



## MOTIVATION

- Unlike Autoregressive models, Diffusion Models
  - ✓ Support **parallel generation**.
  - ✓ Have potential to improve **long-term planning, controllable generation, and sampling speed**.

However, until recently Diffusion Models  
 ✗ Exhibit a **performance gap** relative to AR models.

## NOTATION

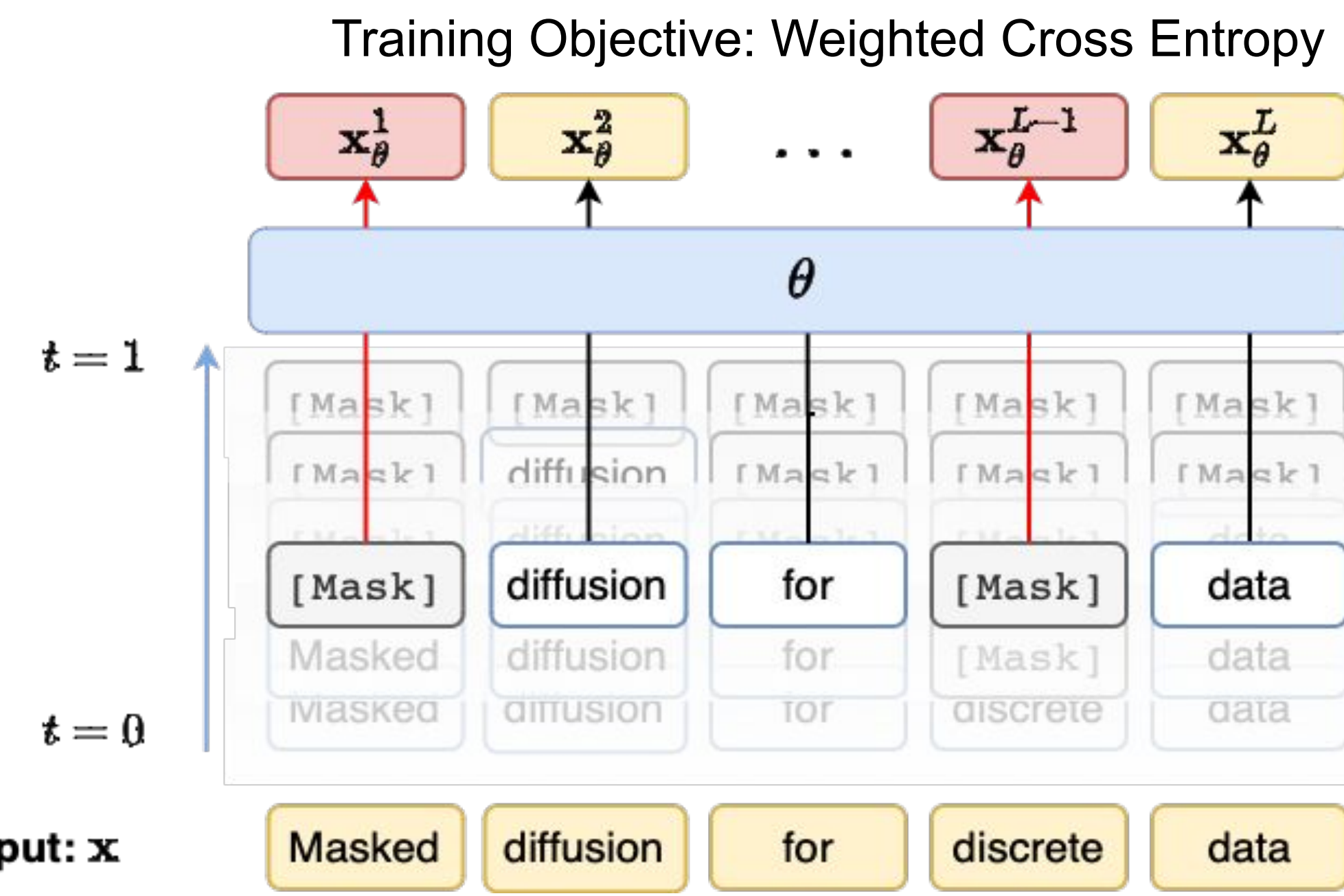
- $\mathcal{V} \in \{\mathbf{x} \in \{0,1\}^K : \sum_{i=1}^K \mathbf{x}_i = 1\}$
- $\Delta^K$ : Probability simplex
- $\ell$ : Index in a sequence
- $L$ : Sequence length
- $\mathbf{x} \in \mathcal{V}^L$ : input sequence of length  $L$
- $\mathbf{x}^\ell \in \mathcal{V}$ : Input token at index  $\ell$
- $\bar{\mathbf{m}} \in \mathcal{V}$ : mask token
- $\mathbf{m}$ : Sequence of  $L$  mask tokens
- $\alpha_t \in [0, 1]$ : Signal level
- $t \in [0, 1]$ : Diffusion time step
- $T$ : Total number of diffusion steps
- $\mathbf{x}_\theta(\mathbf{z}_t) : \mathcal{V}^L \rightarrow (\Delta^K)^L$ : Denoising model
- $\mathbf{z}_t \in \mathcal{V}^L$ : Masked input
- $\langle \mathbf{a}, \mathbf{b} \rangle$ : Dot product between 2 vectors  $\mathbf{a}$  and  $\mathbf{b}$

### $\mathbf{x}$

Many years later, as he faced the firing squad, Colonel Aureliano Buendía was to remember that distant afternoon when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the bank of a river of clear water that ran along a bed of polished stones, which were white and enormous, like prehistoric eggs.

### $\mathbf{z}_t$

Many years later, -- he faced --- firing squad, ----- Aureliano Buendía was to remember that ----- afternoon when his father took him to discover ice. At that time Macondo was a village of twenty adobe houses, built on the ---- of a river of clear water that ran ----- a bed of polished stones, which were ----- --- enormous, like prehistoric eggs.

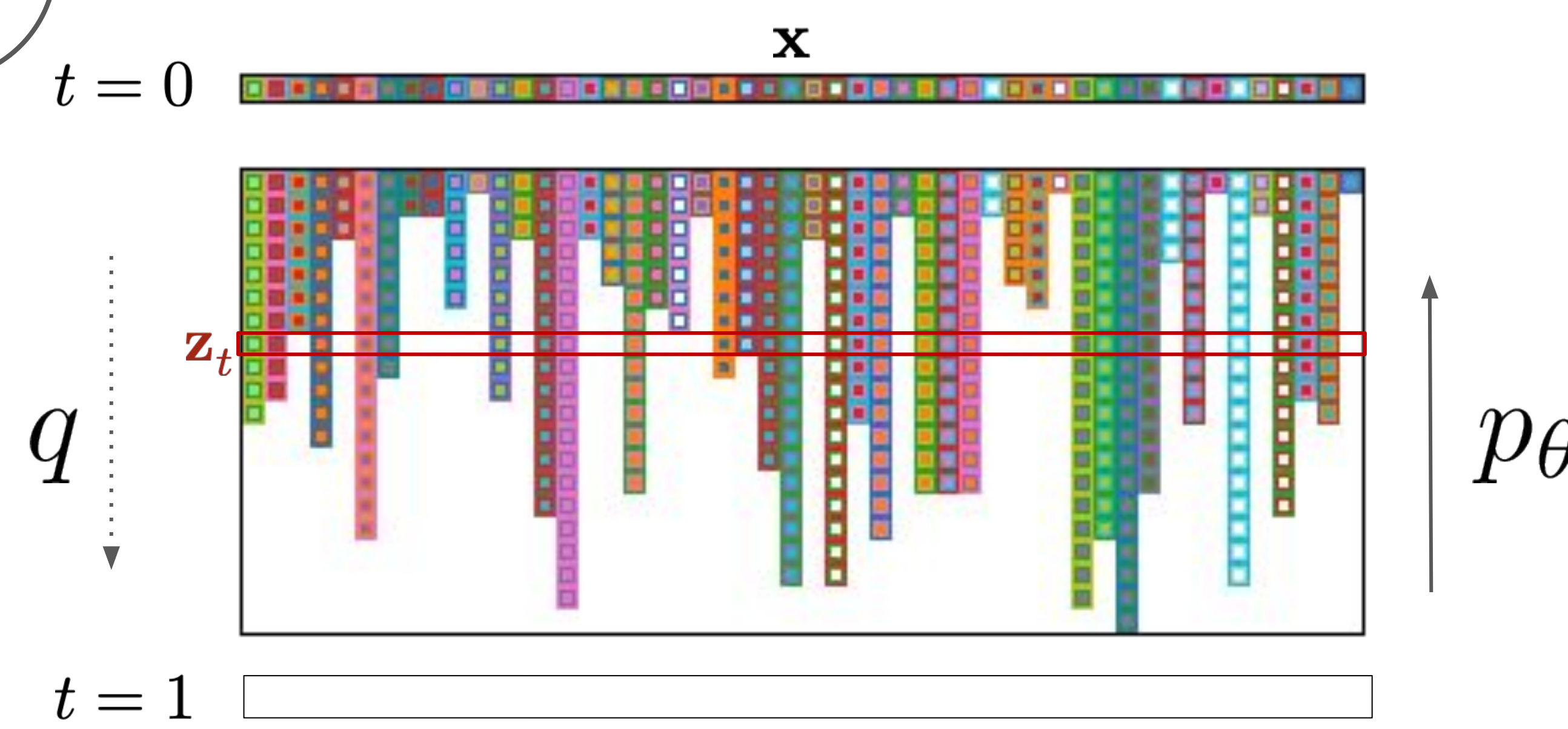


## Simplified Masked Diffusion LM

- **Simple** interpolating masking noise processes
- Masking rate is **random**, Objective is a **variational lower bound**
- Admits fast **ancestral sampling**
- Novel substitution based parameterization
- Diffusion Objective is a **simple average** of MLM losses
- **Improved implementation** relative to previous masked diffusion

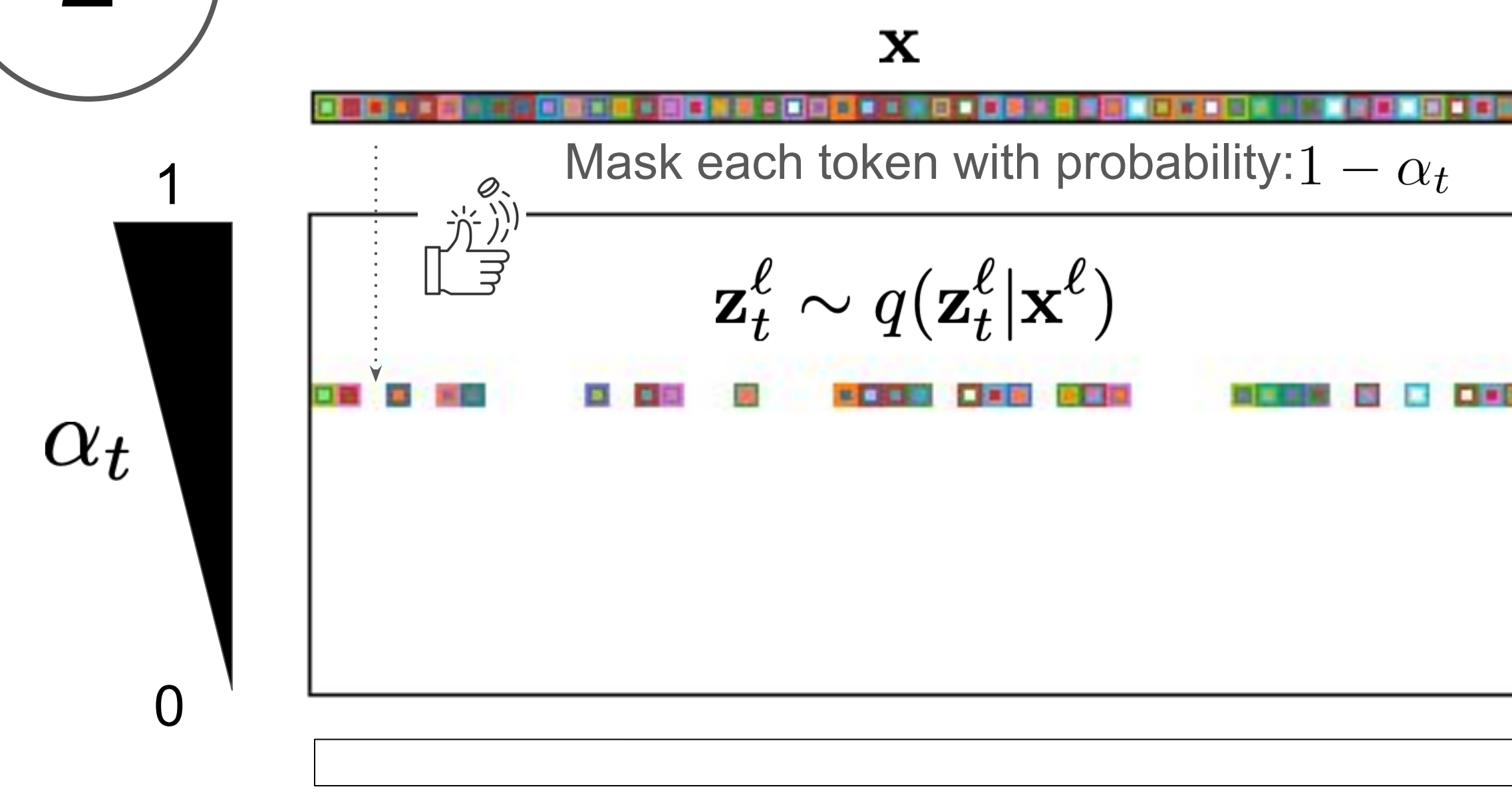
## MASKED DIFFUSION LANGUAGE MODEL (MDLM)

### 1 Masked Diffusion



We consider **simplified diffusion processes** that interpolate between the clean data,  $\mathbf{x}$ , and the prior  $\mathbf{m}$

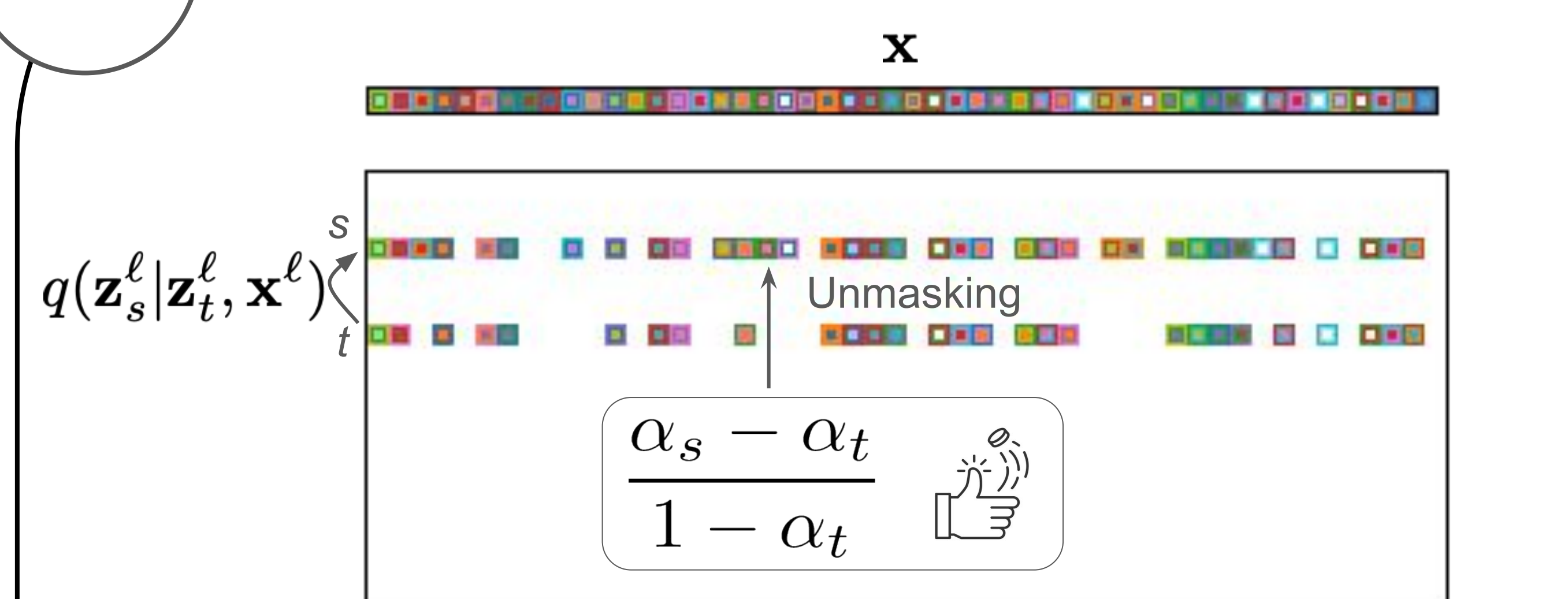
### 2 Forward Process



A diffusion process defines latents  $\mathbf{z}_t^\ell$  for each token  $\mathbf{x}^\ell$  for  $t \in \{0, 1/T, \dots, 1\}$  using markov forward process inducing the marginals:

$$q(\mathbf{z}_t^\ell | \mathbf{x}^\ell) = \text{Cat}(\mathbf{z}_t^\ell; \alpha_t \mathbf{x}^\ell + (1 - \alpha_t) \bar{\mathbf{m}})$$

### 3 Reverse Process



- The true reverse posterior for a token for a timestep  $s < t$  is given as:

$$q(\mathbf{z}_s^\ell | \mathbf{z}_t^\ell, \mathbf{x}^\ell) = \begin{cases} \text{Cat}(\mathbf{z}_s^\ell; \mathbf{z}_t^\ell) & \mathbf{z}_t^\ell \neq \bar{\mathbf{m}} \\ \text{Cat}(\mathbf{z}_s^\ell; \frac{(1 - \alpha_s)\bar{\mathbf{m}} + (\alpha_s - \alpha_t)\mathbf{x}^\ell}{1 - \alpha_t}) & \mathbf{z}_t^\ell = \bar{\mathbf{m}} \end{cases}$$

- The approximate posterior is given as:

$$p_\theta(\mathbf{z}_s^\ell | \mathbf{z}_t) = q(\mathbf{z}_s^\ell | \mathbf{z}_t^\ell, \mathbf{x}^\ell = \mathbf{x}_\theta^\ell(\mathbf{z}_t))$$

where  $\mathbf{x}_\theta(\mathbf{z}_t)$  is a BERT style denoising model.

### 4 Training Objective

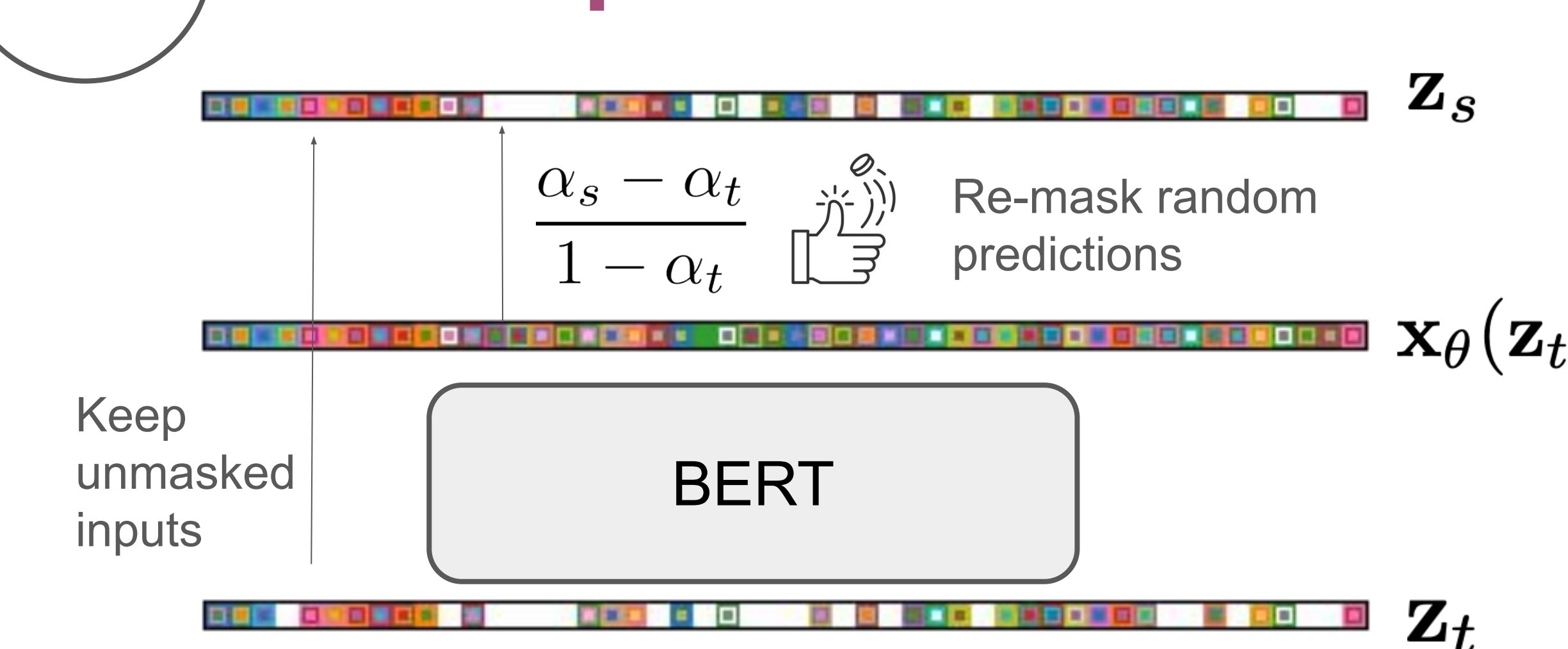
Denoising model  $\mathbf{x}_\theta(\mathbf{z}_t)$  uses **SUBSTITUTION** based parameterization:

- **Zero Masking Probabilities:**  $\langle \mathbf{x}_\theta^\ell(\mathbf{z}_t), \bar{\mathbf{m}} \rangle = 0$
- **Carry Over Unmasking:**  $\mathbf{x}_\theta^\ell(\mathbf{z}_t) = \mathbf{z}_t^\ell$  if  $\mathbf{z}_t^\ell \neq \bar{\mathbf{m}}$

$T \rightarrow \infty$  + **SUBS** parameterization yields the following **simplified Negative Evidence Lower Bound (N-ELBO)**:

$$\mathbb{E}_q \frac{\alpha'_t}{1 - \alpha_t} \sum_{\ell} \log \langle \mathbf{x}_\theta^\ell(\mathbf{z}_t), \mathbf{x}^\ell \rangle$$

### 5 Sample Generation



Given a noisy sample  $\mathbf{z}_t$ , we construct the clean input  $\mathbf{x}_\theta(\mathbf{z}_t)$ . We then re-mask each token independently with a specified probability, ensuring that tokens unmasked in  $\mathbf{z}_t$  are not re-masked

## RESULTS

### MDLM Perplexity Approaches AR

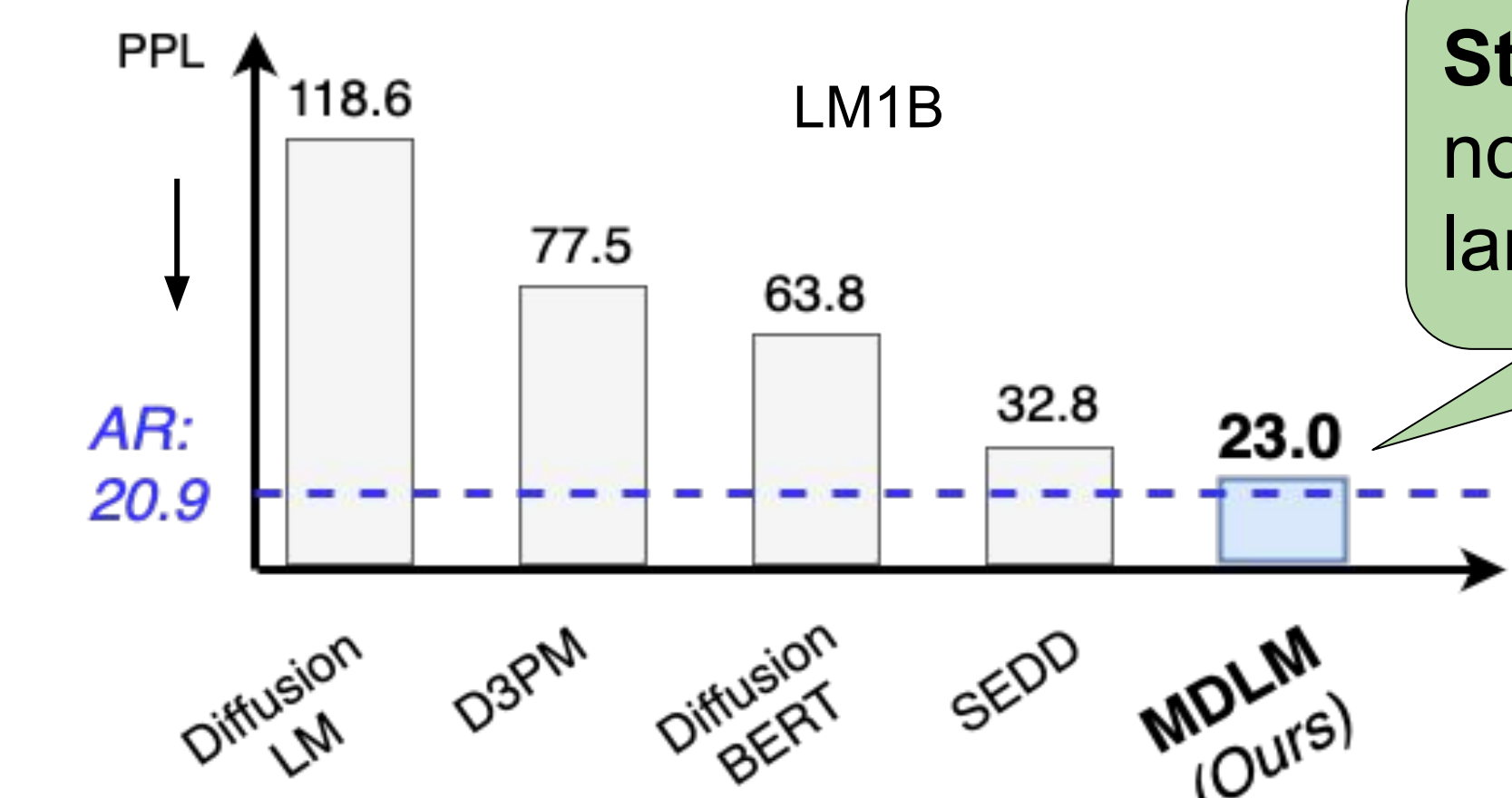


Table 3: Zero-shot validation perplexities ( $\downarrow$ ) of models trained for 524B tokens on OWT. All perplexities for diffusion models are upper bounds.

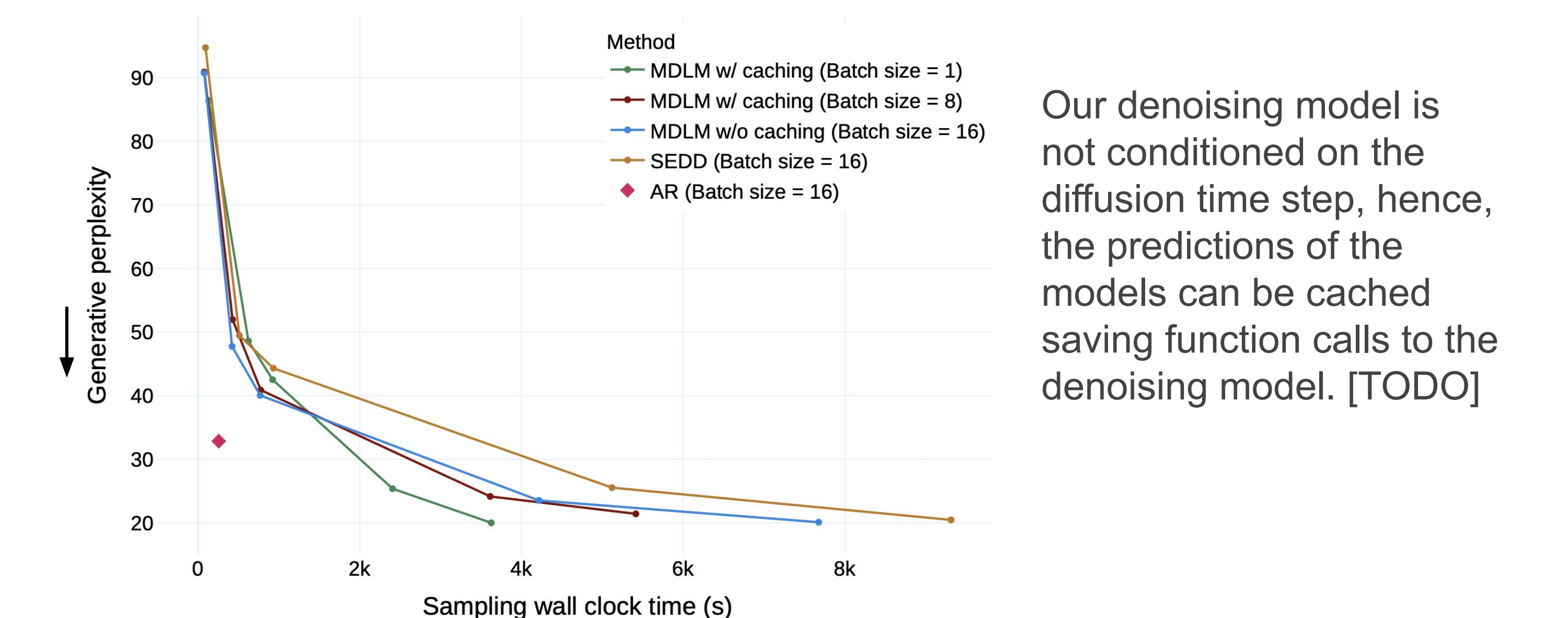
	PTB	Wikitext	LM1B	Lambada	AG News	Pubmed	Arxiv
AR (Retrained)	<b>82.05</b>	<b>25.75</b>	<b>51.25</b>	51.28	<b>52.09</b>	49.01	41.73
SEDD (Retrained)	100.09	34.28	68.20	49.86	62.09	44.53	38.48
MDLM (Ours)	95.26	32.83	67.01	<b>47.52</b>	61.15	<b>41.89</b>	<b>37.37</b>

### Semi-AR Text Generation: Generating Sequences of Arbitrary Length

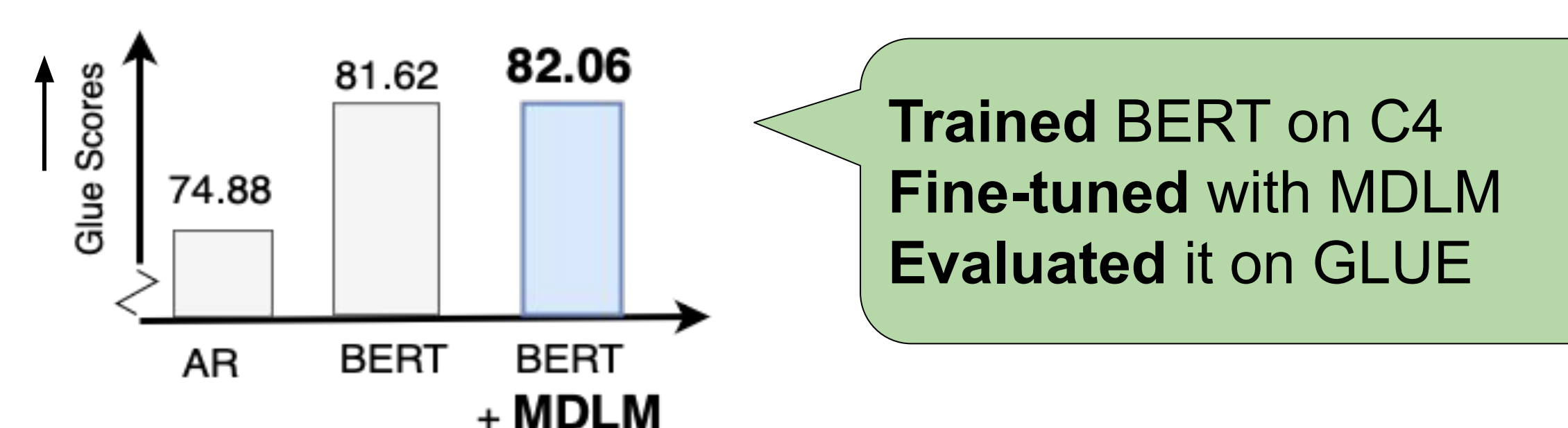
Table 5: Semi-AR generative perplexity (Gen. PPL;  $\downarrow$ ) for sequences of 2048 tokens.

	Gen. PPL ( $\downarrow$ )	Sec/Seq ( $\downarrow$ )
SSD-LM	35.43	2473.9
MDLM (Ours)	<b>27.18</b>	<b>89.3</b>

### Faster Sampler



### MDLM Preserves Representation Learning Capabilities



## REFERENCES

- [1] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. Advances in Neural Information Processing Systems, 34:17981–17993, 2021.
- [2] Zhengfu He, Tianxiang Sun, Kuanning Wang, Xuanjing Huang, and Xipeng Qiu. Diffusion-bert: Improving generative masked language models with diffusion models. arXiv preprint arXiv:2211.15029, 2022.
- [3] Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. Advances in Neural Information Processing Systems, 35:4328–4343, 2022.
- [4] Aaron Lou, Chenlin Meng, and Stefano Ermon. Discrete diffusion language modeling by estimating the ratios of the data distribution. arXiv preprint arXiv:2310.16834, 2023.

Video Tutorial, Code, and more ...

