditional model $\mathbf{z}_\theta(\mathbf{z}_t, \mathbf{x}, t)$, where $\mathbf{x}$ is the source condition. Similar to the previous practice of using diffusion models for conditional generation in vision (Rombach et al., 2021), the diffusion process can be kept unchanged, the same as Eqn. 1. And the length of the target sequences is decided by predicting the length difference between the source and the target.

## 3 The Pitfall of Discreteness: The Noise Scale Matters

In this section, we dive deep into the current weaknesses of diffusion models for conditional sequence learning and find that the noise scale matters, which accordingly motivates our proposal for improved training and inference.

**Settings.** We begin with the vanilla conditional diffusion model modified from DiffusionLM (Li et al., 2022) as described in §2. We follow the original paper of DiffusionLM to apply the sqrt schedule (i.e., $\sigma(t) = t^{0.25}$) to arrange noise scales for training and sampling. We use IWSLT14 DE→EN (Cettolo et al., 2012) machine translation benchmark for evaluation. We also include CMLM (Ghazvininejad et al., 2019) as a baseline for comparison, which is a strong conditional sequence generative model that generates sequence by iterative refinement similar to diffusion models but in discrete tokens space.

---

[1]DiffusionLM adds tiny noise to the embeddings to form $\mathbf{z}_0$ (i.e. $\mathbf{z}_0 \sim \mathcal{N}(\text{EMB}(\mathbf{y}), \sigma_0\mathbf{I})$). We empirically find this unnecessary and letting $\mathbf{z}_0$ follow a Dirac distribution makes training more efficient.
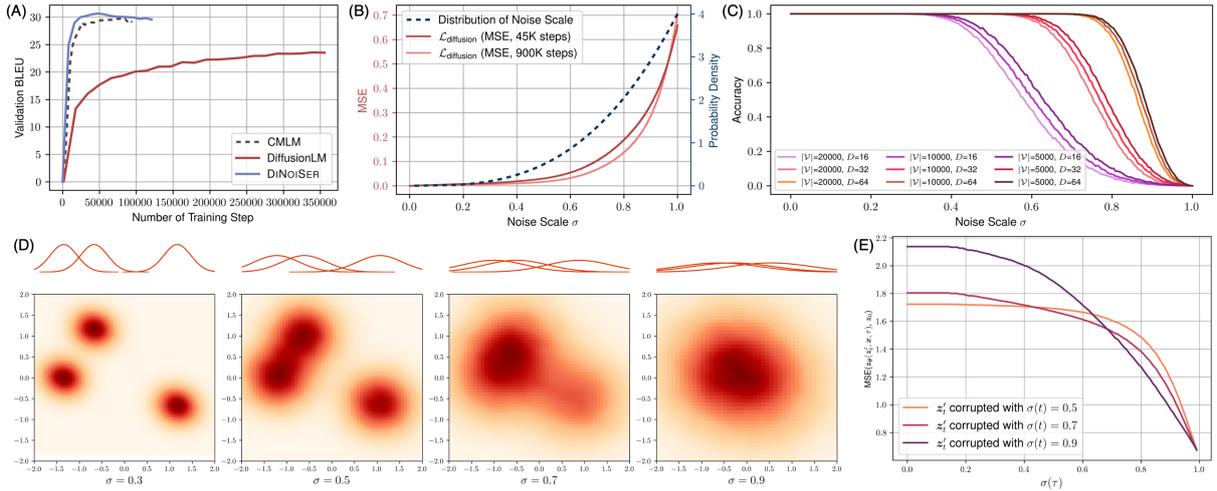
Figure 1: *Preliminary study.* **(A)** The validation BLEU of different models on IWSLT14 DE→EN at different training steps. **(B)** Diffusion loss of DiffusionLM on the validation set of IWSLT14 DE→EN at different noise scales and the distribution of noise scale sampled during training. **(C)** The accuracy of predicting $\mathbf{z}_0$ from $\mathbf{z}_t$ by finding the nearest neighbor for $\mathbf{z}_t$ with different noise scales, vocabulary sizes $|\mathcal{V}|$, and dimensions $D$. **(D)** An illustrative example of the distributions of $\mathbf{z}_t$ of three data points corrupted with different noise scales as in Eqn. 1, where for small noise scales, a large proportion of the embedding space between modes (associated with tokens) remains vacant. **(E)** The tendency of whether the model prediction is more influenced by the source or target side information when fed with timestep correspond to different noise scales. In addition to the source condition $\mathbf{x}$, we feed the model with $\mathbf{z}'_t = \mathbf{z}_t(\mathbf{y}', t)$ that is corrupted with a timestep-dependent noise $\sigma(t)$ from a negative $\mathbf{y}'$, which is different from the original (positive sample of) target sequence $\mathbf{y}$. We compare the similarity between the model prediction $\mathbf{z}_\theta(\mathbf{z}'_t, \mathbf{x}, \tau)$ to the embedding of ground-truth $\mathbf{z}_0(\mathbf{y})$, and study to what degree the model prediction is governed by the source condition $\mathbf{x}$ (via the embedding of the ground-truth $\mathbf{z}_0(\mathbf{y})$ as the proxy), or the target information $\mathbf{y}'$ (via $\mathbf{z}'_t$) otherwise.

**Observations.**   Here are our findings.

O1. ***DiffusionLM still falls short of conditional sequence learning.*** Fig. 1(A) shows the validation performance of the two models at different training steps, in which the performance of DiffusionLM still lags behind CMLM by a large margin, even taking many more steps before convergence. This shows that the performance and training efficiency of the vanilla diffusion-based sequence learner remain unsatisfactory.

O2. ***Diffusion losses at small noise scales are unexpectedly small.*** DiffusionLM uniformly samples timesteps hence the corresponding noise scales during training. As shown in Fig. 1(B), we find that the magnitudes of diffusion losses approach almost zero for small noise scales, indicating that it is quite trivial to recover the corrupted embeddings under such circumstances. We conjecture that, combined with the illustrated example in Fig. 1(D), this is because there remain highly discrete modes for the embedding density such that any corrupted embedding is very likely to lie in a

region with a small radius around the original token embedding. As a consequence, the more the modes of embeddings separate from each other the smaller the diffusion loss, which adheres to the following observation.

O3. ***It becomes increasingly harder for the diffusion process to eliminate discreteness while the dimension of the embedding space scales up.*** Fig. 1(C) shows a surprisingly high accuracy of recovering corrupted embeddings that can be easily achieved by simply seeking the nearest neighbor when embedding dimensions enlarge, even at considerably large noise scales. This reveals that scaling embedding space leads to more severe discreteness, namely a curse of dimensionality.

O4. ***On condition learning: larger noise scales calibrate diffusion models in taking into account more source conditional information during inference.*** We have seen that recovering embeddings corrupted with small noise scales is easy (O2), and if modes distribute separately, even finding the nearest neighbor is enough (O3). In Fig. 1(C), as the noise scale

Table 1: Illustration of hallucinations of vanilla DiffusionLM on IWSLT14 DE→EN translation task, along with DINOISER's predictions, where the vanilla DiffusionLM generates inexplicable expressions that are irrelevant to the source condition whose meaning dramatically differs from the groud-truth target.

| | |
|---|---|
| **Source** | Mit welchen worten würden sie ban beschreiben? |
| **Reference** | What are the words you would use to describe ban? |
| **DiffusionLM** | In which words would you save ban? |
| **DINOISER** | In which words would you describe ban? |

decreases, the prediction accuracy by finding the nearest neighbor increases and achieves almost 100% under a threshold, which can be learned trivially even with little to no source conditions. This results in the hallucination as shown in Tab. 1, which is an unexpected consequence for conditional sequence generative models to yield output loyal to the input condition. To mitigate this, we quantitatively study the influence of noise scales on conditional reliance. As shown in Fig 1(E), we find that as the noise scales are larger, the model can predict more faithfully to source conditions.

**Concluding remarks.** We summarize conclusions from the aforementioned observations along with suggestions for more plausible diffused conditional sequence learning:

C1. ***We should not train on too small noise scales to circumvent the pitfall of discreteness.*** Both O2 and O4 show the negative influences of small noise scales on training that it leads to a not smooth embedding space with vast regions of low density between modes associated with tokens (O2). These regions can inevitably be sampled during inference[2], thereby giving rise to error accumulation. Besides, it also impedes conditional learning (O4). To remedy this, probably a simple way is to eliminate the chance of training with small noise scales.

C2. ***We need to determine the noise schedule according to the dimensionality of the embedding space.*** Fitting more complex datasets usually requires larger embedding dimensions. O3 indicates the criterion to distinguish large

___
[2] Consider that a token $\alpha$ is translated into token $A$ or $a$ with 50% each. Without extra information, the optimal prediction for translating $\alpha$ is the center of the embedding of $A$ and $a$. This is because minimizing its training objective (*i.e.*, $\frac{1}{2}\|\mathbf{z}_\theta - \mathbf{z}_A(0)\|_2^2 + \frac{1}{2}\|\mathbf{z}_\theta - \mathbf{z}_a(0)\|_2^2$) results in $\mathbf{z}_\theta = \frac{\mathbf{z}_A(0)+\mathbf{z}_a(0)}{2}$. The prediction exactly falls in the blank area that lies between embeddings.

and small noise scales depends on the embeddings hence the complexity of the datasets. However, existing methods employ a fixed noise schedule for all embedding dimensions, which lacks scalability. This, therefore, demands a task-specific noise schedule to accommodate diverse datasets.

C3. ***We could expose the model to larger noise scales for better source conditions leverage.*** O4 suggests that the more corrupted the embeddings, the more difficult for the model to recover, thereby necessitating more reliance on source conditions. Accordingly, we may encourage trained diffusion models to care more about source conditions for free by *post-hoc* manipulating the noise to large ones.

# 4    DINOISER

Provided the observations and postulates we discussed in §3, we accordingly propose DINOISER, a simple yet effective method that improves **di**ffusion models by manipulating **noi**ses for conditional **se**quence lea**r**ning. The general principle of DINOISER is to determine the best-suited noise scales for both training and inference for conditional sequence generation. In a nutshell, as for training, we propose to eliminate the chance of training diffused sequence learners with small-scale noises so as to circumvent the aforementioned pitfall of discreteness in embedding space (§4.1). As for sampling, we propose a new effective sampler to amplify the impact of source conditions on the model prediction, where timesteps corresponding to large noise scales are always fed into the model (§4.2). We now dive deep into the details of DINOISER.

## 4.1    Noise Scale Clipping: Manipulating Noises for Counter-Discreteness Training

Recall that C1 and C2 in §3 demonstrate that small noises can barely help "discrete" embeddings populate the entire continuous space, and also undermine conditional learning. A simple yet effective way to mitigate this is to encourage training diffusion models with sufficiently large noise scales. As such, we propose *noise scale clipping*, where we bound the minimum noise scale $\sigma_{\min}$ for training such that only timesteps satisfying $\sigma(t) \geq \sigma_{\min}$ could be sampled, which is decided adaptively as the model learn progresses.

To start with, we can eliminate the scaling effect of $\alpha(t)$ in the forward diffusion process for each
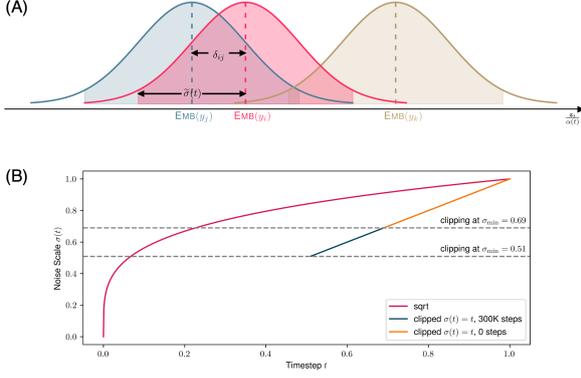
Figure 2: **(A)** Illustration of the proposed noise scale clipping. To remedy the pitfall of discreteness, we propose to ensure a sufficiently large minimum "overlap" between corrupted embeddings. As shown in this example, such a goal of counter-discreteness is achieved by bounding the standard deviation $\tilde{\sigma}(t)$ of $\text{EMB}(y_i)$ by $\delta_{ij}$ the "distance" to its nearest neighbor (*i.e.*, $y_j$). **(B)** Comparison between `sqrt` noise schedule and our noise schedule $\sigma(t) = t$ manipulated with the proposed noise scale clipping.

token embedding by rewriting Eqn. 1 into:

$$
\frac{\mathbf{z}_t[i]}{\alpha(t)} = \mathbf{z}_0 + \frac{\sigma(t)}{\alpha(t)}\epsilon_t \Rightarrow \frac{\mathbf{z}_t[i]}{\alpha(t)} \sim \mathcal{N}\left(\mathbf{z}_0[i], \frac{\sigma^2(t)}{1-\sigma^2(t)}\mathbf{I}\right)
$$
$$
\Rightarrow \frac{\mathbf{z}_t[i]}{\alpha(t)} \sim \mathcal{N}\left(\text{EMB}(y_i), \tilde{\sigma}^2(t)\mathbf{I}\right) \quad (7)
$$

As illustrated in Fig. 2(A), there, intuitively, should exist a sufficiently large number $\delta^*$ measuring the minimum "overlap" between the distributions of two corrupted embeddings under the Gaussian perturbation kernel with a standard deviation of $\tilde{\sigma}(t) = \frac{\sigma(t)}{\sqrt{1-\sigma^2(t)}}$. To this end, we let $\delta^2$ be the minimum amount of variation of added noise, defined as the average squared L2-distances between the embeddings and their nearest neighbor, normalized by the dimension of embeddings (according to C2 in §3):

$$
(\delta^*)^2 = \frac{1}{|\mathcal{V}|}\sum_{i=1}^{|\mathcal{V}|} \min_{1\le j\ne i\le|\mathcal{V}|} \delta_{ij}^2
$$
$$
= \frac{1}{|\mathcal{V}|}\sum_{i=1}^{|\mathcal{V}|} \min_{1\le j\ne i\le|\mathcal{V}|} \frac{1}{D}\|\text{EMB}(y_i) - \text{EMB}(y_j)\|_2^2. \quad (8)
$$

We now define the noise scale clipping[3] as follows:

---

[3]**Our goal can be motivated through the lens of optimal transport**. That it to say, we aim to determine the minimum cumulative cost $\delta^2 = \sum_{i=1}^{L} \mathbf{T}_{ij}\delta_{ij}^2$, where $L$ is the sequence length, by finding the optimal transportation $\mathbf{T}$ of moving the perturbed embeddings at timestep $t$, *i.e.* $\frac{\mathbf{z}_t[i]}{\alpha(t)} \sim \mathcal{N}\left(\text{EMB}(y_i), \tilde{\sigma}^2(t)\mathbf{I}\right), \quad i \in \mathcal{V}$, such that the corrupted embedding $\frac{\mathbf{z}_t[i]}{\alpha(t)}$, if gets noised by a Gaussian deviation

**Definition** (The noise scale clipping). *Let $\mathcal{V}$ be the target vocabulary with corresponding embeddings in $D$-dimensional space $\forall y_i \in \mathcal{V} : \text{EMB}(y_i) \in \mathbb{R}^D$, the noise scale clipping is performed so that the noise scale $\sigma(t)$ always satisfies:*

$$
\tilde{\sigma}^2(t) = \frac{\sigma^2(t)}{1-\sigma^2(t)} \ge (\delta^*)^2 \Rightarrow \frac{\sigma_{\min}^2}{1-\sigma_{\min}^2} = (\delta^*)^2, \quad (9)
$$

*the clipping threshold $\sigma_{\min}$ is whereby derived when the equality in Eqn. 9 holds, such that*

$$
\sigma_{\min} = \left(\frac{|\mathcal{V}|\cdot D}{\sum_{i=1}^{|\mathcal{V}|} \min\limits_{1\le j\ne i\le|\mathcal{V}|} \|\text{EMB}(y_i)-\text{EMB}(y_j)\|_2^2} + 1\right)^{-\frac{1}{2}}, \quad (10)
$$

*which is obtained by substituting Eqn. 8 into the R.H.S of Eqn. 9.*

---

satisfying $\tilde{\sigma}(t) >= \delta^*$, cannot be discriminated from those originate from different embeddings, otherwise a smaller $\tilde{\sigma}(t)$ will lead to trivial reconstruction to the original one as a consequence of the pitfall of discreteness (O2 & O3 in §3). As a result, $\delta^*$ serves as a minimum clipping threshold of noise scale for effective training of sequence diffusion models. This can closely relate to the minimum Word Mover Distance, a Wasserstein metric introduced in Kusner et al. (2015):

$$
\delta^2 = \min_{\mathbf{T}} \sum_{i=1}^{L} \mathbf{T}_{ij}\delta_{ij}^2, \quad \text{s.t.} \sum_j \mathbf{T}_{ij} = d_i,
$$

where $\mathbf{T} \in \mathbb{R}^{L\times L}$ is a (sparse) stochastic matrix, where $\mathbf{T}_{ij}$ denotes *how much* of a word $i$ travels to word $j$, subject to the flow consistency equality $\sum_j \mathbf{T}_{ij} = d_i$, with $d_i$ representing the "amount" of a word $i$ appearing in token embedding space (we treat $d_i = 1$). According to the Eqn. (2) in Kusner et al. (2015), under mild conditions, the optimal solution $\mathbf{T}^*$ is for each token $i$ to move all its probability mass to the most similar token $j$ *w.r.t.* a certain measure of their embedding distances, $\delta_{ij} = \|\frac{\mathbf{z}_t[i]}{\alpha(t)} - \text{EMB}(y_j)\|_2$:

$$
\mathbf{T}_{ij}^* = \begin{cases} d_i & \text{if } j = \arg\min_{1\le j\ne i\le L} \delta_{ij}^2 \\ 0 & \text{otherwise} \end{cases}.
$$

As a result, the final minimum transportation cost becomes:

$$
(\delta^*)^2 = \sum_{i=1}^{L} \mathbf{T}_{ij}^* \delta_{ij}^2 = \sum_{i=1}^{L} d_i \cdot (\delta_{ij}^2)^* = \sum_{i=1}^{L} \min_{1\le j\ne i\le L} \delta_{ij}^2
$$
$$
= \sum_{i=1}^{L} \min_{1\le j\ne i\le L}\left[\left\|\frac{\mathbf{z}_t[i]}{\alpha(t)} - \text{EMB}(y_j)\right\|_2\right]^2
$$

A too-small noise scales result in that the nearest neighbors of the corrupted embeddings are exactly their origins, thus

$$
(\delta^*)^2 = \quad = \sum_{i=1}^{L}\left\|\frac{\mathbf{z}_t[i]}{\alpha(t)} - \text{EMB}(y_i)\right\|_2^2 = \sum_{i=1}^{L}(\tilde{\sigma}(t)\epsilon_\mathbf{i})^2,
$$

where $\epsilon_\mathbf{i}$ are standard Gaussian noises. The above results contain no model parameters, indicating that a diffusion model, which learns to minimize the Wasserstein distance between prediction and target distribution (Kwon et al., 2022), can not learn from those mildly perturbed samples. From this perspective, our noise clipping tries to avoid training on these unhelpful samples.

**Algorithm 1** Training with DINOISER

**Input** Training dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})\}$.
**Output** Optimized parameters $\theta$.

1: **repeat**
2:     Sample $\mathbf{x}, \mathbf{y}$ from the dataset $\mathcal{D}$ and embed $\mathbf{y}$ into $\mathbf{z}_0$
3:     $t \sim \mathcal{U}(\sigma^{-1}(\sigma_{\min}), 1)$, where $\sigma_{\min}$ is from Eqn. 10
4:     Sample $\mathbf{z}_t$ with Gaussian reparameterization (Eqn. 1)
5:     Take gradient descent step on
$$\nabla_\theta \left[ -\log p_\theta(\mathbf{y}|\mathbf{z}_0) + \|\mathbf{z}_\theta(\mathbf{z}_t, \mathbf{x}, t) - \mathbf{z}_0\|_2^2 \right]$$
6: **until** converged

As illustrated in Fig. 2(B), the clipping threshold $\sigma_{\min}$ is estimated *adaptively* during training, depending on how properly the model learns the embeddings up to the minimum pair-wise distances within the vocabulary. In each training step, we first estimate the clipping threshold $\sigma_{\min}$ with Eqn. 10, then sample timesteps among $t$ that satisfies $\sigma(t) > \sigma_{\min}$. In practice, one can first estimate the noise scale threshold $\sigma_{\min}$ and then turn it into the timestep threshold $t_{\min} = \sigma^{-1}(\sigma_{\min})$ in general. In this work, we select $\sigma(t) = t$ as the noise scheduler to simplify this procedure[4].

As a result, the updated diffusion loss with an enlarged minimum timestep threshold (thus an increased minimum noise scale) in the final training objective (modified from Eqn. 6) now becomes:

$$\mathcal{L}'_{\text{diffusion}}(\mathbf{y}) = \mathop{\mathbb{E}}_{t \sim \mathcal{U}(t_{\min}, 1), \epsilon_t} \left[ \|\mathbf{z}_\theta(\mathbf{z}_t, \mathbf{x}, t) - \mathbf{z}_0\|_2^2 \right].$$

We provide pseudocodes regarding how to manipulate noises in training as such in Alg. 1.

### 4.2 CEDI: Manipulating Noises for Condition-Enhanced Sampling

Based on C3 in §3, we suppose the model relies more on the source conditions when the input noise scale is large. This implies that we may make the model generate prediction more faithful to source conditions by feeding timesteps corresponding to large noise scales to the model. Fig. 3 shows a synthesis experiment similar to Fig. 1(E), wherein the predictions using a large timestep 0.995 (namely, a larger noise scale) are closer to the embedding

---

[4]This can be done since the effects of different noise schedules are theoretically interchangeable up to different weight factors under the simplified training objective (Ho et al., 2020) we adopted (see Appendix C). We also provide empirical comparisons between different schedules in Tab. 4.
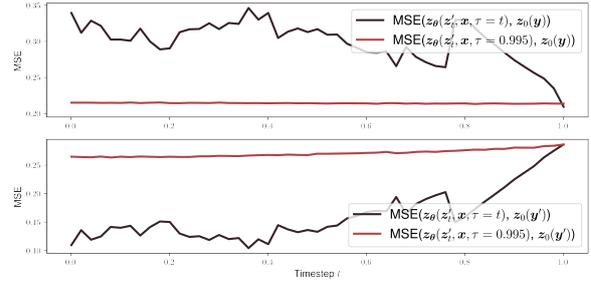


Figure 3: A synthesis experiment where the model is asked to predict with current timestep $\tau = t$ and an alternative larger timestep $\tau = 0.995$, respectively. We compare the MSE between the model prediction $\mathbf{z}_\theta(\mathbf{z}'_t, \mathbf{x}, \tau)$ to the embedding of ground-truth $\mathbf{z}_0(\mathbf{y})$ (top) and negative sample $\mathbf{z}_0(\mathbf{y}')$ (bottom) respectively, and study to which target the model prediction assimilates, the original $\mathbf{y}$ or the negative one $\mathbf{y}'$, hence should most likely be governed by the source or the target information.

of the original target $\mathbf{y}$, while more distant to that of the misleading $\mathbf{y}'$, reiterating that the model relies more on the source condition $\mathbf{x}$ when being exposed to a larger noise due to manipulation in inference.

Accordingly, we propose a **c**ondition-**e**nhanced **d**enoi**s**er (CEDI) for sampling. CEDI always feeds a large $t$ to the model $\mathbf{z}_\theta$ to encourage the model to make use of the source condition. In practice, we largely follow the framework of DDIM solver (Song et al., 2020a) but pick two sets of timesteps. In the first set $\{t_i\}_{i=0}^M$, timesteps decrease uniformly from $t_0 = 1$ to $t_M \approx 0$ as normal. As for the other set $\{\tau_i\}_{i=0}^M$, $\tau_i$s decrease uniformly from $\tau_0 = 1$ to a large time $\tau_M \gg 0$[5]. When making predictions, we assign timesteps from the second set to the model. By replacing corresponding timesteps in the framework of DDIM (Eqn. 4 and Eqn. 5) with $\tau_i$, we generate our predictions by iteratively computing

$$\hat{\mathbf{z}}_{t_i} = \alpha(t_i)\mathbf{z}_\theta(\hat{\mathbf{z}}_{t_{i-1}}, \mathbf{x}, \tau_{i-1}) + \sigma(t_i)\epsilon_\theta(\hat{\mathbf{z}}_{t_{i-1}}, \mathbf{x}, \tau_{i-1}),$$

where the predicted noise is also updated as

$$\epsilon_\theta(\hat{\mathbf{z}}_{t_{i-1}}, \mathbf{x}, \tau_{i-1}) = \frac{\hat{\mathbf{z}}_{t_{i-1}} - \alpha(\tau_{i-1})\mathbf{z}_\theta(\hat{\mathbf{z}}_{t_{i-1}}, \mathbf{x}, \tau_{i-1})}{\sigma(\tau_{i-1})}.$$

We also demonstrate how CEDI works in Alg. 2.

## 5 Experiment

We conduct experiments to verify the effectiveness of the DINOISER and study its characteristics.

---

[5]Empirically, we find that $\tau_M$ satisfying $\sigma(\tau_M) = 0.99$ (*i.e.*, $\tau_M = 0.99$ for $\sigma(t) = t$ and $\tau_M = 0.9606$ for Li et al. (2022)'s `sqrt` schedule $\sigma(t) = t^{0.25}$) generally works well.

**Algorithm 2** Sampling with DiNoiser

**Input** Source condition $\mathbf{x}$; number of steps $M$; model parameters $\theta$.

**Output** Predicted target $\hat{\mathbf{y}}$.

1: Uniformly discretize $[T, 1]$ into $M{+}1$ steps $\{t_i\}_{i=0}^M$ in descend order ($T \approx 0$)
2: Uniformly discretize $[\mathcal{T}, 1]$ into $M + 1$ steps $\{\tau_i\}_{i=0}^M$ in descend order ($\mathcal{T} \gg 0$)
3: Sample $\hat{\mathbf{z}}_{t_0} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
4: **for** $i = 1$ to $M$ **do**
5: $\quad \hat{\mathbf{z}}_0 \leftarrow \mathbf{z}_\theta(\hat{\mathbf{z}}_{t_{i-1}}, \mathbf{x}, \tau_{i-1})$
6: $\quad \hat{\epsilon} \leftarrow \frac{\hat{\mathbf{z}}_{t_{i-1}} - \alpha(\tau_{i-1})\hat{\mathbf{z}}_0}{\sigma(\tau_{i-1})}$
7: $\quad \hat{\mathbf{z}}_{t_i} \leftarrow \alpha(t_i)\hat{\mathbf{z}}_0 + \sigma(t_i)\hat{\epsilon}$
8: **end for**
9: Map $\hat{\mathbf{z}}_{t_M}$ to $\hat{\mathbf{y}}$ with the embeddings

## 5.1 Experimental Setup

**Tasks and Datasets.** We mainly experiment on machine translation, a well-established benchmark task for conditional sequence learning. We consider IWSLT14 DE↔EN (160K pairs), WMT14 EN↔DE (4.0M pairs), and WMT14 EN↔RO (610K pairs), six translation tasks with variant sizes of training data. Additionally, we experiment on two of the datasets introduced by DiffuSeq (Gong et al., 2022), including Wiki (Jiang et al., 2020) for text simplification and QQP[6] for paraphrasing.

**Baselines.** We include three groups of baselines for machine translation: (1) The autoregressive Transformer (Vaswani et al., 2017); (2) The CMLM (Ghazvininejad et al., 2019), an iterative-based non-autoregressive model for conditional sequence learning. (3) Previous diffusion-based sequence generative models, including the vanilla design that simply extends the original DiffusionLM (Li et al., 2022) with an additional condition encoder, and the other recently proposed improved methods CDCD (continuous diffusion for categorical data, Dieleman et al., 2022), DiffuSeq (Gong et al., 2022), SeqDiffuSeq (Yuan et al., 2022) and Difformer (Gao et al., 2022). For text simplification and paraphrasing, we compare our method with DiffuSeq (Gong et al., 2022).

**Metrics.** We primarily report SacreBLEU[7] (Post, 2018) for machine translation, following

CDCD (Dieleman et al., 2022). For text simplification and paraphrasing, we follow DiffuSeq to employ sentence-level BLEU under the tokenizer of `BERT-BASE-UNCASED`.

**Implementation.** All our implementations are based on `Transformer-BASE` (Vaswani et al., 2017) for all datasets except IWSLT14. For IWSLT14, we use a smaller architecture that has 4 attention heads and 1024-dimensional feedforward layers. The embedding dimension for the diffusion model is 16 on IWSLT14 and 64 on the others. In the implementation of our method, we follow recent advances and apply self-conditioning techniques (Dieleman et al., 2022; Chen et al., 2022; Strudel et al., 2022). Besides, following previous practice in non-autoregressive machine translation, we train our model both with and without knowledge distillation[8] (KD, Kim and Rush, 2016; Zhou et al., 2020).

During inference, for machine translation, we apply beam search in the autoregressive Transformer with beam size 5. Correspondingly, we use length beam 5 in the non-autoregressive models, except for CDCD and DiffuSeq since they vary the target lengths by predicting paddings instead of length predictions. For text simplification and paraphrasing, we report results with various length beams as length prediction on these tasks is more challenging and less studied. For all the diffusion-based methods, we follow previous work (Li et al., 2022; Gong et al., 2022; Dieleman et al., 2022) and apply Minimum Bayes-Risk (MBR) decoding (Kumar and Byrne, 2004). For both DiffusionLM and our model, we perform sampling with 20 steps.

We implement DiffusionLM and DiNoiser upon `fairseq` (Ott et al., 2019), and also train Transformer and CMLM baselines using `fairseq`. The training batch size is 128K for WMT14/WMT16, and 32K for the others. For more details, please refer to §B.

## 5.2 Main Results

The results of machine translation and the other two tasks are in Tab. 2 and Tab. 3, respectively.

---

[6] https://www.kaggle.com/c/quora-question-pairs

[7] The signature is `nrefs:1|case:mixed|eff:no|tok:intl|smooth:exp|version:2.3.1` if the target language is German, and `nrefs:1|case:mixed|eff:no|tok:13a|smooth:exp|version:2.3.1` for others.

[8] Non-autoregressive sequence learning models typically struggle with learning multimodal distributions (Gu et al., 2018). For this reason, a common technique to improve their performance is to apply knowledge distillation, which simplifies the target distribution by replacing target samples with predictions from an autoregressive teacher model.

Table 2: Comparison in **SacreBLEU** on machine translation tasks. "LB": the size of the length beam search. "MBR": the number of candidates for each length beam to apply Minimum Bayes-Risk decoding. "**KD**": results are obtained with knowledge distillation (KD, Kim and Rush, 2016; Zhou et al., 2020). Provided that KD is common and effective practice in non-autoregressive (NAR) machine translation, though not the focus of this study, we also provide further experiments with KD for reference. The best NAR results with and without KD are in **bold** and the second best ones are underlined. We report 95% confidential interval for our method computed with `compare-mt` (Neubig et al., 2019).
†: how CMLM originally selects candidates with different lengths differs from the MBR decoding we used for diffusion models, and we thus include its results with MBR decoding for fair comparisons. ‡: the results are quoted from Dieleman et al. (2022). ‡: the results are quoted from Gao et al. (2022) while the results of the rest datasets are missing in the original paper, for which we obtain through their opensource code. Note that the results of DiffuSeq and SeqDiffuSeq are presented in tokenized BLEU as reported in Gao et al. (2022), and we encourage readers to check the original papers for more details.

| Methods | IWSLT14 | | WMT14 | | WMT16 | |
|---|---|---|---|---|---|---|
| | DE→EN | EN→DE | DE→EN | EN→DE | RO→EN | EN→RO |
| Transformer (Vaswani et al., 2017) (AR, beam = 5) | 33.61 | 28.30 | 30.55 | 26.85 | 33.08 | 32.86 |
| CMLM (Ghazvininejad et al., 2019) (NAR, LB = 5) | 29.41 | 24.33 | 28.71 | 23.22 | 31.13 | **31.26** |
| CMLM (Ghazvininejad et al., 2019) (NAR, LB = 5, MBR=1†) | 29.32 | 24.34 | 28.43 | 23.09 | 31.07 | 30.92 |
| DiffusionLM (Li et al., 2022) (LB = 5, MBR = 1) | 26.61 | 20.29 | 17.31 | 15.33 | 28.61 | 27.01 |
| DiffusionLM (Li et al., 2022) (LB = 5, MBR = 10) | 29.11 | 22.91 | 19.69 | 17.41 | 30.17 | 29.39 |
| CDCD (Dieleman et al., 2022) (MBR = 10) | - | - | 25.40‡ | 19.70‡ | - | - |
| CDCD (Dieleman et al., 2022) (MBR = 100) | - | - | 26.00‡ | 20.00‡ | - | - |
| Difformer (Gao et al., 2022) (LB × MBR = 20) | 28.01 | 23.31 | 25.30 | 23.80‡ | 29.37 | 29.20 |
| DINOISER (LB = 5, MBR = 1) | $31.29_{0.67}$ | $25.55_{0.65}$ | $28.83_{0.92}$ | $24.25_{0.86}$ | $31.14_{1.13}$ | $30.93_{1.12}$ |
| DINOISER (LB = 5, MBR = 10) | $\mathbf{31.61}_{0.67}$ | $\underline{25.70}_{0.62}$ | $\mathbf{29.05}_{0.92}$ | $\underline{24.26}_{0.84}$ | $\underline{31.22}_{1.15}$ | $\underline{31.08}_{1.12}$ |
| DINOISER (LB = 10, MBR = 5) | $\underline{31.44}_{0.68}$ | $\mathbf{26.14}_{0.65}$ | $\underline{29.01}_{0.88}$ | $\mathbf{24.62}_{0.88}$ | $\mathbf{31.24}_{1.12}$ | $31.03_{1.13}$ |
| DiffuSeq (Gong et al., 2022) (**KD**, LB×MBR = 10) | - | - | - | 15.37‡ | - | 25.45‡ |
| SeqDiffuSeq‡ (Yuan et al., 2022) (**KD**, LB×MBR = 10) | - | - | - | 17.14‡ | - | 26.17‡ |
| DINOISER (**KD**, LB = 10, MBR = 5) | - | - | $\mathbf{30.30}_{0.94}$ | $\mathbf{25.88}_{0.95}$ | $33.13_{1.20}$ | $32.84_{1.16}$ |

**Overall performance.** Our DINOISER demonstrates effectiveness on all selected conditional sequence learning tasks, which we summarize into the following three aspects:

- DINOISER achieves state-of-the-art results among diffusion-based models on one of the representative conditional sequence generation tasks, *i.e.*, machine translation, where DINOISER outperforms the vanilla design of DiffusionLM, as well as the previous strongest approaches such as CDCD (Dieleman et al., 2022) and Difformer (Gao et al., 2022) by a large margin (Tab. 2). For DiffuSeq and SeqDiffuSeq, although their reported tokenized BLEUs are not strictly comparable to our SacreBLEU results due to the difference in tokenizers, our performance is far above them by over 4 BLEU score and even more if we involve knowledge distillation, which supports our superiority over them.
- DINOISER demonstrates strong competitiveness in conditional sequence learning. It even surpasses CMLM (Ghazvininejad et al., 2019) on almost all the experimented machine translation datasets (Tab. 2). Provided

that CMLM is one of the leading approaches among NAR sequence learners, the performance DINOISER achieves can be considered quite competitive.

- DINOISER is generic to various conditional sequence learning tasks. Results on Tab. 3 shows that DINOISER also works well in tasks other than machine translation, surpassing previously proposed DiffuSeq.

In addition to the overall performance, DINOISER also demonstrates several nice properties. We elaborate on them as follows:

**Scalability.** As shown in Tab. 2, DiffusionLM seems more challenging to accommodate larger datasets like WMT14 than smaller ones (*e.g.*, IWSLT14). This verifies the curse of scalability problem discussed in §3. In contrast, DINOISER surpasses CMLM on almost all large- and small-scale scenarios, which indicates that DINOISER indeed greatly improves the scalability of diffusion-based sequence learners. This advantage of DINOISER could help facilitate further research of large-scale real-world applications of diffused sequence generative models.

Table 3: **Sentence-level BLEU** of our method and Diffuseq on Wiki (text simplification) and QQP (paraphrasing). "NFE": number of function evaluations, measuring the total number of forward passes to the model for each prediction. The results of DiffuSeq are quoted from Gong et al. (2022).

| Methods | Steps | LB | MBR | NFE | Wiki | QQP |
|---------|-------|-----|-----|------|------|-----|
| DiffuSeq | 2000 | - | 10 | 20000 | 36.22 | 24.13 |
| DiNoiser | 20 | 10 | 1 | 200 | $35.36_{1.63}$ | $\mathbf{26.07}_{1.24}$ |
| DiNoiser | 20 | 20 | 1 | 400 | $\mathbf{36.94}_{1.95}$ | $25.42_{1.46}$ |
| DiNoiser | 20 | 20 | 5 | 2000 | $36.88_{1.77}$ | $25.57_{1.29}$ |

**Sampling efficiency.** Given the sizes of length beam (LB) and MBR decoding shown in Tab. 2, DiNoiser surpasses or closely approaches CMLM even when MBR=1, while the vanilla DiffusionLM heavily relies on a large number of candidates for MBR decoding. Besides, DiNoiser necessitates much fewer NFEs to achieve strong performance, *e.g.*, only 20 steps, resulting in only 1% to 10% computational costs and latency compared to previous works (Gong et al., 2022; Li et al., 2022; Dieleman et al., 2022). This manifests that DiNoiser is more accurate yet efficient compared to previous diffusion-based sequence learning models.

## 5.3 Effect of Noise Scale Clipping for Training

**Ablation study on noise clipping.** We compare models trained with different settings in Tab. 4 to study the effect of our training strategy. We find that the proposed noise scale clipping consistently helps improve the model performance. Replacing the noise schedule in DiNoiser (*i.e.*, $\sigma(t) = t$) with the `sqrt` schedule proposed by Li et al. (2022) has negligible influence on the final performance. This is expected since we only set the noise schedule as $\sigma(t) = t$ for convenience. The difference is that the improvement by noise clipping is relatively smaller when using the `sqrt` schedule. This is because the noise scale of the sqrt schedule increases quickly in small timesteps (Fig. 2B). When $t = 0.2$, $\sigma_{sqrt}(t) \approx 0.67$, which is close to our initial clipping threshold. This suggests the success of the sqrt schedule may also be partly explained as involving more large-scale noises.

**On the effect of different noise scale clipping thresholds.** To understand how different noise clipping thresholds affect the model performance, we compare our model trained with adaptive noise threshold and different fixed noise thresholds in Fig. 4. Results show that the performance degrades if the clipping threshold is either too small or too large. Our proposed strategy adaptively finds the



Figure 4: SacreBLEU on IWSLT14 DE→EN with our models trained with adaptive *vs.* different fixed noise clipping thresholds. We sample results with MBR=5 and oracle length. The star marker (⋆) stands for the clipping threshold of our final checkpoint trained with adaptive clipping threshold.

balance well. It finds the ideal noise clipping threshold and achieves strong performance.

## 5.4 Effect of Condition Enhancement for Sampling

We compare the performance between different denoisers, *i.e.*, DDIM and the proposed CeDi, in Tab. 4. In a nutshell, CeDi impressively outperforms DDIM, especially for small MBR candidate sizes. We also notice that DDIM performs particularly unsatisfactorily when the model is trained without noise scale clipping. However, for these models, CeDi can still produce a relatively good performance (over 20.50 for MBR=1) that even surpasses well-designed CDCD (20.00 for MBR=100). What's more, we highlight two critical characteristics of CeDi as follows:

**CeDi indeed better leverage source conditions for inference.** Recall that we propose the CeDi with the purpose of encouraging the model to make better use of source conditions for prediction (§4.2). To investigate whether the denoiser achieves this, we apply Layer-wise Relevant Propagation (LRP, Bach et al., 2015; Voita et al., 2021) to measure the relative contribution of the source condition to the model prediction. As shown in Fig. 5(B), we compare the source contribution of our CeDi and the DDIM solver along the sampling iterations. CeDi maintains a high source contribution, while the source contribution of CeDi is unsatisfactory in the first few steps, which demonstrates that sampling with our CeDi does leverage the source condition more. Correspondingly, as shown in Fig. 5(A), the prediction accuracy of CeDi increases steadily, while the performance of DDIM fails to improve

Table 4: Ablation Study on WMT14 EN→DE. All the results are in SacreBLEU scores with $LB = 5$.

| Training Settings | DDIM (MBR = 1) | CEDI (MBR = 1) | DDIM (MBR = 10) | CEDI (MBR = 10) |
|---|---|---|---|---|
| Ours [final] | 19.23 | 24.25 | 22.12 | 24.26 |
| w/o self-conditioning | 20.37 | 23.03 | 22.58 | 23.14 |
| w/o noise scale clipping | 7.95 | 21.16 | 11.51 | 21.40 |
| w/o self-conditioning, w/o noise scale clipping | 11.30 | 20.86 | 14.84 | 21.47 |
| w/ sqrt noise schedule | 20.11 | 24.13 | 22.83 | 24.07 |
| w/ sqrt noise schedule, w/o noise scale clipping | 16.68 | 23.22 | 20.46 | 23.40 |



Figure 5: The difference between CEDI and DDIM solver over steps. (A) The prediction accuracy at each step, measured with SacreBLEU. (B) The proportion of source contribution to the prediction in Layer-wise Relevant Propagation (LRP) at each step.

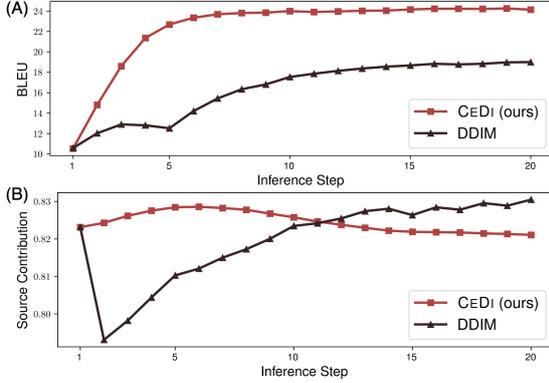Table 5: Results of multilingual machine translation ({DE,RO,PT-BR,NL}↔EN). "BILINGUAL": integrated results of separate models of every language pair. "MULTILING.": results from a unified multilingual model. We employ langdetect (Shuyo, 2010) to infer the language of generated sequences for computing the language accuracy.

| Settings | Methods | SacreBLEU (Lang Acc %) | |
|---|---|---|---|
| | | BILINGUAL | MULTILING. |
| many-to-one {DE,RO,PT-BR,NL} →EN | CMLM | 33.85 | 35.23 |
| | DINOISER (DDIM, MBR=1) | 32.40 | 33.43 |
| | DINOISER (DDIM, MBR=10) | 33.73 | 35.48 |
| | DINOISER (MBR=1) | 34.57 | 35.26 |
| | DINOISER (MBR=10) | 34.74 | 35.66 |
| one-to-many EN→ {DE,RO,PT-BR,NL} | CMLM | 28.10 | 30.55 (94.73) |
| | DINOISER (DDIM, MBR=1) | 27.44 | 17.95 (89.77) |
| | DINOISER (DDIM, MBR=10) | 28.54 | 18.57 (89.53) |
| | DINOISER (MBR=1) | 28.72 | 30.52 (95.31) |
| | DINOISER (MBR=10) | 28.81 | 30.67 (95.42) |

at the beginning of the iterations, suggesting correlations between higher source contribution and higher performance improvement. Among all iterations, the first few steps establish the foundation for the overall performance. Although DDIM improves its performance in later iterations, it still falls behind our CEDI. This suggests the effectiveness of increasing the source contribution, especially at the beginning of the sampling process.

To further show the strength of CEDI in capturing source conditions, we explore more complex conditional sequence learning scenarios. We simulate this under two multilingual translation settings, *i.e.* many-to-one and one-to-many translation. In the many-to-one scenario, a unified model needs to translate source sentences in multiple different source languages to English counterparts, requiring the model to handle complicated source conditions. On the other hand, the one-to-many setting simulates a multi-conditional scenario, requiring the model to recognize the target language as another crucial condition to capture the target distribution. To this end, we construct a dataset by combining four language pairs of IWSLT14 translation benchmark, *i.e.*, EN↔DE, EN↔RO, EN↔NL, and EN↔PT-BR. In the one-to-many translation, we append language tokens to the source sequences to incor-

porate the target language as a condition. We also include a baseline in which the models are trained separately for each language pair for comparison.

**CEDI can handle sequence generation from complex and multiple conditions.** As shown in Tab. 5, DINOISER works well in multilingual settings, showing its strong capability in modeling conditions, *i.e.* a complex multimodal condition (source sentences of four languages in many-to-one), and multiple conditions (source sentence in English as well as the identity of target languages). Particularly, CEDI shows huge advantages over DDIM in the multilingual setting of one-to-many translation. In this case, the language accuracy of DDIM is much lower than that of CEDI, suggesting DDIM has trouble capturing the given condition, namely the language identity in this one-to-many scenario. In contrast, DINOISER augmented with CEDI yields satisfactory predictions with high language accuracy, exhibiting superiority in working with multiple conditions. This is also consistent with our findings from the qualitative examples as shown in Tab. 6, where DDIM may fail to capture the language condition and generate non-sense articles shared across languages, while the full DINOISER produces fluent and satisfactory results.

Table 6: A quanlitative example for one-to-many translation. The source contains both the English sentence to be translated and the target language. We compare generation results of CMLM, DINOISER but sampling with DDIM instead of CEDI, and the complete DINOISER.

| Source | (target language: ro) something as dramatic as our identity has now become a matter of choice, as this slide is meant to indicate. |
|---|---|
| Reference | ceva atât de important ca identitatea noastră a devenit acum o problemă de alegere, și această tranziție are rolul de a arăta ast |
| CMLM | ceva la fel de dramatic ca identitatea noastră a devenit o problemă de alegere, cum acest slide este făcut să arate. |
| DINOISER (w/ DDIM) | ceva de de de de de de de de de a de de de de de de de de de de de de de de. |
| DINOISER (w/ CEDI) | ceva atât de dramatic ca identitatea noastră, a devenit acum o problemă de alegere, așa cum se înseamnă să indice acest imagine |

## 6 Related Work

**Non-autoregressive Sequence Generative Models.** Non-autoregressive sequence learning (NAR) was first proposed by Gu et al. (2018) as an alternative to its autoregressive counterpart. It generates target tokens in parallel, either fully NAR (Gu et al., 2018) or up to a mild number of iterations (Ghazvininejad et al., 2019), liberating sequence modeling from the constraint of a predefined order (Qian et al., 2022). With recent efforts, NAR shows great potential in the applications of various domains, including language (Qian et al., 2021; Qi et al., 2021; Huang et al., 2022c; Qian et al., 2022), speech (Kim et al., 2021), proteins (Zheng et al., 2023a; Wang et al., 2024), and molecules (Hoogeboom et al., 2022). Different from more commonly-used autoregressive (AR) models (Sutskever et al., 2014), NAR models assume conditional independence among the output tokens. Such an assumption risks ignoring the target dependencies (Ren et al., 2020; Huang et al., 2022b) and leads to the multi-modality problem (Gu et al., 2018). As a result, the vanilla fully NAR model has inferior generation quality. Some of the later improvements alleviate the strong assumption by reformulating NAR formulation under iterative refinement (Lee et al., 2018; Gu et al., 2019; Ghazvininejad et al., 2019, 2020; Huang et al., 2022a,d; Zheng et al., 2023b; Ye et al., 2023), which iteratively takes as input the previously generated sequence, which serves as an intermediate random variable, to produce the tokens of its refined or denoised version in parallel until convergence or the budget of maximum iterations run out. Some recent advances herein follow the idea of discrete diffusion (Sohl-Dickstein et al., 2015; Austin et al., 2021) and formalize iterative refinement as Markov processes (Savinov et al., 2021; He et al., 2022; Reid et al., 2022). Although both are named after diffusion models, these works operate on discrete state space, whereas our focus, continuous diffusion models accommodate the continuous (embedding) space of discrete tokens.

**Diffusion Models for Sequence Learning.** Continuous diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020; Song et al., 2020b) gained first success in generating high-quality images. Recently, DiffusionLM (Li et al., 2022) successfully adapted them to sequence learning and proposed the DiffusionLM, the first diffusion-based sequence generative model with a special focus on controllable text generation. Later improvements to the diffusion-based sequence generative models are mainly categorized threefold. The first line includes novel components for diffusion modeling, such as the partial diffusion process proposed by DiffuSeq (Gong et al., 2022), self-conditioning techniques introduced by Strudel et al. (2022), and the adaptive noise schedule of Yuan et al. (2022). The second line applies diffusion models to the latent space of specific pretrained language models (Lovelace et al., 2022). And the third tries to incorporate conventional practice in discrete token prediction. For instance, CDCD (Dieleman et al., 2022), Difformer (Gao et al., 2022) and SSD (Han et al., 2022) incorporate the cross-entropy objectives in training. For the application of diffusion-based models for sequence learning, previous work found their advantages in controllable generation (Yu et al., 2022; Liu et al., 2022; Li et al., 2022), and generating diverse sequences (Gong et al., 2022). GENIE (Lin et al., 2022) demonstrates that diffusion-based sequence generative models can benefit from large-scale self-supervised pretraining. While almost all these works mainly focus on the training phrase of diffusion-based sequence generative models, our study emphasizes both training and inference.

## 7 Conclusion

In this paper, we shed light on the crucial role of noise schedules in diffusion models for conditional sequence learning by systematic empirical study. Motivated by our findings, we propose DINOISER to determine the best-suited noise scales for both training and inference. As a result, DINOISER

makes training more effective and also enables the model to better utilize source conditions for prediction, thereby leading to considerable performance improvements. We expect that our study can help facilitate further research on diffusion models to empower various applications in NLP.

## Acknowledgement

We would like to thank the anonymous reviewers and editors for their invaluable feedback.

## References

Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne van den Berg. 2021. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993.

Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixelwise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one*, 10(7):e0130140.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Fan Bao, Shen Nie, Kaiwen Xue, Chongxuan Li, Shi Pu, Yaole Wang, Gang Yue, Yue Cao, Hang Su, and Jun Zhu. 2023. One transformer fits all distributions in multi-modal diffusion at scale.

Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. Wit3: Web inventory of transcribed and translated talks. In *Conference of european association for machine translation*, pages 261–268.

Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. 2022. Analog bits: Generating discrete data using diffusion models with self-conditioning. *arXiv preprint arXiv:2208.04202*.

Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin, Pierre H Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, et al. 2022. Continuous diffusion for categorical data. *arXiv preprint arXiv:2211.15089*.

Zhujin Gao, Junliang Guo, Xu Tan, Yongxin Zhu, Fang Zhang, Jiang Bian, and Linli Xu. 2022. Difformer: Empowering diffusion model on embedding space for text generation. *arXiv preprint arXiv:2212.09412*.

Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. 2019. Mask-predict: Parallel decoding of conditional masked language models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6112–6121.

Marjan Ghazvininejad, Omer Levy, and Luke Zettlemoyer. 2020. Semi-autoregressive training improves mask-predict decoding. *arXiv preprint arXiv:2001.08785*.

Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and LingPeng Kong. 2022. Diffuseq: Sequence to sequence text generation with diffusion models. *arXiv preprint arXiv:2210.08933*.

Jiatao Gu, James Bradbury, Caiming Xiong, Victor OK Li, and Richard Socher. 2018. Non-autoregressive neural machine translation. In *International Conference on Learning Representations*.

Jiatao Gu and Xiang Kong. 2021. Fully non-autoregressive neural machine translation: Tricks of the trade. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 120–133.

Jiatao Gu, Changhan Wang, and Junbo Zhao. 2019. Levenshtein transformer. *Advances in Neural Information Processing Systems*, 32.

Xiaochuang Han, Sachin Kumar, and Yulia Tsvetkov. 2022. Ssd-lm: Semi-autoregressive simplex-based diffusion language model for text generation and modular control. *arXiv preprint arXiv:2210.17432*.

Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. 2022. Diffusionbert: Improving generative masked language models with diffusion models. *arXiv preprint arXiv:2211.15029*.

Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko,

Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. 2022. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pages 6840–6851. Curran Associates, Inc.

Emiel Hoogeboom, Vıctor Garcia Satorras, Clément Vignac, and Max Welling. 2022. Equivariant diffusion for molecule generation in 3d. In *International Conference on Machine Learning*, pages 8867–8887. PMLR.

Chenyang Huang, Hao Zhou, Osmar R Zaïane, Lili Mou, and Lei Li. 2022a. Non-autoregressive translation with layer-wise prediction and deep supervision. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10776–10784.

Fei Huang, Tianhua Tao, Hao Zhou, Lei Li, and Minlie Huang. 2022b. On the learning of non-autoregressive transformers. In *International Conference on Machine Learning*, pages 9356–9376. PMLR.

Fei Huang, Hao Zhou, Yang Liu, Hang Li, and Minlie Huang. 2022c. Directed acyclic transformer for non-autoregressive machine translation. In *International Conference on Machine Learning*, pages 9410–9428. PMLR.

Xiao Shi Huang, Felipe Perez, and Maksims Volkovs. 2022d. Improving non-autoregressive translation models without distillation. In *International Conference on Learning Representations*.

Chao Jiang, Mounica Maddela, Wuwei Lan, Yang Zhong, and Wei Xu. 2020. Neural crf model for sentence alignment in text simplification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7943–7960.

Jaehyeon Kim, Jungil Kong, and Juhee Son. 2021. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *International Conference on Machine Learning*, pages 5530–5540. PMLR.

Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327.

Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.

Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. 2020. Diffwave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*.

Shankar Kumar and William Byrne. 2004. Minimum bayes-risk decoding for statistical machine translation. Technical report, JOHNS HOPKINS UNIV BALTIMORE MD CENTER FOR LANGUAGE AND SPEECH PROCESSING (CLSP).

Matt Kusner, Yu Sun, Nicholas Kolkin, and Kilian Weinberger. 2015. From word embeddings to document distances. In *International conference on machine learning*, pages 957–966. PMLR.

Dohyun Kwon, Ying Fan, and Kangwook Lee. 2022. Score-based generative modeling secretly minimizes the wasserstein distance. *Advances in Neural Information Processing Systems*, 35:20205–20217.

Jason Lee, Elman Mansimov, and Kyunghyun Cho. 2018. Deterministic non-autoregressive neural sequence modeling by iterative refinement. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1173–1182.

Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori Hashimoto. 2022. Diffusion-lm improves controllable text generation. *ArXiv*, abs/2205.14217.

Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Weizhu Chen, and Nan Duan. 2022. Genie: Large scale pretraining for text generation with diffusion model. *arXiv preprint arXiv:2212.11685*.

Guangyi Liu, Zeyu Feng, Yuan Gao, Zichao Yang, Xiaodan Liang, Junwei Bao, Xiaodong He, Shuguang Cui, Zhen Li, and Zhiting Hu. 2022.

Composable text controls in latent space with odes. *arXiv preprint arXiv:2208.00638*.

Justin Lovelace, Varsha Kishore, Chao Wan, Eliot Shekhtman, and Kilian Weinberger. 2022. Latent diffusion for language generation. *arXiv preprint arXiv:2212.09462*.

Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2022. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *arXiv preprint arXiv:2206.00927*.

Graham Neubig, Zi-Yi Dou, Junjie Hu, Paul Michel, Danish Pruthi, Xinyi Wang, and John Wieting. 2019. compare-mt: A tool for holistic comparison of language generation systems. *CoRR*, abs/1903.07926.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. *arXiv preprint arXiv:1904.01038*.

Matt Post. 2018. A call for clarity in reporting bleu scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191.

Aaditya Prakash, Sadid A. Hasan, Kathy Lee, Vivek V. Datla, Ashequl Qadir, Joey Liu, and Oladimeji Farri. 2016. Neural paraphrase generation with stacked residual LSTM networks. *CoRR*, abs/1610.03098.

Weizhen Qi, Yeyun Gong, Jian Jiao, Yu Yan, Weizhu Chen, Dayiheng Liu, Kewen Tang, Houqiang Li, Jiusheng Chen, Ruofei Zhang, Ming Zhou, and Nan Duan. 2021. Bang: Bridging autoregressive and non-autoregressive generation with large scale pretraining. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8630–8639. PMLR.

Lihua Qian, Mingxuan Wang, Yang Liu, and Hao Zhou. 2022. Diff-glat: Diffusion glancing transformer for parallel sequence to sequence learning. *arXiv preprint arXiv:2212.10240*.

Lihua Qian, Yi Zhou, Zaixiang Zheng, Yaoming Zhu, Zehui Lin, Jiangtao Feng, Shanbo Cheng,

Lei Li, Mingxuan Wang, and Hao Zhou. 2021. The volctrans glat system: Non-autoregressive translation meets wmt21. *WMT 2021*, page 187.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J Liu, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67.

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*.

Machel Reid, Vincent Josua Hellendoorn, and Graham Neubig. 2022. Diffuser: Diffusion via edit-based reconstruction. In *International Conference on Learning Representations*.

Yi Ren, Jinglin Liu, Xu Tan, Zhou Zhao, Sheng Zhao, and Tie-Yan Liu. 2020. A study of non-autoregressive model for sequence generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 149–159.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2021. High-resolution image synthesis with latent diffusion models.

Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 379–389, Lisbon, Portugal. Association for Computational Linguistics.

Nikolay Savinov, Junyoung Chung, Mikolaj Binkowski, Erich Elsen, and Aaron van den Oord. 2021. Step-unrolled denoising autoencoders for text generation. In *International Conference on Learning Representations*.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725.

Nakatani Shuyo. 2010. Language detection library for java.

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. 2015. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2256–2265, Lille, France. PMLR.

Jiaming Song, Chenlin Meng, and Stefano Ermon. 2020a. Denoising diffusion implicit models. In *International Conference on Learning Representations*.

Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. 2020b. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*.

Robin Strudel, Corentin Tallec, Florent Altché, Yilun Du, Yaroslav Ganin, Arthur Mensch, Will Grathwohl, Nikolay Savinov, Sander Dieleman, Laurent Sifre, et al. 2022. Self-conditioned embedding diffusion for text generation. *arXiv preprint arXiv:2211.04236*.

Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27.

Arash Vahdat, Karsten Kreis, and Jan Kautz. 2021. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Elena Voita, Rico Sennrich, and Ivan Titov. 2021. Analyzing the source and target contributions to predictions in neural machine translation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1126–1140.

Xinyou Wang, Zaixiang Zheng, Fei Ye, Dongyu Xue, Shujian Huang, and Quanquan Gu. 2024. Diffusion language models are versatile protein learners. *arXiv preprint arXiv:2402.18567*.

Antoine Wehenkel and Gilles Louppe. 2021. Diffusion priors in variational autoencoders. In *ICML Workshop on Invertible Neural Networks, Normalizing Flows, and Explicit Likelihood Models*.

Sam Wiseman, Stuart M Shieber, and Alexander M Rush. 2017. Challenges in data-to-document generation. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2253–2263.

Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Quanquan Gu. 2023. Diffusion language models can perform many tasks with scaling and instruction-finetuning. *arXiv preprint arXiv:2308.12219*.

P Yu, S Xie, X Ma, B Jia, B Pang, R Gao, Y Zhu, S-C Zhu, and YN Wu. 2022. Latent diffusion energy-based model for interpretable text modeling. In *International Conference on Machine Learning (ICML 2022)*.

Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Fei Huang, and Songfang Huang. 2022. Seqdiffuseq: Text diffusion with encoder-decoder transformers. *arXiv preprint arXiv:2212.10325*.

Zaixiang Zheng, Yifan Deng, Dongyu Xue, Yi Zhou, Fei YE, and Quanquan Gu. 2023a. Structure-informed language models are protein designers. *arXiv preprint arXiv:2302.01649*.

Zaixiang Zheng, Yi Zhou, and Hao Zhou. 2023b. Deep equilibrium non-autoregressive sequence learning. In *Findings of the Association for Computational Linguistics: ACL 2023*.

Chunting Zhou, Graham Neubig, and Jiatao Gu. 2020. Understanding knowledge distillation in non-autoregressive machine translation. In *International Conference on Learning Representations*.

## A  More Results on Machine Translation

To provide references for further study and comparisons, we report more results on machine translation.

**Knowledge distillation (KD).**  A common practice to improve the performance of non-autoregressive machine translation is knowledge distillation (Kim and Rush, 2016; Zhou et al., 2020). We report the performance of our method trained on distilled data of WMT14 and WMT16 on Tab. 7. The result shows that the performance gap between our method and the autoregressive transformer is small when knowledge distillation is used. This suggests that our method achieves performance that satisfies the needs of applications.

Table 7: Model performances on machine translation with knowledge distillation. The results of transformer are from raw data, while DINOISER is trained on distilled data. The performances are measured with **SacreBLEU**.

| Methods | WMT14 | | WMT16 | |
| --- | --- | --- | --- | --- |
| | DE→EN | EN→DE | RO→EN | EN→RO |
| Transformer | 30.55 | 26.85 | 33.08 | 32.86 |
| DINOISER (LB=5, MBR=1) | 30.13 | 25.70 | 32.96 | 32.58 |
| DINOISER (LB=5, MBR=10) | 30.12 | 25.90 | 33.04 | 32.57 |
| DINOISER (LB=10, MBR=5) | 30.30 | 25.88 | 33.13 | 32.84 |

**Evaluation with tokenized BLEU.**  Some of the previous studies in machine translation reported tokenized BLEU, despite inconsistent tokenizers (other than the standard Moses tokenizer) they might use. To help conveniently compare DINOISER to them, we also report the performance of DINOISER with tokenized BLEU[9] in Tab. 8.

Table 8: Tokenized BLEU of our method on machine translation datasets. We use the `moses` tokenizer for all the texts. "LB": the size of length beam. "MBR": the number of candidates for each length beam to apply Minimum Bayes-Risk decoding. +KD means the results are obtained with knowledge distillation.

| Methods | IWSLT14 | | WMT14 | | WMT16 | |
| --- | --- | --- | --- | --- | --- | --- |
| | DE→EN | EN→DE | DE→EN | EN→DE | RO→EN | EN→RO |
| DINOISER (LB=5, MBR=1) | 32.23 | 25.54 | 29.35 | 24.43 | 31.21 | 31.18 |
| DINOISER (LB=5, MBR=10) | 32.48 | 25.68 | 29.53 | 24.45 | 31.39 | 31.29 |
| DINOISER (LB=10, MBR=5) | 32.25 | 25.99 | 29.40 | 24.48 | 31.50 | 31.27 |
| DINOISER + KD (LB=5, MBR=1) | - | - | 30.64 | 26.08 | 33.21 | 32.57 |
| DINOISER + KD (LB=5, MBR=10) | - | - | 30.62 | 26.29 | 33.29 | 32.59 |
| DINOISER + KD (LB=10, MBR=5) | - | - | 30.76 | 26.04 | 33.40 | 32.89 |

## B  Implementation Details

All our implementations are based on `Transforme-base` (Vaswani et al., 2017) for all datasets except IWSLT14. For IWSLT14, we use a smaller architecture that has 4 attention heads and 1024-dimensional feedforward layers. The embedding dimension for the diffusion model is 16 on IWSLT14 and 64 on the others. In the implementation of our method, we follow recent advances and apply self-conditioning techniques (Dieleman et al., 2022; Chen et al., 2022; Strudel et al., 2022). Besides, following previous practice in non-autoregressive machine translation, we train our model both with and without knowledge distillation (KD, Kim and Rush, 2016; Zhou et al., 2020).

During inference, we apply beam search in the autoregressive Transformer with beam size 5 for machine translation. Correspondingly, we use length beam 5 in the non-autoregressive models, except for CDCD and DiffuSeq since they vary the target lengths by predicting paddings instead of length predictions. For text simplification and paraphrasing, we report results with various length beams as length prediction on

---
[9]https://github.com/alvations/sacremoses

these tasks is more challenging and less studied. For all the diffusion-based methods, we follow previous work (Li et al., 2022; Gong et al., 2022; Dieleman et al., 2022) and apply Minimum Bayes-Risk (MBR) decoding (Kumar and Byrne, 2004). For both DiffusionLM and our model, we sample with 20 steps.

We implement DiffusionLM and DINOISER upon `fairseq` (Ott et al., 2019), and also train Transformer and CMLM baselines using `fairseq`. For data preprocessing, we follow the instruction in `fairseq` for IWSLT14[10] and use the preprocessed data by (Gu and Kong, 2021) for WMT14 and WMT16[11]. For Wiki and QQP, we use preprocessed data provided by DiffuSeq[12] and tokenized them with byte-pair encoding (BPE, Sennrich et al., 2016). The training batch size is 128K for WMT14/WMT16, and 32K for the others. We empirically find checkpoint averaging unnecessary for our method and have not applied it in all our implementations.

## C Relationship Between Different Noise Schedules and Time Samplers

Generally, the training objective of diffusion models can be expressed as

$$\mathbb{E}_{t\sim r(t),\epsilon\sim\mathcal{N}(\mathbf{0},\mathbf{I})}[w(t)\|\mathbf{z}_\theta - \mathbf{z}(0)\|_2^2],$$

where $\mathbf{z}_\theta$ is the model prediction $\mathbf{z}_\theta(\sqrt{1 - \sigma^2(t)}\mathbf{z}(0) + \sigma(t)\epsilon, t)$ for short.

The above expectation over timesteps can be rewritten into the expectation of noise scales as follows.

$$\mathbb{E}_{t\sim r(t),\epsilon}[w(t)\|\mathbf{z}_\theta - \mathbf{z}(0)\|_2^2]$$
$$=\mathbb{E}_\epsilon[\int_0^1 r(t)w(t)\|\mathbf{z}_\theta - \mathbf{z}(0)\|_2^2\mathrm{d}t]$$
$$=\mathbb{E}_\epsilon[\int_0^1 \hat{r}(\sigma)\hat{w}(\sigma)\|\mathbf{z}_\theta - \mathbf{z}(0)\|_2^2\frac{\mathrm{d}t}{\mathrm{d}\sigma}\mathrm{d}\sigma]$$
$$=\mathbb{E}_{\sigma\sim U(0,1),\epsilon}[w'(\sigma)\|\mathbf{z}_\theta - \mathbf{z}(0)\|_2^2],$$

where $w'(\sigma) = w(\sigma^{-1})r(\sigma^{-1})\frac{\mathrm{d}t}{\mathrm{d}\sigma}$. Therefore, training with different noise schedules and different time samplers is interchangeable by applying different weighting functions.

## D Effect of Sizes of Length Beam and MBR

DINOISER can leverage both length beam search and MBR decoding to produce diverse candidates for selection. The results on all of the evaluated datasets (Tab. 2 and Tab. 3) demonstrate that the method can gain its performance by properly adjusting the two hyperparameters for sampling. In particular, we search on various combinations of length beams and MBRs and evaluate the corresponding performance of DINOISER on the validation set of WMT14 EN→DE, shown in Fig. 6. The model performance rises first and then drops down as the length beam increases. And for each length beam, we can further boost the performance of DINOISER with MBR > 1, suggesting that the effects of the two factors are complementary.

Using both length beam search and MBR decoding also brings benefits to DINOISER over those only involving one of them. Compared to CMLM which decodes deterministically, DINOISER is able to sample multiple sentences for each length beam, providing more diverse candidates. Compared to CDCD, which predicts paddings to generate sentences of various lengths and whose sampling efficiency is restricted by maximum target length, DINOISER's use of length beams allows more fine-grained control of the computational budget.

---

[10]https://github.com/facebookresearch/fairseq/tree/main/examples/translation
[11]https://github.com/shawnkx/Fully-NAT
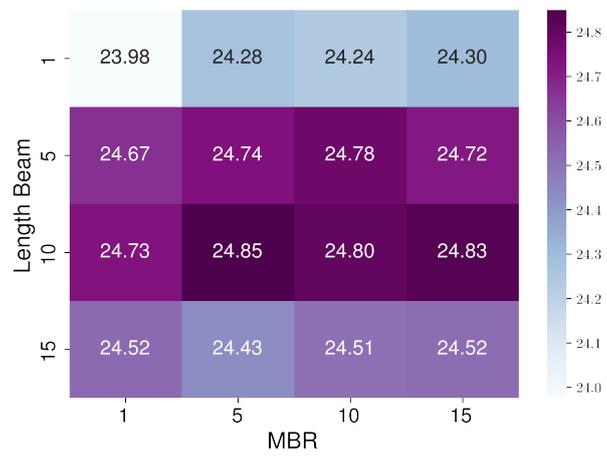[12]https://github.com/Shark-NLP/DiffuSeq

Figure 6: SacreBLEU on the validation set of WMT14 EN→DE with different length beams and MBR sizes.