

---

# AR-DIFFUSION: Auto-Regressive Diffusion Model for Text Generation

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Diffusion models have gained significant attention in the realm of image generation  
2 due to their exceptional performance. Their success has been recently expanded  
3 to text generation via generating all tokens within a sequence concurrently. How-  
4 ever, natural language exhibits a far more pronounced sequential dependency in  
5 comparison to images, and the majority of existing language models are trained  
6 with a left-to-right auto-regressive approach. To account for the inherent sequen-  
7 tial characteristic of natural language, we introduce Auto-Regressive Diffusion  
8 (AR-DIFFUSION). AR-DIFFUSION ensures that the generation of tokens on the  
9 right depends on the generated ones on the left, a mechanism achieved through  
10 employing a dynamic number of denoising steps that vary based on token position.  
11 This results in tokens on the left undergoing fewer denoising steps than those on  
12 the right, thereby enabling them to generate earlier and subsequently influence  
13 the generation of tokens on the right. In a series of experiments on various text  
14 generation tasks, including text summarization, machine translation, and common  
15 sense generation, AR-DIFFUSION clearly demonstrated its superiority over existing  
16 diffusion language models and that it can be  $100\times \sim 600\times$  faster when achieving  
17 comparable results. Our code will be publicly available.

## 18 1 Introduction

19 Text generation is a fundamental task within the field of natural language processing (NLP). Pre-  
20 trained language models like GPT-4 [OpenAI, 2023], LLaMA [Touvron et al., 2023], and Alpaca  
21 [Taori et al., 2023] have garnered significant attention with their ability to generate fluent and human-  
22 like textual content. These models utilize the auto-regressive (AR) Transformer decoders [Vaswani  
23 et al., 2017] to emit generated tokens one-by-one in sequential order from left to right. By leveraging  
24 the power of position dependency, AR models are able to enhance the naturalness, coherence, and  
25 adherence to human language conventions in the generated text [Brown et al., 2020].

26 Recent studies have shown the remarkable performance of diffusion models in image generation [Ho  
27 et al., 2020], motivating researchers to extend diffusion to text generation [Li et al., 2022a, Gong  
28 et al., 2022, Dieleman et al., 2022, Yuan et al., 2022, Ye et al., 2023]. By introducing timestep, these  
29 methods progressively regulate the interpolation between the original tokens and Gaussian noise, then  
30 iteratively denoise for text generation. At each timestep, the diffusion-based text generator predicts  
31 all tokens simultaneously following Non-Auto-Regression (NAR) [Lewis et al., 2020, Qi et al., 2020,  
32 2021, Li et al., 2022b], leading to faster decoding speed compared to AR. However, it also inherits  
33 the drawback of NAR, namely the sacrifice of inter-token position dependency [Li et al., 2022c] and  
34 the drop of generation performance [Bao et al., 2021].

35 To conduct a comprehensive analysis, we introduce a two-dimensional coordinate system to track the  
36 diffusion timestep of tokens  $f(\cdot)$  positioned at various locations. As illustrated in Figure 1, the system

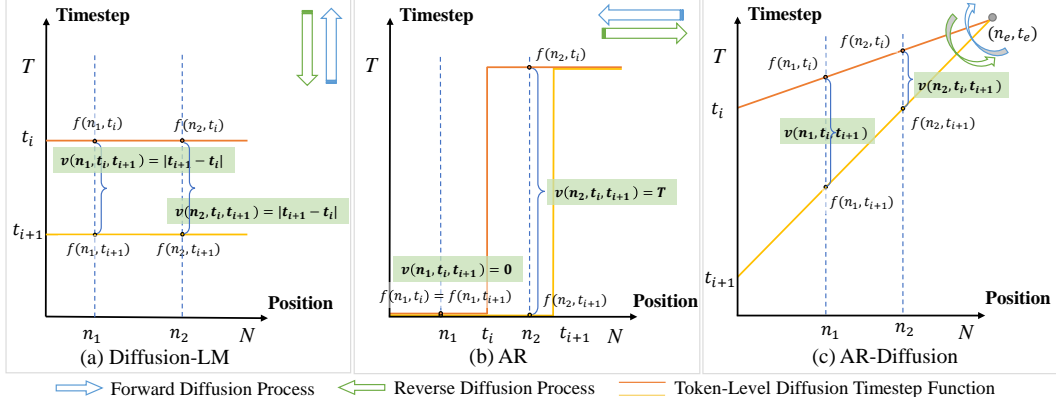


Figure 1: Model behaviors illustrated on a two-dimensional coordinate system, where the horizontal axis stands for the position and the vertical axis represents the diffusion timestep. In the inference stage, different models will behave differently. (a) For the typical Diffusion-LM [Li et al., 2022a], each token share the identical movement speed  $v(n_1, t_i, t_{i+1}) = v(n_2, t_i, t_{i+1}) = |t_{i+1} - t_i|$ . (b) For AR from the perspective of diffusion models, the tokens have two states based on the degree of interpolation between the original tokens and Gaussian noise: to be decoded (at timestep  $t = T$ ) and already decoded (at timestep  $t = 0$ ). Specifically, we have  $v(n_1, t_i, t_{i+1}) = 0$  and  $v(n_2, t_i, t_{i+1}) = T$ . (c) In AR-DIFFUSION,  $(n_e, t_e)$  is the coordinate of anchor point. Tokens in different positions exhibit varying movement speeds, such as  $v(n_1, t_i, t_{i+1}) > v(n_2, t_i, t_{i+1})$  when  $n_1 < n_2$ .

37 assigns the token position  $n \in [1, N]$  to the horizontal axis and the diffusion timestep  $t \in [0, T]$  to  
38 the vertical axis. Diffusion-LM [Li et al., 2022a], which is followed by existing diffusion-based text  
39 generation models, is shown in Figure 1(a). It assigns a uniform timestep  $t$  to all tokens. In contrast,  
40 tokens in the AR model depicted in Figure 1(b) exhibit distinct timesteps within a generation step ( $t_i$ ).  
41 For instance, the already decoded token at position  $n_1$  has a timestep of 0, while the to-be-decoded  
42 token at position  $n_2$  has a timestep of  $T$ . This approach effectively captures the sequential dependency.  
43 Motivated by this observation, we introduce AR-DIFFUSION, an auto-regressive diffusion method,  
44 for the disparity in token positions and the principle of sequential token identification.

45 In AR-DIFFUSION, we propose a **multi-level diffusion strategy** that includes both sentence-level  
46 and token-level diffusion. We randomly choose a sentence-level timestep  $t$ , and assign **dynamic**  
47 **movement speeds**  $v(\cdot)$  by determining position-sensitive token-level timestep  $f(n, t)$  for each token.  
48 This enables tokens at the left of a sentence to undergo faster movement from random Gaussian noise  
49 to token embedding, while those at the right of the sentence experience slower movement to better  
50 utilize information from previously denoised tokens. During inference, to reduce the significant  
51 number of inference steps (e.g., 2,000) required in Diffusion-LM [Li et al., 2022a], SeqDiffSeq [Yuan  
52 et al., 2022] and GENIE [Lin et al., 2023], we introduce a skipping mechanism that collaborates with  
53 the multi-level diffusion strategy to accelerate the process.

54 Experimental results across various text generation tasks, such as text summarization, machine  
55 translation, and common sense generation, have consistently demonstrated that AR-DIFFUSION  
56 surpasses existing text diffusion models, including AR methods in terms of both quality and diversity.  
57 Moreover, our verification reveals that AR-DIFFUSION requires fewer resources during decoding  
58 while maintaining superior performance. It achieves  $100\times$  faster than SeqDiffSeq [Yuan et al., 2022]  
59 in machine translation and  $600\times$  faster than GENIE [Lin et al., 2023] in text summarization while  
60 delivering comparable results. Furthermore, it demonstrates promising results even in a challenging  
61 scenario where decoding is limited to only two steps.

## 62 2 Preliminary

### 63 2.1 Conditional Generative Language Models

64 In the field of natural language generation, conditional generative models are commonly implemented  
65 using either auto-regressive (AR) or non-auto-regressive (NAR) methods. In AR [Vaswani et al.,

2017], tokens on the right are predicted based on visible left tokens. The likelihood is given by  $p_{\text{AR}}(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^N p(y_i|y_{1:i-1}; \mathbf{x})$ , where  $y_i$  denotes the  $i$ -th token of  $\mathbf{y}$ . On the other hand, NAR [Gu et al., 2017] assumes conditional independence among tokens and generates them uniformly without distinction during decoding, resulting in the likelihood  $p_{\text{NAR}}(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^N p(y_i|\mathbf{x})$ . This parallel generation approach is of lower quality compared to AR, although it offers a substantial speed advantage.

## 2.2 Diffusion Models for Text Generation

Recently, Li et al. [2022a] propose a natural language generation model based on the diffusion process, which is typically divided into a forward noising process and a reverse denoising process.

Specifically, the forward process is a fixed linear Gaussian model, which gradually perturbs the random variable  $\mathbf{z}_0$  until it becomes the standard Gaussian distribution. This can be formalized as:

$$q(\mathbf{z}_t | \mathbf{z}_0; \mathbf{x}) = N(\mathbf{z}_t; \bar{\alpha}_t \mathbf{z}_0; (1 - \bar{\alpha}_t) \mathbf{I}); \quad (1)$$

where,  $\bar{\alpha}_t = \prod_{j=1}^t \alpha_j$ , and  $\alpha_j$  is a coefficient that monotonically decreases with timestep  $t$ ,  $\mathbf{z}_t$  is the latent state at timestep  $t$ .

The reverse process is to initiate from standard Gaussian noise and progressively utilize the denoising transition  $p_{\theta}(\mathbf{z}_{t-1} | \mathbf{z}_t; \mathbf{x})$  for generation.

$$p(\mathbf{z}_{t-1} | \mathbf{z}_t; \mathbf{x}) = N(\mathbf{z}_{t-1}; \mu_{\theta}(\mathbf{z}_t; t; \mathbf{x}); \Sigma_{\theta}(\mathbf{z}_t; t; \mathbf{x})); \quad (2)$$

where the mean  $\mu_{\theta}$  and variance  $\Sigma_{\theta}$  are learned from the model. In particular, we follow Li et al. [2022a]’s approach of using predefined variance without trainable parameters.

To extend the continuous diffusion process to discrete text generation, Li et al. [2022a] introduce an additional Markov transition from the discrete tokens  $\mathbf{y}$  to the latent variable  $\mathbf{z}_0$ . In practice, we add an embedding step  $q(\mathbf{z}_0 | \mathbf{y}) = \mathcal{N}(\mathbf{z}_0; \text{Emb}(\mathbf{y}), (1 - \alpha_0) \mathbf{I})$  in the forward process, and use a trainable rounding step which is parametrized by  $p_{\theta}(\mathbf{y} | \mathbf{z}_0; \mathbf{x}) = \prod_{i=1}^N p_{\theta}(y_i | z_0^i; \mathbf{x})$  in the reverse process. In each timestep, we utilize an encoder-decoder model  $\mathbf{g}_{\theta}(\mathbf{z}_t; t; \mathbf{x})$  to approximate  $\mathbf{z}_0$  [Lin et al., 2023] in a NAR manner and then estimate  $\mu_{\theta}(\mathbf{z}_t; t; \mathbf{x})$ .

In consequence, combined with maximizing the evidence lower bound (ELBO) of  $\log p_{\theta}(\mathbf{y}|\mathbf{x})$ , our training objective of the conditional diffusion language model is:

$$L = \mathbb{E}_{q(\mathbf{z}_{0:T}|\mathbf{y})} \left[ \log p(\mathbf{y} | \mathbf{z}_0; \mathbf{x}) + \sum_{t=1}^T k \mathbf{z}_0 \cdot \mathbf{g}(\mathbf{z}_t; t; \mathbf{x}) k^2 \right]; \quad (3)$$

## 3 Methodology

### 3.1 Multi-Level Diffusion

In the typical diffusion process, every token in the text sequence has the same diffusion timestep. In order to leverage the sequential nature of language, we enable tokens to have different diffusion timesteps during the forward and reverse pass. To accomplish this, we propose a multi-level diffusion strategy that includes both sentence-level and token-level diffusion. Firstly, at the sentence-level, we follow Diffusion-LM [Li et al., 2022a] to randomly select a timestep  $t$ . Secondly, at the token-level, we incorporate positional information  $n \in [1, N]$  based on the sentence-level timestep to regulate the diffusion timestep for the current token. The procedure is illustrated as:

$$\mathbf{z}_t = (\mathbf{z}_{f(1;t)}^1; \mathbf{z}_{f(2;t)}^2; \dots; \mathbf{z}_{f(N;t)}^N); \quad (4)$$

where  $N$  is the given target sentence length,  $\mathbf{z}_t$  is the sentence representation at timestep<sup>1</sup>  $t$ ,  $\mathbf{z}_{f(n;t)}^n$  is the latent representation for the  $n$ -th token at sentence-level timestep  $t$ , and  $f(n, t)$  is a token-level timestep function that denotes the token-level diffusion timestep determined by token position  $n$  and sentence-level timestep  $t$ .

We visualize the token-level timestep  $(n, f(n, t))$  onto a two-dimensional coordinate system as Figure 1, which takes the token **position** as the horizontal axis and the sentence-level **timestep** as the

<sup>1</sup>Please note that if we talk about a ‘‘timestep’’ without explicitly indicating that it is for token-level, it should be for sentence-level.

105 vertical axis. Furthermore, to provide a more profound description of the characteristics of movement,  
 106 we define the speed of movement as the following equation.

$$v(n; t_i; t_{i+1}) = f(n; t_{i+1}) - f(n; t_i); \quad (5)$$

107 where  $t_i$  and  $t_{i+1}$  are the start and end sentence-level diffusion timesteps. It can be observed that  
 108 tokens in Diffusion-LM share the same movement speed, while those in AR exhibit different speeds.

### 109 3.2 Token-Level Diffusion with Dynamic Movement Speed

110 Based on the speed of movement, we propose a fundamental principle, dynamic movement speed,  
 111 for designing the token-level diffusion timestep function  $f(n, t)$  to take advantage of AR in diffusion.  
 112 Specifically, elements on the left side of a sentence undergo higher movement speed from random  
 113 Gaussian noise to token embedding, while those on the right side experience lower movement speed,  
 114 thereby they can be generated in the later sentence-level timestep and utilize information from  
 115 previously generated tokens more effectively.

---

#### Algorithm 1 Training Process of AR-DIFFUSION.

---

**Input:** Dataset  $f(\mathbf{x}; \mathbf{y})g$ , maximum timestep number  $T$  and maximum target length  $N$ .

**Output:** Optimized model parameters .

1: Define an anchor point  $(n_e; t_e)^2$ .

2: **repeat**

3: Sample  $(\mathbf{x}; \mathbf{y})$  from the dataset and embed  $\mathbf{y}$  into  $\mathbf{z}_0$ .

4: Sample a sentence-level timestep  $t$  from the interval  $[0; N + T]$ , then the start point is determined by the following equation:

$$(n_s; t_s) = (\text{clip}(N - t; 0; N); \text{clip}(t - N; 0; T)) \quad (6)$$

5: Use the point-slope linear function to determine the token-level timestep  $f(n; t)$  in position  $n$ :

$$f(n; t) = \text{clip}\left(\frac{t_e - t_s}{n_e - n_s}(n - n_s) + t_s; 0; T\right) \quad (7)$$

6: Sample  $\mathbf{z}_{f(n;t)}^n$  for each  $n$  in different positions with Gaussian reparameterization.

7: According to equation (3) and equation (9), employ gradient descent to optimize the objective:

$$\min \left[ -\log p(\mathbf{y} | \mathbf{z}_0; \mathbf{x}) + \sum_{n=1}^N \|\mathbf{g}(\mathbf{z}_{f(n;t)}^n; f(n; t); \mathbf{x}) - \mathbf{z}_0\|^2 \right] \quad (8)$$

8: **until** converged

---

116 Following the guidance of the principle, we develop a token-level diffusion strategy with the  
 117 linear function, which is shown in Figure 1(c). In particular, the procedure is illustrated in  
 118 Algorithm 1, where  $\text{clip}(x, \min, \max)$  function is to clamp all elements in  $x$  into the range  
 119  $[\min, \max]$ . Specifically, in the forward process of diffusion, the start point goes to the left from  
 120  $(N, 0)$  to  $(0, 0)$  along the horizontal axis and then moves up to  $(0, T)$  along the vertical axis.  
 121 Therefore, the entire range of sentence-level timestep is extended to  $[0, N + T]$ .

122 In the reverse diffusion process, the multi-level diffusion follows the formula:

$$\mathbf{g}(\mathbf{z}_t; t; \mathbf{x}) = \mathbf{g}\left(\left(\mathbf{z}_{f(1;t)}^1; f(1; t)\right); \left(\mathbf{z}_{f(2;t)}^2; f(2; t)\right); \dots; \left(\mathbf{z}_{f(N;t)}^N; f(N; t)\right); \mathbf{x}\right); \quad (9)$$

123 where  $\mathbf{g}(\mathbf{z}_{f(n;t)}^n, f(n, t); \mathbf{x})$  denotes the  $n$ -th element.

### 124 3.3 Inference with Skipping

125 Typically, the generation process needs to go through all the sentence-level timesteps from  $T + N$   
 126 to 0. To reduce the decoding time, we introduce a skipping mechanism that allows us to traverse  
 127 a subset of timesteps.

---

<sup>2</sup>In particular, the anchor point is set as  $(2 - N; T)$  in our implementation. The impact of different choices of the anchor point is discussed in supplementary material D.

---

**Algorithm 2** Inference Process of AR-DIFFUSION with the Skipping Mechanism.

---

**Input:** Source condition  $\mathbf{x}$ , number of decoding steps  $M$  and model parameters  $\theta$ .

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

- 1: Define an anchor point  $(n_e; t_e)$ .
- 2: Uniformly select a decreasing sequence of timesteps  $\hat{t}_i, g_{i=0}^M$  ranging from  $T + N$  to 0.
- 3: Sample  $\mathbf{z}_{t_0} \sim N(\mathbf{0}; \mathbf{I})$ .
- 4: **for**  $j = 0$  to  $M - 1$  **do**
- 5:   Calculate the start point  $(n_s; t_s)$  using equation (6).
- 6:   Based on the current sentence-level inference steps  $t_i$  and the next one  $t_{i+1}$ , assign token-level timesteps  $f(n; t_i)$  and  $f(n; t_{i+1})$  to token in position  $n$  using equation (7).
- 7:   Reverse sample  $\mathbf{z}_{t_{i+1}} = (\mathbf{z}_{f(1; t_{i+1})}^1; \mathbf{z}_{f(2; t_{i+1})}^2; \dots; \mathbf{z}_{f(N; t_{i+1})}^N)$  from  $p(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}; \mathbf{x})$  with the following formulas:

$$p(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}; \mathbf{x}) = \prod_{n=1}^N p(\mathbf{z}_{f(n; t_{i+1})}^n | \mathbf{z}_{f(n; t_i)}^n; \mathbf{x}) \quad (10)$$

$$p(\mathbf{z}_{f(n; t_{i+1})}^n | \mathbf{z}_{f(n; t_i)}^n; \mathbf{x}) = N(\mathbf{z}_{f(n; t_{i+1})}^n; \mathbf{z}_{f(n; t_i)}^n + \mathbf{g}(\mathbf{z}_{f(n; t_i)}^n; f(n; t_i); \mathbf{x}); \mathbf{I}) \quad (11)$$

- 8: **end for**
  - 9: Map  $\mathbf{z}_{t_M}$  to the nearest embedding  $\hat{\mathbf{y}}$ .
- 

128 In practice, we propose an algorithm for the inference, illustrated in Algorithm 2.

$$= \frac{\sqrt{\frac{f(n; t_i)}{f(n; t_{i+1})}} (1 - \frac{f(n; t_{i+1})}{f(n; t_i)})}{1 - \frac{f(n; t_{i+1})}{f(n; t_i)}}; = \frac{\sqrt{\frac{f(n; t_{i+1})}{f(n; t_i)}} (1 - \frac{f(n; t_i)}{f(n; t_{i+1})})}{1 - \frac{f(n; t_i)}{f(n; t_{i+1})}}; = \frac{(1 - \frac{f(n; t_{i+1})}{f(n; t_i)}) (1 - \frac{f(n; t_i)}{f(n; t_{i+1})})}{1 - \frac{f(n; t_i)}{f(n; t_{i+1})}} \quad (12)$$

129 In equation (10), the conditional distribution of  $\mathbf{z}_{t_{i+1}}$  is inferred by  $p(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}; \mathbf{x})$ , and then we  
 130 decompose it by positions due to the independent forward process of elements at different positions.  
 131 From equation (11) to equation (12), we establish the relationship between tokens at different  
 132 timesteps, and the detailed derivation can be found in supplementary material E.

## 133 4 Experiments

### 134 4.1 Tasks and Datasets

135 **Text Summarization** This task involves taking a long document as input and generating a concise  
 136 sentence as output. This requires models with the ability to identify important content and rewrite it  
 137 in a condensed form. In our experiments, we use the publicly available XSUM [Narayan et al., 2018]  
 138 and CNN/DAILYMAIL Hermann et al. [2015] on GLGE<sup>3</sup>, which is also named as GLGE-Easy.

139 **Machine Translation** Translation is a widely used sequence-to-sequence task. The input is a  
 140 sequence of words in the source language, and the output is a sequence of corresponding words in the  
 141 target language. We choose the IWSLT 2014 dataset and the data processing method is to follow the  
 142 scripts provided by fairseq<sup>4</sup>.

143 **Common Sense Generation** In this task, the model is provided with a concept set consisting of  
 144 objects and actions as input. The objective is to generate a sentence that incorporates these concepts  
 145 and describes a realistic scenario. We use COMMONGEN<sup>5</sup> dataset for evaluation.

### 146 4.2 Main Results

147 The results on different datasets are shown in Table 1, Table 2, Table 3 and Table 4. The best  
 148 result is **bolded** and the second-best result is underlined. “ $k$ ” indicates the number of generated

<sup>3</sup><https://microsoft.github.io/glge/>

<sup>4</sup><https://github.com/facebookresearch/fairseq/tree/main/examples/translation>

<sup>5</sup><https://inl.ucsc.edu/CommonGen/>

149 candidate samples<sup>6</sup>. It can be seen from the results in each table that AR-DIFFUSION achieves the  
 150 best performance.

151 During the inference process, we utilize **20** inference steps and employ Minimum Bayes Risk (MBR)  
 152 [Kumar and Byrne, 2004] decoding to select the best sample, following [Li et al., 2022a]. We choose  
 153 MBR instead of the selection approach in GENIE, as GENIE picks up the best sample by calculating  
 154 the maximum score for each generated one using ground truth, which introduces unfairness. To  
 155 ensure a fair comparison, we re-implement GENIE using our configuration and perform inference  
 156 with 20 steps. More experimental details can be found in the supplementary material B.

Table 1: Results on XSUM test set. The results of NAR and Semi-NAR are from Qi et al. [2021], and the results of AR are from GLGE [Liu et al., 2021].

Methods	Pattern	ROUGE-1	ROUGE-2	ROUGE-L
NAT [Gu et al., 2017]	NAR	24.0	3.9	20.3
iNAT [Lee et al., 2018]		24.0	4.0	20.4
CMLM [Ghazvininejad et al., 2019]		23.8	3.6	20.2
LevT [Gu et al., 2019]		24.8	4.2	20.9
InsT [Stern et al., 2019]	Semi-NAR	17.7	5.2	16.1
iNAT [Lee et al., 2018]		27.0	6.9	22.4
CMLM [Ghazvininejad et al., 2019]		29.1	7.7	23.0
LevT [Gu et al., 2019]		25.3	7.4	21.5
LSTM [Greff et al., 2017]	AR <sup>7</sup>	25.1	6.9	19.9
Transformer [Vaswani et al., 2017]		30.5	<u>10.4</u>	24.2
GENIE [Lin et al., 2023] ( $k = 50$ )	Diffusion	29.3	8.3	21.9
AR-DIFFUSION ( $k = 50$ )		<u>31.7</u>	10.1	<u>24.7</u>
AR-DIFFUSION ( $k = 500$ )		<b>32.2</b>	<b>10.6</b>	<b>25.2</b>

Table 2: Results on CNN/DAILYMAIL test set. The results of AR are from GLGE Liu et al. [2021].

Methods	Pattern	ROUGE-1	ROUGE-2	ROUGE-L
LSTM [Greff et al., 2017]	AR	37.3	15.7	34.4
Transformer [Vaswani et al., 2017]		39.5	<u>16.7</u>	36.7
GENIE [Lin et al., 2023] ( $k = 50$ )	Diffusion	34.4	12.8	32.1
AR-DIFFUSION ( $k = 50$ )		<u>39.6</u>	16.3	<u>37.1</u>
AR-DIFFUSION ( $k = 500$ )		<b>40.2</b>	<b>17.1</b>	<b>37.7</b>

157 **Text Summarization** The results presented in Table 1 and Table 2 clearly demonstrate that AR-  
 158 DIFFUSION outperforms the existing NAR and Semi-NAR approaches across all metrics. Moreover,  
 159 AR-DIFFUSION consistently achieves significant improvements over GENIE in terms of all indicators.  
 160 Furthermore, in comparison to Transformer, AR-DIFFUSION outperforms it on both ROUGE-1 and  
 161 ROUGE-L, while achieving comparable performance in terms of ROUGE-2. Notably, when the sample  
 162 number is 500, AR-DIFFUSION demonstrates superiority over Transformer across all the measures.

163 **Machine Translation** Table 3 presents the BLEU score implemented by SeqDiffuSeq setting<sup>8</sup>.  
 164 AR-DIFFUSION outperforms the non-auto-regressive CNAT in greedy search for a single sample, and  
 165 achieves a substantial gain. Moreover, the BLEU score of AR-DIFFUSION surpasses GENIE by a large  
 166 margin and shows a slightly better performance than the AR Transformer. Specially, AR-DIFFUSION  
 167 achieves a more powerful result at  $k = 500$ .

168 **Common Sense Generation** As depicted in Table 4, AR-DIFFUSION achieves superior perfor-  
 169 mance compared to the current AR, NAR, and other diffusion methods across all the metrics on the  
 170 COMMONGEN dataset.

<sup>6</sup>The relationship between sample number and results is discussed in supplementary material D.

<sup>7</sup>Notably, although AR’s beam search has a small beam, the search space may be larger than 50 or even 500.

<sup>8</sup>We also report SacreBLEU in supplementary material C to compare with DINOISER.

Table 3: Results on IWSLT14 DE-EN test set following the setting of SEQDIFFUSeq. “NFE” indicates the Number of Function Evaluations [Ye et al., 2023].

Methods	Pattern	BLEU	Steps	NFE (Steps $k$ )
Transformer [Vaswani et al., 2017]	AR	34.74	-	-
CNAT [Bao et al., 2021]	NAR	29.81	-	-
SeqDiffuSeq [Yuan et al., 2022] ( $k = 1$ )	Diffusion	29.83	2,000	2,000 (2,000 1)
AR-DIFFUSION ( $k = 1$ )		30.19	20	20 (20 1)
GENIE [Lin et al., 2023] ( $k = 50$ )	Diffusion	30.08	20	1,000 (20 50)
AR-DIFFUSION ( $k = 50$ )		34.95	20	1,000 (20 50)
AR-DIFFUSION ( $k = 500$ )		<b>35.62</b>	20	10,000 (20 500)

Table 4: Results on COMMONGEN dev set. Results of NAR and AR are from Lin et al. [2020].

Methods	Pattern	ROUGE-2/L		BLEU-3/4		METEOR	SPICE
bRNN-CopyNet [Gu et al., 2016]	AR	9.23	30.57	13.60	7.80	17.40	16.90
Trans-CopyNet [Lin et al., 2020]		11.08	32.57	17.20	10.60	18.80	18.00
MeanPooling-CopyNet [Lin et al., 2020]		11.36	34.63	14.80	8.90	19.20	20.20
LevT [Gu et al., 2019]	NAR	12.22	<u>35.42</u>	<u>23.10</u>	<u>15.00</u>	22.10	21.40
ConstLeven [Susanto et al., 2020]		<u>13.47</u>	35.19	21.30	12.30	<b>25.00</b>	<u>23.20</u>
GENIE [Lin et al., 2023] ( $k = 50$ )	Diffusion	12.89	35.21	22.00	13.30	<u>24.30</u>	23.00
AR-DIFFUSION ( $k = 50$ )		<b>13.93</b>	<b>37.36</b>	<b>25.60</b>	<b>16.40</b>	<b>25.00</b>	<b>24.20</b>

### 171 4.3 Inference Efficiency

172 First, we use the number of function evaluations (NFE) as a measure to compare inference effi-  
 173 ciency [Ye et al., 2023] in machine translation. From Table 3, it is evident that even when the NFE  
 174 is reduced to 1% of SeqDiffuSeq (equivalent to  $100\times$  faster), AR-DIFFUSION still outperforms  
 175 SeqDiffuSeq. Moreover, increasing the number of generated candidate samples ( $k = 500$ ) leads to  
 176 further performance improvements, albeit with increased time consumption.

177 Second, we conduct experiments with an **extremely limited number of inference steps** (2 and 3)<sup>9</sup>  
 178 and compare the performance with that of GENIE in XSUM. The results are presented in Table 5.  
 179 When reducing the number of steps to 2, GENIE experiences a significant decline, with an average  
 180 score of 4.20 in the AVG Drop column, while AR-DIFFUSION exhibits a comparatively smaller  
 181 decrease of 1.34. Furthermore, with 3 steps, although the performance deterioration of GENIE is  
 182 somewhat reduced, the average score still shows a decline of 2.81. In contrast, AR-DIFFUSION  
 183 maintains a high performance level, with an average score differing from the 20-step result by only  
 184 0.64. Notably, the results of AR-DIFFUSION at 3 steps are comparable to the results of GENIE  
 185 at 2,000 steps. Therefore, compared to GENIE, the inference speed of AR-DIFFUSION can be  
 186 accelerated by up to  $600\times$ .

### 187 4.4 Analysis

188 **Diversity of Samples** Diversity is a key advantage of diffusion models. To measure the diversity  
 189 of generated samples, We adopt the SELF-BLEU [Zhu et al., 2018] metric, in which a lower score  
 190 indicates higher diversity. In Lin et al. [2023], various sampling methods were applied to the pre-  
 191 trained auto-regressive model BART<sup>10</sup>. As shown in Table 6, AR-DIFFUSION achieves significantly  
 192 higher diversity compared to the auto-regressive model. Furthermore, the diversity can be comparable  
 193 to GENIE with a better performance.

194 **Ablation Study** To demonstrate the effectiveness of our proposed method, we perform ablation  
 195 experiments on the XSUM dataset. Our results show that both our proposed multi-level diffusion and  
 196 skipping mechanism are essential for achieving the high performance of AR-DIFFUSION.

<sup>9</sup>The time consumed by each step in the inference process is exactly the same.

<sup>10</sup>Specific details are written in supplementary material F.

Table 5: Experimental results of GENIE and AR-DIFFUSION with inference steps of **2** and **3** on XSUM test set. Take  $k = 10$  to apply the MBR decoding strategy. (·) indicates the **drop** score compared to the 20-step.

Methods	Steps	NFE	ROUGE-1	ROUGE-2	ROUGE-L	AVG Drop
GENIE	2,000	20,000	30.36	8.78	23.31	-
	20	200	28.33	7.46	21.15	-
	3	30	25.03 (-3.30)	5.32 (-2.14)	18.17 (-2.98)	2.81
	2	20	23.45 (-4.88)	3.95 (-3.51)	16.94 (-4.21)	4.20
AR-DIFFUSION	20	200	30.99	9.32	23.95	-
	3	30	30.23 (-0.76)	8.68 (-0.64)	23.43 (-0.52)	<b>0.64</b>
	2	20	29.28 (-1.71)	7.99 (-1.33)	22.98 (-0.97)	<b>1.34</b>

Table 6: Diversity of **10** generated samples on XSUM test set and average of **10** results. The results of BART and GENIE are quoted from Lin et al. [2023].

Methods	BART						GENIE	AR-DIFFUSION
Sampling	Greedy Search	Beam Search	Diverse Beam Search	Typical Sample	Top-k Sample	Nucleus Sample	Diffusion	
SELF-BLEU ↓	100.0	93.4	75.6	76.9	80.2	79.1	29.3	30.4

197 Maintaining the skipping inference method, we remove the token-level diffusion during the training  
 198 process, which degenerates to GENIE w/ skipping. The comparison results are shown in Figure 2(a).  
 199 It can be observed that after removing, the AVG-ROUGE score is greatly lower after 2 steps.

200 The performance of applying our proposed skipping mechanism and DDIM [Song et al., 2021] to  
 201 AR-DIFFUSION is shown in Figure 2(b). The results demonstrate that the skipping mechanism  
 202 consistently outperforms DDIM in various inference steps. Additionally, the skipping mechanism  
 203 can be easily applied to GENIE. As depicted in Figure 2(c), DDIM suffers a significant drop in  
 204 performance when the number of inference steps is less than 40. In contrast, the skipping mechanism  
 205 consistently maintains good performance across all inference steps.

206 **Case Study** By mapping the state to the token with the highest logits, we visualize the intermediate  
 207 states of AR-DIFFUSION. As depicted in Figure 3, AR-DIFFUSION undergoes a denoising process,  
 208 transforming the random Gaussian noise into a coherent sentence over 20 steps, and we present 5 of  
 209 them. With the progression of each timestep, compared to the tokens on the right side of the sentence,  
 210 the tokens on the left side demonstrate faster determination and a rapid increase in the corresponding  
 211 logits. This behavior is consistent with our principle of dynamic movement speed from left to right.

## 212 5 Related Work

213 **AR and NAR Language Models** AR models have been the dominant approach for text generation  
 214 [OpenAI, 2023, Touvron et al., 2023, Dong et al., 2023], but their token-by-token generation nature  
 215 often leads to unsatisfactory inference speed. To address this issue, NAR models have been developed  
 216 in recent years. The NAR method is initially proposed by Gu et al. [2017], its objective is generate the  
 217 entire output sequence in parallel, thereby improving generation speed and efficiency. Subsequently,  
 218 LevT [Gu et al., 2019] adopts insertion and deletion to address the lack of flexibility in NAR  
 219 generation, CMLM [Ghazvininejad et al., 2019] utilizes a masked language model to improve the  
 220 quality of NAR generation through a constant number of iterations, and CNAT [Bao et al., 2021]  
 221 introduces latent variables to represent the category information of the target word to make full use  
 222 of the latent representation. However, these NAR methods are hard to model inter-token position  
 223 dependency and deficient in generation performance.

224 **Continuous Text Diffusion** The application of diffusion models to continuous text space is first  
 225 introduced by Li et al. [2022a]. Through the embedding and rounding processes, the direct integration



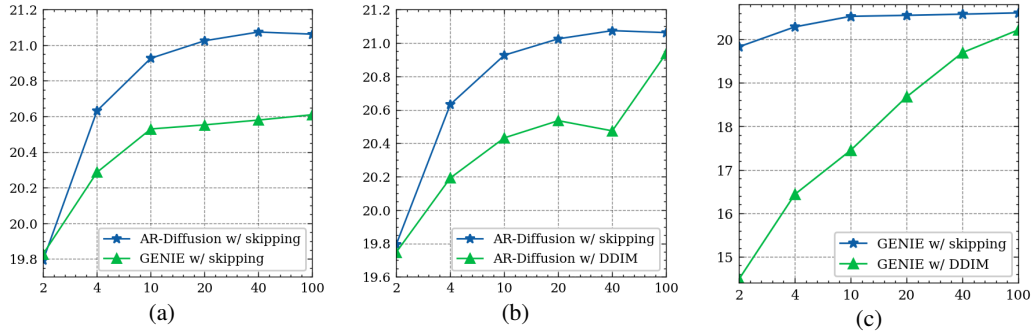


Figure 2: Ablation experiments on XSUM test set and taking  $k = 5$ . The horizontal axis is the number of inference steps and the vertical axis is  $\text{AVG-ROUGE} = (\text{ROUGE-1} + \text{ROUGE-2} + \text{ROUGE-L}) / 3$ .

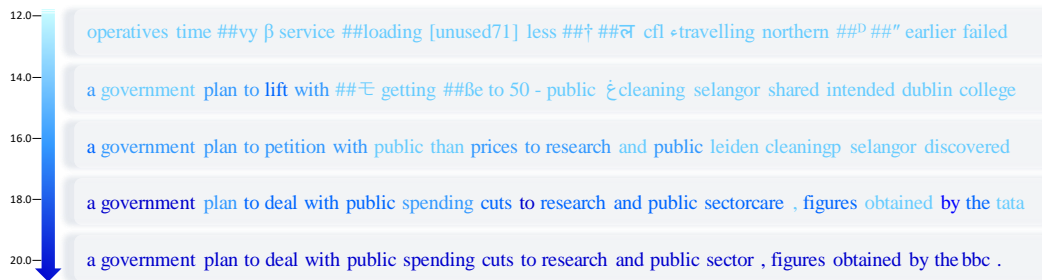


Figure 3: The intermediate state of AR-DIFFUSION gradually generating real text from a standard Gaussian noise through 20 steps. The brightness of the color represents the magnitude of the logits, with darker colors indicating larger logits. More cases are shown in the supplementary materials G.

226 of continuous noise into word embeddings was accomplished. After that, more people attempt to  
 227 adopt continuous text diffusion model to solve sequence-to-sequence tasks. DiffuSeq [Gong et al.,  
 228 2022] divides the input into two parts, utilizing one part as a condition, and perturbs the other part  
 229 with noise. CDCD [Dieleman et al., 2022] proposes score interpolation and time warping to allow  
 230 diffusion model and Euclidean embedding to share the same loss function for training. SeqDiffuSeq  
 231 [Yuan et al., 2022], GENIE [Lin et al., 2023] and DINOISER [Ye et al., 2023] incorporate diffusion  
 232 model into the encoder-decoder structure through cross-attention mechanisms.

233 It is important to highlight the differences between our method and both ARDMs [Hoogeboom  
 234 et al., 2022] and TimeGrad [Rasul et al., 2021], despite the common references to autoregression  
 235 and diffusion in all these. ARDMs employ an order-agnostic technique, leveraging masking and  
 236 prediction for generation in arbitrary orders. On the other hand, TimeGrad integrates RNN and  
 237 diffusion to model the conditional distribution of future steps of multivariate time series. In contrast,  
 238 our research focuses on implementing the diffusion process within a continuous embedding space,  
 239 with the primary aim of generating text in a left-to-right sequence.

## 240 6 Conclusion

241 This paper introduces AR-DIFFUSION, which exhibits AR-like generation behavior but enables  
 242 efficient parallel decoding. Embracing the inherent sequential nature of language, we propose a multi-  
 243 level diffusion model, consisting of sentence-level and token-level components, to assign dynamic  
 244 movement speeds to tokens. Consequently, compared to those on the right, the left tokens undergo  
 245 fewer denoising steps and generate earlier to subsequently influence the later ones. Furthermore,  
 246 we introduce a skipping mechanism to facilitate parallel generation within the multi-level diffusion  
 247 framework. The experimental results across various tasks demonstrate that AR-DIFFUSION surpasses  
 248 existing diffusion models in terms of quality while maintaining diversity. Additionally, compared  
 249 to existing diffusion language models, AR-DIFFUSION achieves comparable results while being  
 250  $100\times \sim 600\times$  faster.

## 251 References

- 252 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/arXiv.2303.08774.  
253 URL <https://doi.org/10.48550/arXiv.2303.08774>.
- 254 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
255 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand  
256 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language  
257 models. *CoRR*, abs/2302.13971, 2023. doi: 10.48550/arXiv.2302.13971. URL <https://doi.org/10.48550/arXiv.2302.13971>.
- 259 Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy  
260 Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model.  
261 [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023.
- 262 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
263 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von  
264 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors,  
265 *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,  
266 Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 268 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
269 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
270 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 271 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In  
272 Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-  
273 Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Con-  
274 ference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-  
275 12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html>.
- 277 Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori Hashimoto. Diffusion-  
278 LM improves controllable text generation. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave,  
279 and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022a. URL  
280 <https://openreview.net/forum?id=3s9IrEsjLyk>.
- 281 Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and Lingpeng Kong. Diffuseq: Sequence to  
282 sequence text generation with diffusion models. *CoRR*, abs/2210.08933, 2022. doi: 10.48550/  
283 arXiv.2210.08933. URL <https://doi.org/10.48550/arXiv.2210.08933>.
- 284 Sander Dieleman, Laurent Sartran, Arman Roshannai, Nikolay Savinov, Yaroslav Ganin, Pierre H  
285 Richemond, Arnaud Doucet, Robin Strudel, Chris Dyer, Conor Durkan, et al. Continuous diffusion  
286 for categorical data. *arXiv preprint arXiv:2211.15089*, 2022.
- 287 Hongyi Yuan, Zheng Yuan, Chuanqi Tan, Fei Huang, and Songfang Huang. Seqdiffuseq: Text  
288 diffusion with encoder-decoder transformers, 2022.
- 289 Jiasheng Ye, Zaixiang Zheng, Yu Bao, Lihua Qian, and Mingxuan Wang. Dinoiser: Diffused  
290 conditional sequence learning by manipulating noises, 2023.
- 291 Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy,  
292 Veselin Stoyanov, and Luke Zettlemoyer. BART: denoising sequence-to-sequence pre-training  
293 for natural language generation, translation, and comprehension. In Dan Jurafsky, Joyce Chai,  
294 Natalie Schluter, and Joel R. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the  
295 Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880.  
296 Association for Computational Linguistics, 2020. doi: 10.18653/v1/2020.acl-main.703. URL  
297 <https://doi.org/10.18653/v1/2020.acl-main.703>.
- 298 Weizhen Qi, Yu Yan, Yeyun Gong, Dayiheng Liu, Nan Duan, Jiusheng Chen, Ruofei Zhang, and  
299 Ming Zhou. ProphetNet: Predicting future n-gram for sequence-to-SequencePre-training. In  
300 *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2401–2410,

- 301 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
302 findings-emnlp.217. URL <https://aclanthology.org/2020.findings-emnlp.217>.
- 303 Weizhen Qi, Yeyun Gong, Jian Jiao, Yu Yan, Weizhu Chen, Dayiheng Liu, Kewen Tang, Houqiang  
304 Li, Jiusheng Chen, Ruofei Zhang, et al. Bang: Bridging autoregressive and non-autoregressive  
305 generation with large scale pretraining. In *International Conference on Machine Learning*, pages  
306 8630–8639. PMLR, 2021.
- 307 Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. ELMER: A non-  
308 autoregressive pre-trained language model for efficient and effective text generation. In Yoav  
309 Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on*  
310 *Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab*  
311 *Emirates, December 7-11, 2022*, pages 1044–1058. Association for Computational Linguistics,  
312 2022b. URL <https://aclanthology.org/2022.emnlp-main.68>.
- 313 Yafu Li, Leyang Cui, Yongjing Yin, and Yue Zhang. Multi-granularity optimization for non-  
314 autoregressive translation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceed-*  
315 *ings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022,*  
316 *Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 5073–5084. Association for Com-  
317 putational Linguistics, 2022c. URL <https://aclanthology.org/2022.emnlp-main.339>.
- 318 Yu Bao, Shujian Huang, Tong Xiao, Dongqi Wang, Xinyu Dai, and Jiajun Chen. Non-autoregressive  
319 translation by learning target categorical codes. In *Proceedings of the 2021 Conference of the North*  
320 *American Chapter of the Association for Computational Linguistics: Human Language Technolo-*  
321 *gies*, pages 5749–5759, Online, June 2021. Association for Computational Linguistics. doi: 10.  
322 18653/v1/2021.naacl-main.458. URL <https://aclanthology.org/2021.naacl-main.458>.
- 323 Zhenghao Lin, Yeyun Gong, Yelong Shen, Tong Wu, Zhihao Fan, Chen Lin, Nan Duan, and Weizhu  
324 Chen. Text generation with diffusion language models: A pre-training approach with continuous  
325 paragraph denoise, 2023.
- 326 Jiatao Gu, James Bradbury, Caiming Xiong, Victor OK Li, and Richard Socher. Non-autoregressive  
327 neural machine translation. *arXiv preprint arXiv:1711.02281*, 2017.
- 328 Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don’t give me the details, just the summary!  
329 topic-aware convolutional neural networks for extreme summarization. In *Proceedings of the*  
330 *2018 Conference on Empirical Methods in Natural Language Processing*, pages 1797–1807,  
331 Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi:  
332 10.18653/v1/D18-1206. URL <https://aclanthology.org/D18-1206>.
- 333 Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Su-  
334 leyman, and Phil Blunsom. Teaching machines to read and comprehend. In C. Cortes, N. Lawrence,  
335 D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,  
336 volume 28. Curran Associates, Inc., 2015. URL [https://proceedings.neurips.cc/paper\\_](https://proceedings.neurips.cc/paper_files/paper/2015/file/afdec7005cc9f14302cd0474fd0f3c96-Paper.pdf)  
337 [files/paper/2015/file/afdec7005cc9f14302cd0474fd0f3c96-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/afdec7005cc9f14302cd0474fd0f3c96-Paper.pdf).
- 338 Shankar Kumar and William Byrne. Minimum bayes-risk decoding for statistical machine translation.  
339 Technical report, JOHNS HOPKINS UNIV BALTIMORE MD CENTER FOR LANGUAGE AND  
340 SPEECH PROCESSING (CLSP), 2004.
- 341 Dayiheng Liu, Yu Yan, Yeyun Gong, Weizhen Qi, Hang Zhang, Jian Jiao, Weizhu Chen, Jie Fu,  
342 Linjun Shou, Ming Gong, et al. Glge: A new general language generation evaluation benchmark.  
343 In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 408–420,  
344 2021.
- 345 Jason Lee, Elman Mansimov, and Kyunghyun Cho. Deterministic non-autoregressive neural sequence  
346 modeling by iterative refinement. *arXiv preprint arXiv:1802.06901*, 2018.
- 347 Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. Mask-predict: Parallel  
348 decoding of conditional masked language models. In *Proceedings of the 2019 Conference on*  
349 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*  
350 *on Natural Language Processing (EMNLP-IJCNLP)*, pages 6112–6121, Hong Kong, China,  
351 November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1633. URL  
352 <https://aclanthology.org/D19-1633>.

- 353 Jiatao Gu, Changhan Wang, and Junbo Zhao. Levenshtein transformer. In *Advances in Neural*  
354 *Information Processing Systems*, pages 11181–11191, 2019.
- 355 Mitchell Stern, William Chan, Jamie Kiros, and Jakob Uszkoreit. Insertion transformer: Flexible  
356 sequence generation via insertion operations. *arXiv preprint arXiv:1902.03249*, 2019.
- 357 Klaus Greff, Rupesh Kumar Srivastava, Jan Koutník, Bas R. Steunebrink, and Jürgen Schmidhuber.  
358 LSTM: A search space odyssey. *IEEE Trans. Neural Networks Learn. Syst.*, 28(10):2222–2232,  
359 2017.
- 360 Bill Yuchen Lin, Wangchunshu Zhou, Ming Shen, Pei Zhou, Chandra Bhagavatula, Yejin Choi, and  
361 Xiang Ren. CommonGen: A constrained text generation challenge for generative commonsense  
362 reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages  
363 1823–1840, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/  
364 2020.findings-emnlp.165. URL <https://aclanthology.org/2020.findings-emnlp.165>.
- 365 Jiatao Gu, Zhengdong Lu, Hang Li, and Victor O. K. Li. Incorporating copying mechanism in  
366 sequence-to-sequence learning. In *ACL (1)*. The Association for Computer Linguistics, 2016.
- 367 Raymond Hendy Susanto, Shamil Chollampatt, and Liling Tan. Lexically constrained neural machine  
368 translation with levenshtein transformer. In *ACL*, pages 3536–3543. Association for Computational  
369 Linguistics, 2020.
- 370 Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. Tegygen:  
371 A benchmarking platform for text generation models. In *The 41st international ACM SIGIR*  
372 *conference on research & development in information retrieval*, pages 1097–1100, 2018.
- 373 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *9th*  
374 *International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May*  
375 *3-7, 2021*. OpenReview.net, 2021. URL <https://openreview.net/forum?id=St1giarCHLP>.
- 376 Chenhe Dong, Yinghui Li, Haifan Gong, Miaoxin Chen, Junxin Li, Ying Shen, and Min Yang. A  
377 survey of natural language generation. *ACM Comput. Surv.*, 55(8):173:1–173:38, 2023. doi:  
378 10.1145/3554727. URL <https://doi.org/10.1145/3554727>.
- 379 Emiel Hoogeboom, Alexey A. Gritsenko, Jasmijn Bastings, Ben Poole, Rianne van den Berg,  
380 and Tim Salimans. Autoregressive diffusion models. In *International Conference on Learning*  
381 *Representations*, 2022. URL <https://openreview.net/forum?id=Lm8T39vLDTE>.
- 382 Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. Autoregressive denoising  
383 diffusion models for multivariate probabilistic time series forecasting. In Marina Meila and Tong  
384 Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume  
385 139 of *Proceedings of Machine Learning Research*, pages 8857–8868. PMLR, 18–24 Jul 2021.  
386 URL <https://proceedings.mlr.press/v139/rasul21a.html>.
- 387 Taku Kudo and John Richardson. SentencePiece: A simple and language independent subword  
388 tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference*  
389 *on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71,  
390 Brussels, Belgium, November 2018. Association for Computational Linguistics. doi: 10.18653/v1/  
391 D18-2012. URL <https://aclanthology.org/D18-2012>.
- 392 Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*,  
393 2022.
- 394 Yisheng Xiao, Lijun Wu, Junliang Guo, Juntao Li, Min Zhang, Tao Qin, and Tie-Yan Liu. A survey on  
395 non-autoregressive generation for neural machine translation and beyond. *CoRR*, abs/2204.09269,  
396 2022.
- 397 Ashwin K. Vijayakumar, Michael Cogswell, Ramprasaath R. Selvaraju, Qing Sun, Stefan Lee,  
398 David J. Crandall, and Dhruv Batra. Diverse beam search: Decoding diverse solutions from neural  
399 sequence models. *CoRR*, abs/1610.02424, 2016.
- 400 Clara Meister, Tiago Pimentel, Gian Wihler, and Ryan Cotterell. Typical decoding for natural language  
401 generation. *CoRR*, abs/2202.00666, 2022.

- 402 Angela Fan, Mike Lewis, and Yann N. Dauphin. Hierarchical neural story generation. In *ACL (1)*,  
403 pages 889–898. Association for Computational Linguistics, 2018.
- 404 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text  
405 degeneration. In *ICLR*. OpenReview.net, 2020.

## 406 A Limitation

407 A primary limitation of our work lies in the requirement of generating a large number of candidate  
408 samples for optimal performance. As an illustration in Table 2 of CNN/DAILYMAIL dataset, AR-  
409 DIFFUSION ( $k = 50$ ) achieves a 0.8 lower ROUGE-2 score compared to AR-DIFFUSION ( $k = 500$ ).  
410 We anticipate exploring more efficient sampling strategies to minimize the number of generated  
411 samples without performance drop.

## 412 B Experimental Details

413 **Model Setup** Our model configuration is implemented based on Transformer-base [Vaswani  
414 et al., 2017]. In particular, For XSUM and CNN/DAILYMAIL, we set the diffusion embedding  
415 dimension to 128. For IWSLT14, we use 64-dimensional diffusion embedding, 4 attention heads  
416 and 1024-dimensional feed-forward layers. For COMMONGEN, we adopt 64-dimensional diffusion  
417 embedding, 8 attention heads and 512-dimensional feed-forward layers.

418 **Training and Inference** In the training phase, we employ a square-root noise schedule and 2,000  
419 diffusion steps [Li et al., 2022a]. Specially, we use the tokenizer and vocabulary constructed by Byte  
420 Pair Encoding (BPE)<sup>11</sup> [Kudo and Richardson, 2018] for translation tasks. For other tasks, we adopt  
421 the tokenizer and vocabulary of bert-base-uncased.

422 **Baselines** We set four groups of baselines:

423 • NAR: NAT [Gu et al., 2017], iNAT [Lee et al., 2018], CMLM [Ghazvininejad et al., 2019], LevT  
424 [Gu et al., 2019] and CNAT [Bao et al., 2021];

425 • Semi-NAR: InsT [Stern et al., 2019], iNAT [Lee et al., 2018], CMLM [Ghazvininejad et al., 2019]  
426 and LevT [Gu et al., 2019];

427 • AR: bRNN [Gu et al., 2016], LSTM [Greff et al., 2017] and Transformer [Vaswani et al., 2017];

428 • Diffusion: DiffusionLM [Li et al., 2022a], CDCD [Dieleman et al., 2022], SeqDiffuSeq [Yuan  
429 et al., 2022], DINOISER [Ye et al., 2023] and GENIE [Lin et al., 2023].

430 **Metrics** We follow the approach of Qi et al. [2020]<sup>12</sup> to evaluate the **ROUGE-1/2/L** of the  
431 summarization task. For the evaluation of translation tasks, we adopt the setting of SeqDiffuSeq [Yuan  
432 et al., 2022] to report BLEU score. In addition, we also calculate the SacreBLEU score according  
433 to the setting of DINOISER [Ye et al., 2023] for comparison. For COMMONGEN, we employ  
434 ROUGE-2/L, BLEU-3/4, METEOR and SPICE under the evaluation methods of Lin et al. [2020]<sup>13</sup>.

435 **Training Parameters** Our training parameters on different datasets are shown in Table 7. Our  
436 linear schedule warm up steps is  $4,000 \times N_{gc}$ , where  $N_{gc}$  denotes gradient accumulation number. In  
437 addition, we use the AdamW (weight decay = 0.0) optimizer and dropout is 0.2. All experiments are  
438 implemented on 8 Tesla V100-32G. It takes about 20 hours to train XSUM and CNN/DAILYMAIL,  
439 about 5 hours to train IWSLT14, and about 2 hours to train COMMENGEN.

## 440 C More Results on IWSLT14

441 In Table 8 we give the SacreBLEU score according to the setting of DINOISER. AR-DIFFUSION  
442 has notable improvements over non-auto-regressive CMLM. Besides, AR-DIFFUSION achieves  
443 excellent performance among text diffusion models for both EN DE and DE EN tasks. Specifically,  
444 AR-DIFFUSION is far superior to GENIE and comparable to the newly proposed DINOISER at  $n =$   
445 50. Nevertheless, the performance is stronger than DINOISER when  $k = 500$ <sup>14</sup>.

<sup>11</sup>We train bpe on the training set, and follow the vocabulary size of fairseq, IWSLT14 is set to 10,000.

<sup>12</sup>[https://github.com/microsoft/ProphetNet/tree/master/GLGE\\_baselines](https://github.com/microsoft/ProphetNet/tree/master/GLGE_baselines)

<sup>13</sup>[https://github.com/INL-USC/CommonGen/tree/master/evaluation/Traditional/eval\\_metrics](https://github.com/INL-USC/CommonGen/tree/master/evaluation/Traditional/eval_metrics)

<sup>14</sup>DINOISER has shown in their Figure 4 that their method is not better with a larger  $k$ .

Table 7: Training Parameter Settings. Batch Size = mini batch size  $\times N_{gc} \times$  GPU number, Optimized Steps = total steps /  $N_{gc}$ , and  $N_{gc}$  is gradient accumulation number.

Dataset	Lr & Schedule	Batch Size	Optimized Steps	Target Length
XSUM	8e-4 & Cosine	128 $\times$ 3 $\times$ 8	80,000 / 3	50
CNN/DAILYMAIL	8e-4 & Cosine	80 $\times$ 5 $\times$ 8	100,000 / 5	180
IWSLT14 DE EN	2e-3 & Cosine	192 $\times$ 2 $\times$ 8	160,000 / 2	90
IWSLT14 EN DE	1.8e-3 & Cosine	768 $\times$ 1 $\times$ 8	60,000	90
COMMONGEN	3e-4 & Constant	384 $\times$ 1 $\times$ 8	40,000	54

Table 8: SacreBLEU on the IWSLT14 test set. This result follows the setting of DiNOISER.

Methods	IWSLT14	
	DE $\rightarrow$ EN	EN $\rightarrow$ DE
Transformer (AR, beam = 5) [Vaswani et al., 2017]	33.61	28.30
CMLM (NAR, $k = 5$ ) [Ghazvininejad et al., 2019]	29.41	24.34
DiffusionLM ( $k = 50$ ) [Li et al., 2022a]	29.11	22.91
DINOISER ( $k = 50$ ) [Ye et al., 2023]	31.61	<u>26.14</u>
GENIE ( $k = 50$ ) [Lin et al., 2023]	29.45	23.89
AR-DIFFUSION ( $k = 50$ )	<u>31.80</u>	26.01
AR-DIFFUSION ( $k = 500$ )	<b>32.35</b>	<b>26.51</b>

## 446 D Impact of Minimum Bayes Risk and Anchor Point

447 **Minimum Bayes Risk** To investigate the relationship between the number of generated candidate  
 448 samples ( $k$ ) and the quality of generation, we generate varying numbers of samples, ranging up to  
 449 1,000, on the IWSLT14 De EN test set and present the results in Figure 4. The curve demonstrates  
 450 an initial gain of approximately 0.5 SacreBLEU within the first 200 samples, after which the gain  
 451 becomes insignificant with generating more samples.

452 **Anchor Point** We conduct experiments on AR-DIFFUSION using different anchor points ( $n_e, t_e$ ).  
 453 These anchor points vary in terms of  $n_e$  values, namely  $1.0 \times N$ ,  $2.0 \times N$  and  $3.0 \times N$ , where  $N$   
 454 denotes the target sentence length. Additionally, they share a common  $t_e$  value of  $T$ , which represents  
 455 the total time step of diffusion. We present the results in Table 9, and determine that the best result is  
 456 achieved at  $(n_e, t_e) = (2.0 \times N, T)$ .

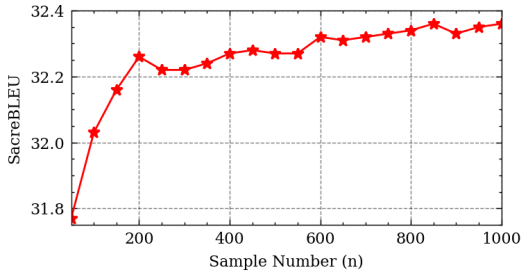


Figure 4: Relationship between the number of candidate samples for applying MBR and SacreBLEU on IWSLT14 DE EN test set.

Table 9: Effect of anchor point at different positions on the IWSLT14 DE EN test set. “N” indicates the target sequence length and “T” represents the total time step of diffusion.

$n_e$	$t_e$	SacreBLEU
1.0	$N$ $T$	31.23
2.0	$N$ $T$	<b>31.80</b>
3.0	$N$ $T$	31.58

## 457 E Proof of Inference with Skipping

458 During the inference process, skipping strategy requires the model  $g$  to infer the state  $Z_{t_{i+1}}^{\eta_2}$  at a  
 459 far-off timestep  $t_{i+1}$  compared to the current state  $Z_{t_i}^{\eta_2}$ , where  $t_{i+1} \ll t_i$ . In our model, due to the

460 dynamic speed setting, token  $\mathbf{z}_{t_{i+1}}^{n_1}$  with smaller timestep  $t_{i+1} \leq t_i$ , which is closer to  $t_{i+1}$ , and  
 461 positions  $n_1 \leq n_2$  can provide stronger auxiliary information than  $\mathbf{z}_{t_i}^{n_1}$ . This reduces the difficulty of  
 462 inferring states for tokens in the end, making our multi-level diffusion model particularly suitable for  
 463 accelerating the generation process.

464 Through maximizing the evidence lower bound (ELBO) of  $p(\mathbf{z}_0)$ , the training object is equivalent to  
 465 minimize the divergence between  $q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{z}_0)$  and  $p(\mathbf{z}_{t-1} | \mathbf{z}_t)$  following [Luo, 2022].

466 By converting the joint probability distribution into a conditional probability distribution, we obtain  
 467 the following formula for  $q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}, \mathbf{z}_0)$ .

$$\begin{aligned}
 q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}, \mathbf{z}_0) &= q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i-1}, \mathbf{z}_{t_i}, \mathbf{z}_0) q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}, \mathbf{z}_0) \\
 &= q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i-1}, \mathbf{z}_0) q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}, \mathbf{z}_0) \\
 &= q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i-2}, \mathbf{z}_0) q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i-1}, \mathbf{z}_0) q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}, \mathbf{z}_0) \\
 &= \prod_{k=1}^{t_i - t_{i+1}} q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i-k}, \mathbf{z}_0)
 \end{aligned} \tag{13}$$

468 Similarly, we reach the same conclusion regarding  $p(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i})$ .

469 Based on equation (13), which consists of  $q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{z}_0)$ , and the interchangeability between  
 470  $q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{z}_0)$  and  $p(\mathbf{z}_{t-1} | \mathbf{z}_t)$ , we can decompose  $q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}, \mathbf{z}_0)$  by incorporating  $\mathbf{z}_{t_i}$  and  $\mathbf{z}_0$ ,  
 471 and utilize our estimated  $\mathbf{z}_0$  to determine the expression of  $p(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i})$ .

$$q(\mathbf{z}_{t_{i+1}} | \mathbf{z}_{t_i}, \mathbf{z}_0) = \prod_{n=1}^N q(z_{f(n;t_{i+1})}^n | z_{f(n;t_i)}^n, z_0^n) \tag{14}$$

472 Next, we obtain the explicit expression  $q(z_{f(n;t_{i+1})}^n | z_{f(n;t_i)}^n, z_0^n)$  through linear interpolation  
 473 between  $z_{f(n;t_i)}^n$  and  $z_0^n$ .

$$\begin{aligned}
 q(z_{f(n;t_{i+1})}^n | z_{f(n;t_i)}^n, z_0^n) &= \frac{q(z_{f(n;t_i)}^n | z_{f(n;t_{i+1})}^n, z_0^n) q(z_{f(n;t_{i+1})}^n | z_0^n)}{q(z_{f(n;t_i)}^n | z_0^n)} \\
 &= \frac{N\left(z_{f(n;t_i)}^n; \sqrt{\frac{f(n;t_i)}{f(n;t_{i+1})}} z_{f(n;t_{i+1})}^n; \left(1 - \frac{f(n;t_i)}{f(n;t_{i+1})}\right) l\right) N\left(z_{f(n;t_{i+1})}^n; \sqrt{f(n;t_{i+1})} z_0^n; \left(1 - \frac{f(n;t_{i+1})}{f(n;t_i)}\right) l\right)}{N\left(z_{f(n;t_i)}^n; \sqrt{\frac{f(n;t_i)}{f(n;t_i)}} z_0^n; \left(1 - \frac{f(n;t_i)}{f(n;t_i)}\right) l\right)} \\
 &\propto \exp \left\{ \frac{\left(z_{f(n;t_i)}^n - \sqrt{\frac{f(n;t_i)}{f(n;t_{i+1})}} z_{f(n;t_{i+1})}^n\right)^2}{2\left(1 - \frac{f(n;t_i)}{f(n;t_{i+1})}\right)} - \frac{\left(z_{f(n;t_{i+1})}^n - \sqrt{f(n;t_{i+1})} z_0^n\right)^2}{2\left(1 - \frac{f(n;t_{i+1})}{f(n;t_i)}\right)} \right. \\
 &\quad \left. + \frac{\left(z_{f(n;t_i)}^n - \sqrt{\frac{f(n;t_i)}{f(n;t_i)}} z_0^n\right)^2}{2\left(1 - \frac{f(n;t_i)}{f(n;t_i)}\right)} \right\} \\
 &= \exp \left\{ \frac{1}{2\left(1 - \frac{t_i}{t_{i+1}}\right)\left(1 - \frac{t_i}{t_{i+1}}\right)} \left[ \left(z_{f(n;t_{i+1})}^n\right)^2 - 2\left(\frac{\sqrt{\frac{f(n;t_i)}{f(n;t_{i+1})}}\left(1 - \frac{f(n;t_{i+1})}{f(n;t_i)}\right)}{1 - \frac{f(n;t_{i+1})}{f(n;t_i)}} z_{f(n;t_i)}^n\right. \right. \right. \\
 &\quad \left. \left. + \frac{\sqrt{f(n;t_{i+1})}\left(1 - \frac{f(n;t_i)}{f(n;t_{i+1})}\right)}{1 - \frac{f(n;t_{i+1})}{f(n;t_i)}} z_0^n\right) z_{f(n;t_{i+1})}^n \right] \right\} \\
 &\propto N\left(z_{f(n;t_{i+1})}^n; \frac{\sqrt{\frac{f(n;t_i)}{f(n;t_{i+1})}}\left(1 - \frac{f(n;t_{i+1})}{f(n;t_i)}\right)}{1 - \frac{f(n;t_{i+1})}{f(n;t_i)}} z_{f(n;t_i)}^n + \frac{\sqrt{f(n;t_{i+1})}\left(1 - \frac{f(n;t_i)}{f(n;t_{i+1})}\right)}{1 - \frac{f(n;t_{i+1})}{f(n;t_i)}} z_0^n; \right. \\
 &\quad \left. \frac{\left(1 - \frac{t_i}{t_{i+1}}\right)\left(1 - \frac{t_i}{t_{i+1}}\right)}{1 - \frac{t_i}{t_{i+1}}} l\right) \\
 &= N\left(z_{f(n;t_{i+1})}^n; z_{f(n;t_i)}^n + z_0^n; l\right)
 \end{aligned} \tag{15}$$



474 where we have the following notations for simplification.

$$r = \frac{\frac{f(n;t_i)}{f(n;t_{i+1})}(1 - f(n;t_{i+1}))}{1 - f(n;t_i)}; \quad p = \frac{\frac{f(n;t_i)}{f(n;t_{i+1})}(1 - \frac{f(n;t_i)}{f(n;t_{i+1})})}{1 - f(n;t_i)}; \quad = \frac{(1 - f(n;t_i))(1 - f(n;t_{i+1}))}{1 - f(n;t_i)}$$

475 Building upon equation (15), we substitute  $g(z_{f(n;t)}^n; f(n;t); x)$ , yielding the final formula  
 476 for  $p(z_{f(n;t_{i+1}}}^n | z_{f(n;t_i)}^n; x)$  as the following equation.

$$p(z_{f(n;t_{i+1}}}^n | z_{f(n;t_i)}^n; x) = N(z_{f(n;t_{i+1}}}^n; z_{f(n;t_i)}^n + g(z_{f(n;t)}^n; f(n;t); x); 1) \quad (16)$$

## 477 F Different Sampling Methods of BART

478 The sampling methods used by Lin et al. [2023] including Greedy Search, Beam Search Xiao et al.  
 479 [2022], Diverse Beam Search (diversity strength = 0.8) Vijayakumar et al. [2016], Typical Sample  
 480 ( $\alpha = 1.2$ ) Meister et al. [2022], Top-k Sampling ( $k = 50$ ) Fan et al. [2018] and Nucleus Sampling  
 481 (0.92) Holtzman et al. [2020].

482 Specifically, greedy search is to select the token with the highest probability at each step. Beam search  
 483 is to select the largest token from among the beams with higher probability at each step. Diverse beam  
 484 search is to divide the beams into multiple groups at each step and ensure the difference between  
 485 groups by calculating the diversity score between groups. Typical sampling selects samples through a  
 486 discrete random process. Top-k sampling is to randomly select one of the candidate tokens with the  
 487 highest probability at each step. Nucleus sampling is to randomly select one token at each step from  
 488 the candidate tokens whose probability density is greater than

## 489 G More Cases