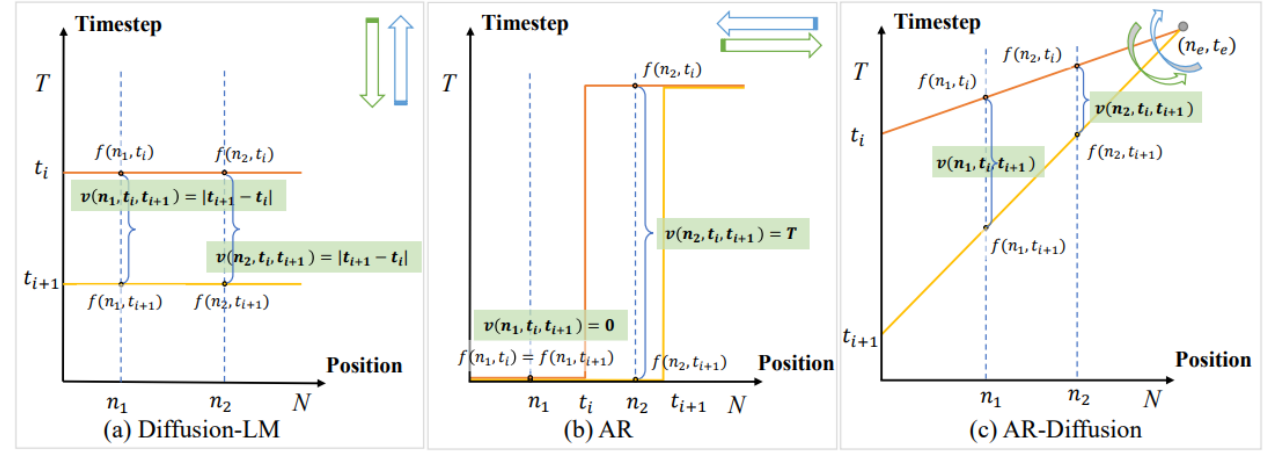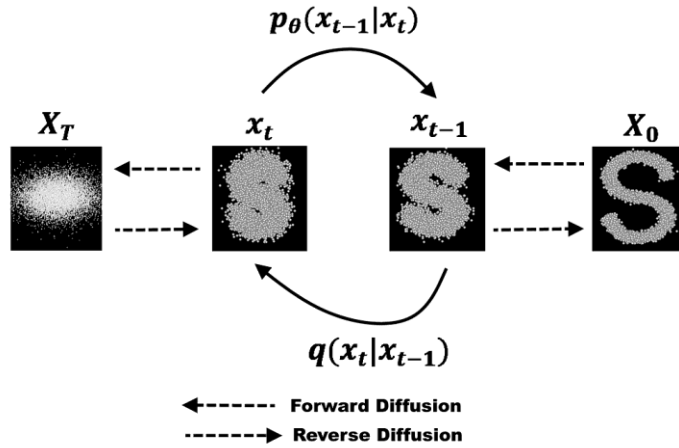# AR-Diffusion: Auto-Regressive Diffusion Model for Text Generation

Tong Wu, Zhihao Fan, Xiao Liu, Yeyun Gong, Yelong Shen, Jian Jiao, Hai-Tao Zheng, Juntao Li, Zhongyu Wei, Jian Guo, Nan Duan, Weizhu Chen

# Diffusion Models for Language Generation

# How AR-Diffusion works?

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

**Input**: Dataset $\{(x, y)\}$, maximum timestep number $T$ and maximum target length $N$.
**Output**: Optimized model parameters $\theta$.

1: Define an anchor point $(n_e, t_e)^5$.
2: **repeat**
3:   Sample $(x, y)$ from the dataset and embed $y$ into $z_0$.
4:   Sample a sentence-level timestep $t$ from the interval $[0, N + T]$, then the start point is determined by the following equation:
$$(n_s, t_s) = (\text{clip}(N - t, 0, N), \text{clip}(t - N, 0, T)) \qquad (6)$$

5:   Use the point-slope linear function to determine the token-level timestep $f(n, t)$ in position $n$:
$$f(n, t) = \text{clip}\left(\frac{t_e - t_s}{n_e - n_s}(n - n_s) + t_s, 0, T\right) \qquad (7)$$

6:   Sample $z_{f(n,t)}^n$ for each $n$ in different positions with Gaussian reparameterization.
7:   According to equation (3) and equation (9), employ gradient descent to optimize the objective:
$$\min_{\theta}\left[-\log p_\theta(y \mid z_0; x) + \sum_{n=1}^{N}\|g_\theta(z_{f(n,t)}^n, f(n, t); x) - z_0\|^2\right] \qquad (8)$$

8: **until** converged

---

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

**Input**: Source condition $x$, number of decoding steps $M$ and model parameters $\theta$.
**Output**: Predicted target embedding $\hat{y}$.

1: Define an anchor point $(n_e, t_e)$.
2: Uniformly select a decreasing sequence of timesteps $\{t_i\}_{i=0}^{M}$ ranging from $T + N$ to 0.
3: Sample $z_{t_0} \sim \mathcal{N}(0, I)$.
4: **for** $i = 0$ 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 $z_{t_{i+1}} = \left(z_{f(1,t_{i+1})}^1, z_{f(2,t_{i+1})}^2, \cdots, z_{f(N,t_{i+1})}^N\right)$ from $p_\theta(z_{t_{i+1}} \mid z_{t_i}; x)$ with the following formulas:
$$p_\theta(z_{t_{i+1}} \mid z_{t_i}; x) = \prod_{n=1}^{N} p_\theta\left(z_{f(n,t_{i+1})}^n \mid z_{f(n,t_i)}^n; x\right) \qquad (10)$$

$$p_\theta\left(z_{f(n,t_{i+1})}^n \mid z_{f(n,t_i)}^n; x\right) \sim \mathcal{N}\left(z_{f(n,t_{i+1})}^n; \lambda z_{f(n,t_i)}^n + \mu g_\theta(z_{f(n,t)}^n, f(n, t); x), \sigma I\right) \qquad (11)$$

8: **end for**
9: Map $z_{t_M}$ to the nearest embedding $\hat{y}$.

# Experimental Results

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] | NAR | 24.0 | 4.0 | 20.4 |
| CMLM [Ghazvininejad et al., 2019] | NAR | 23.8 | 3.6 | 20.2 |
| LevT [Gu et al., 2019] | NAR | 24.8 | 4.2 | 20.9 |
| InsT [Stern et al., 2019] | Semi-NAR | 17.7 | 5.2 | 16.1 |
| iNAT [Lee et al., 2018] | Semi-NAR | 27.0 | 6.9 | 22.4 |
| CMLM [Ghazvininejad et al., 2019] | Semi-NAR | 29.1 | 7.7 | 23.0 |
| LevT [Gu et al., 2019] | Semi-NAR | 25.3 | 7.4 | 21.5 |
| LSTM [Greff et al., 2017] | AR[10] | 25.1 | 6.9 | 19.9 |
| Transformer [Vaswani et al., 2017] | AR[10] | 30.5 | 10.4 | 24.2 |
| GENIE [Lin et al., 2023] ($k = 50$) | Diffusion | 29.3 | 8.3 | 21.9 |
| AR-Diffusion ($k = 50$) | Diffusion | 31.7 | 10.1 | 24.7 |
| AR-Diffusion ($k = 500$) | Diffusion | 32.2 | 10.6 | 25.2 |

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] | AR | 39.5 | 16.7 | 36.7 |
| GENIE [Lin et al., 2023] ($k = 50$) | Diffusion | 34.4 | 12.8 | 32.1 |
| AR-Diffusion ($k = 50$) | Diffusion | 39.6 | 16.3 | 37.1 |
| AR-Diffusion ($k = 500$) | Diffusion | 40.2 | 17.1 | 37.7 |

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 ($k = 1$) | Diffusion | 30.19 | 20 | 20 (20 × 1) |
| GENIE [Lin et al., 2023] ($k = 50$) | Diffusion | 30.08 | 20 | 1,000 (20 × 50) |
| AR-Diffusion ($k = 50$) | Diffusion | 34.95 | 20 | 1,000 (20 × 50) |
| AR-Diffusion ($k = 500$) | Diffusion | 35.62 | 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 | ROUGE-L | BLEU-3 | BLEU-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] | AR | 11.08 | 32.57 | 17.20 | 10.60 | 18.80 | 18.00 |
| MeanPooling-CopyNet [Lin et al., 2020] | AR | 11.36 | 34.63 | 14.80 | 8.90 | 19.20 | 20.20 |
| LevT [Gu et al., 2019] | NAR | 12.22 | 35.42 | 23.10 | 15.00 | 22.10 | 21.40 |
| ConstLeven [Susanto et al., 2020] | NAR | 13.47 | 35.19 | 21.30 | 12.30 | 25.00 | 23.20 |
| GENIE [Lin et al., 2023] ($k = 50$) | Diffusion | 12.89 | 35.21 | 22.00 | 13.30 | 24.30 | 23.00 |
| AR-Diffusion ($k = 50$) | Diffusion | 13.93 | 37.36 | 25.60 | 16.40 | 25.00 | 24.20 |

If you have any questions, please contact
*wu-t21@mails.tsinghua.edu.cn*