

# Diffusion Language Models

Zhihan Huang, Kevin Jiang

Wharton Statistics & Data Science  
University of Pennsylvania

March 6, 2025

# Outline

## 1. Diffusion models and LMs

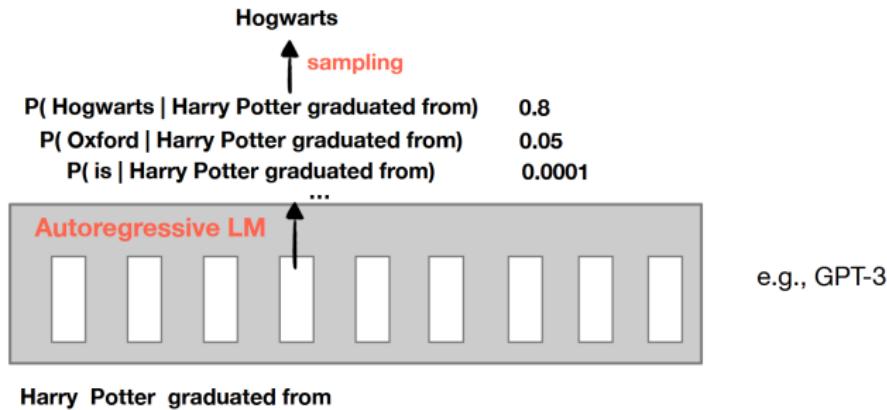
## 2. General Frameworks

- 2.1 Diffusion LM
- 2.2 Masking Diffusion Language Model

## 3. Survey and extensions

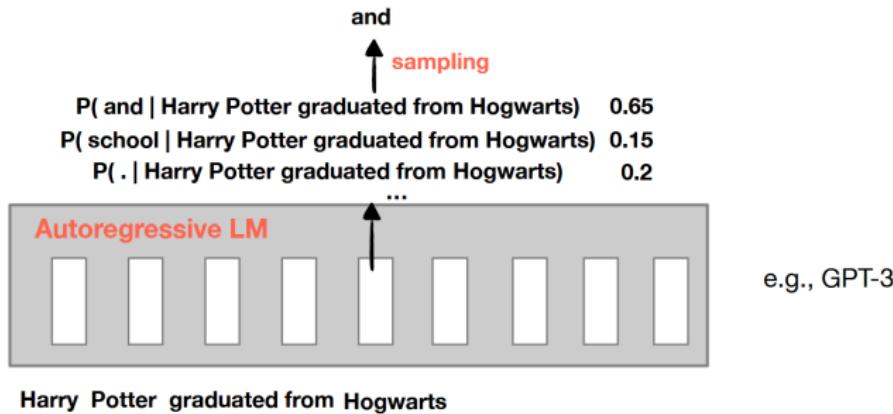
# Autoregressive LMs

Generate text from the underlying distribution:



# Autoregressive LMs

Generate text from the underlying distribution:



Next-token prediction: Generating text from left-to-right, one at a time.

- Parametrize the probability of a sentence via chain rule:
- For each token, compute the next-token distribution.

# Motivation

Some motivation for going beyond next-token prediction:

1. What if we want to generate right-to-left? Or given left and right context, fill in the middle text (inpainting)?

# Motivation

Some motivation for going beyond next-token prediction:

1. What if we want to generate right-to-left? Or given left and right context, fill in the middle text (inpainting)?
2. For Autoregressive models, time complexity is always  $O(L)$  where  $L$  is the length of generated text.
  - Can we do parallel sampling from a language model? Generate the whole sentence at once?

# Motivation

Some motivation for going beyond next-token prediction:

1. What if we want to generate right-to-left? Or given left and right context, fill in the middle text (inpainting)?
2. For Autoregressive models, time complexity is always  $O(L)$  where  $L$  is the length of generated text.
  - Can we do parallel sampling from a language model? Generate the whole sentence at once?

Can we model the text distribution by one joint distribution (instead of a sequential model) and sample from it: **Diffusion Language Models**  
**(Diffusion LMs)**

$$x^{1:L} \sim P_\theta(x)$$

# Diffusion LMs

How to do it?

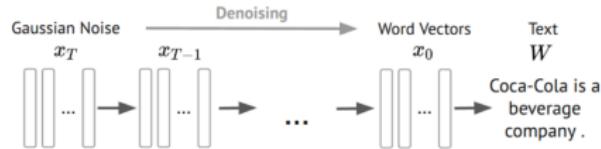
1. Learn an implicit probabilistic model  $P_\theta$  from the data: a very high dimensional distribution in  $\mathbb{R}^{L \cdot d}$  ( $d$  the embedding dimension,  $L$  the length of token)
2. Find an accurate and efficient way to sample from  $P_\theta$ : discrete, low-dimensional manifold, non-smooth?

# Diffusion LMs

How to do it?

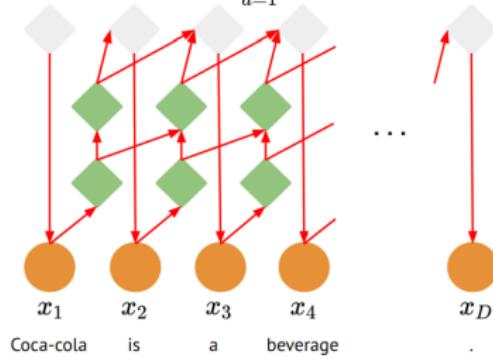
1. Learn an implicit probabilistic model  $P_\theta$  from the data: a very high dimensional distribution in  $\mathbb{R}^{L \cdot d}$  ( $d$  the embedding dimension,  $L$  the length of token)
  2. Find an accurate and efficient way to sample from  $P_\theta$ : discrete, low-dimensional manifold, non-smooth?
- 
- ▶ Use **Diffusion Model** to model and sample from  $P_\theta$
  - ▶ Diffusion LM is more like an experimental model exploring different architectures for tasks — no standard answers
  - ▶ Very different from the autoregressive model — Can diffusion LMs be competitive with autoregressive (AR) models in some tasks?

# Diffusion LMs



(a) Diffusion-based language model

$$p(x) = \prod_{d=1}^D p(x_d | x_{<d})$$



(b) AR language model

## What is a diffusion model?

Diffusion Models are successful in image generation



## What is a diffusion model?

Diffusion Models are successful in image generation

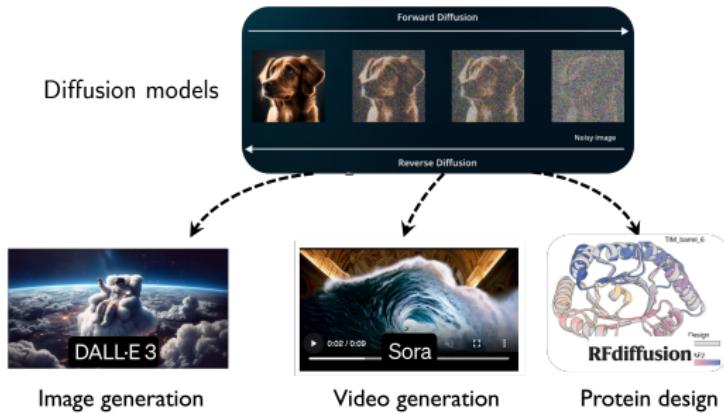


Efficient for conditional generation, e.g., generation with prompts



# What is a diffusion model?

Apply in various tasks across image generation, video generation and some scientific fields



# Diffusion Model

A stochastic/deterministic sampler to sample from a high-dimensional, non-smooth, implicit distribution:

# Diffusion Model

A stochastic/deterministic sampler to sample from a high-dimensional, non-smooth, implicit distribution:

- ▶ Given training data  $X^{\text{train},i} \sim p_{\text{data}}$  ( $1 \leq i \leq N$ ) in Euclidean space
- ▶ Generate rate **new** samples  $X \sim p_{\text{data}}$

# Diffusion Model

A stochastic/deterministic sampler to sample from a high-dimensional, non-smooth, implicit distribution:

- ▶ Given training data  $X^{\text{train},i} \sim p_{\text{data}}$  ( $1 \leq i \leq N$ ) in Euclidean space
- ▶ Generate rate **new** samples  $X \sim p_{\text{data}}$
- ▶ Usually  $p_{\text{data}}$  is a general distribution (density hard to be expressed or may not even exist)

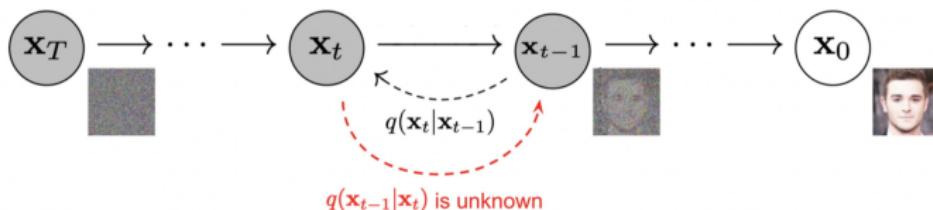
## Diffusion Model

A stochastic/deterministic sampler to sample from a high-dimensional, non-smooth, implicit distribution:

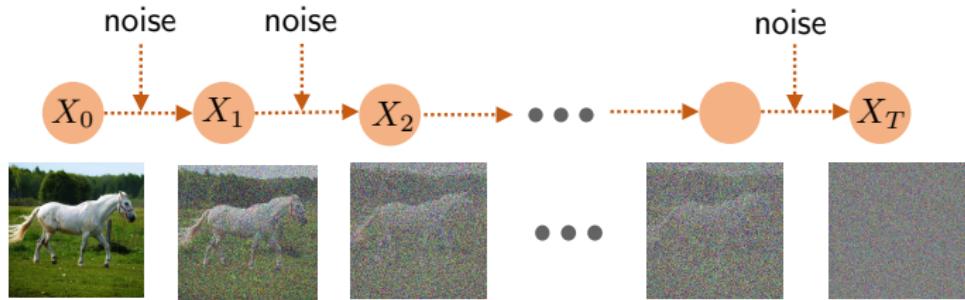
- ▶ Given training data  $X^{\text{train},i} \sim p_{\text{data}}$  ( $1 \leq i \leq N$ ) in Euclidean space
  - ▶ Generate rate **new** samples  $X \sim p_{\text{data}}$
  - ▶ Usually  $p_{\text{data}}$  is a general distribution (density hard to be expressed or may not even exist)

Overcome the challenges?

- Learn the information from  $p_{\text{data}}$  through a diffusion process
  - Sample by reversing this process and introducing new randomness

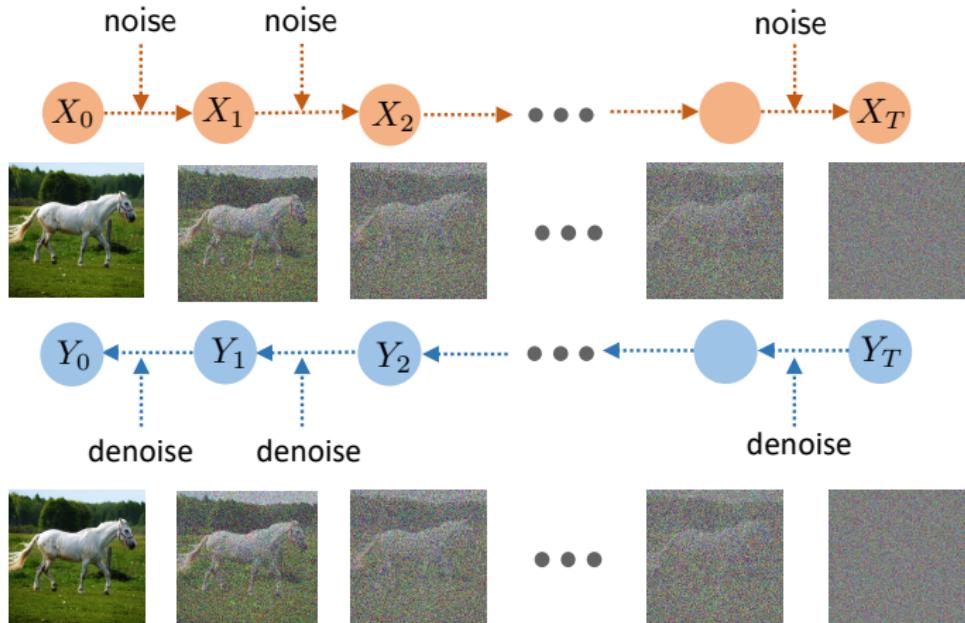


## A Example with Gaussian noise



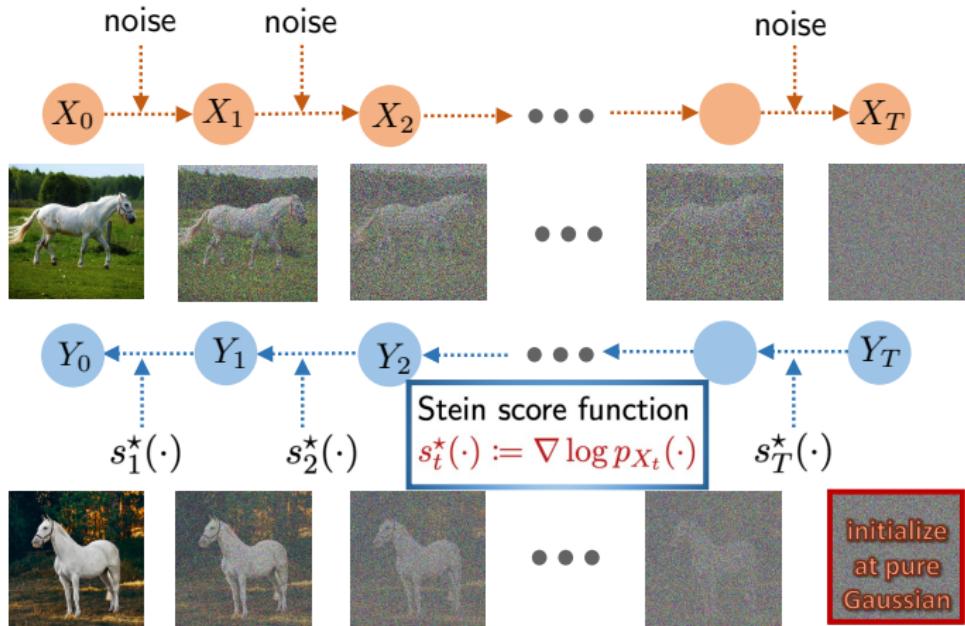
- ▶  $X_T$  is very close to pure Gaussian noise when  $T$  is large

## A Example with Gaussian noise



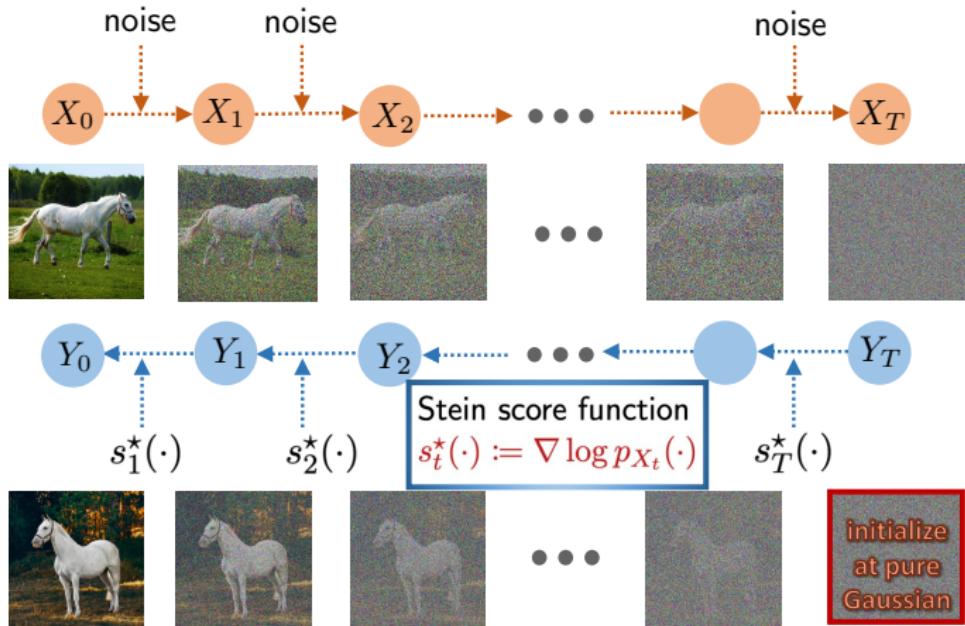
- ▶  $X_T$  is very close to pure Gaussian noise when  $T$  is large
  - ▶ Pretend  $Y_T = X_T = \mathcal{N}(0, I_d)$  and reverse it

## A Example with Gaussian noise



- ▶  $X_T$  is very close to pure Gaussian noise when  $T$  is large
- ▶ Pretend  $Y_T = X_T = \mathcal{N}(0, I_d)$  and reverse it

## A Example with Gaussian noise



- ▶  $X_T$  is very close to pure Gaussian noise when  $T$  is large
- ▶ Pretend  $Y_T = X_T = \mathcal{N}(0, I_d)$  and reverse it

# Diffusion Model

- ▶ Stochastic Differential Equation perspective:

*forward process* :  $X_0 \sim p_{\text{data}},$

$$dX_t = -X_t dt + \sqrt{2} dB_t$$

$$\xrightarrow{\text{discretize}} X_t = \sqrt{1 - \sigma_t^2} X_{t-1} + \sigma_t \mathcal{N}(0, I_d), \quad 1 \leq t \leq T$$

Based on this forward process, we aim to construct a (Markovian) backward process such that

**Goal** :  $X_t \stackrel{d}{\approx} Y_t \quad 1 \leq t \leq T$

*backward process* :  $Y_T \sim \mathcal{N}(0, I_d),$

$$dY_t = (Y_t + 2s_t^*(Y_t))dt + \sqrt{2} dB_t$$

$$\xrightarrow{\text{discretize}} Y_{t-1} = c_t^1 Y_t + c_t^2 s_t^*(Y_t) + c_t^3 \mathcal{N}(0, I_d)$$

$\Rightarrow$  Output  $Y_1$

# Diffusion Model

- ▶ **Variational Inference** perspective: In each backward step, find the best Gaussian approximation of  $X_{t-1}|X_t$  by minimizing the variational lower bound:

$$\begin{aligned} & \min_{\theta} \mathbb{E}[\text{KL}(q(X_{t-1}|X_t) || p_{\theta}(X_{t-1}|X_t))] \\ & \approx \min_{\theta} \mathbb{E}[\text{KL}(q(X_{t-1}|X_t, X_0) || p_{\theta}(X_{t-1}|X_t, X_0))] \end{aligned}$$

1. Notice  $X_{t-1}|X_t \Rightarrow$  intractable  $\Rightarrow X_{t-1}|X_t, X_0$  is Gaussian  $\Rightarrow$  tractable with estimated  $X_0$
2. In practice, minimize the sum of  $KL$  error from  $1 \rightarrow T$ , parametrized by unified  $\theta$

Implicitly learn the score function  $s_{X_t}^*(\cdot) \Leftrightarrow$  the SDE form

# Discrete Diffusion Model

- ▶ When the data is supported on a discrete space, Gaussian noise is not the ideal choice
- ▶ Design a discrete space Markov chain as the forward process, whose the limit is some uniform distribution (easy to sample from)

$$q(X_t) = Q_t q(X_0) \quad \text{for } 1 \leq t \leq T$$

Similar to the continuous case, learn the discrete score function  
 $s_t^*(x) = \mathbf{p}_{X_t}/p_{X_t}(x)$

- ▶ The backward process can be designed as

$$p_\theta(Y_{t-1}) = \tilde{Q}_t p_\theta(Y_t)$$

where  $\tilde{Q}_t$  is a function of  $Q_t$  and  $s_t^*(x)$

# Two Example Frameworks

---

## Diffusion-LM Improves Controllable Text Generation

---

**Xiang Lisa Li**  
Stanford University  
[xlisali@stanford.edu](mailto:xlisali@stanford.edu)

**John Thickstun**  
Stanford University  
[jthickst@stanford.edu](mailto:jthickst@stanford.edu)

**Ishaan Gulrajani**  
Stanford University  
[igul@stanford.edu](mailto:igul@stanford.edu)

**Percy Liang**  
Stanford University  
[pliang@cs.stanford.edu](mailto:pliang@cs.stanford.edu)

**Tatsunori B. Hashimoto**  
Stanford University  
[tashim@stanford.edu](mailto:tashim@stanford.edu)

---

## Simple and Effective Masked Diffusion Language Models

---

**Subham Sekhar Sahoo**  
Cornell Tech, NYC, USA.  
[ssahoo@cs.cornell.edu](mailto:ssahoo@cs.cornell.edu)

**Marianne Arriola**  
Cornell Tech, NYC, USA.  
[ma2238@cornell.edu](mailto:ma2238@cornell.edu)

**Yair Schiff**  
Cornell Tech, NYC, USA.  
[yzs2@cornell.edu](mailto:yzs2@cornell.edu)

**Aaron Gokaslan**  
Cornell Tech, NYC, USA.  
[akg67@cs.cornell.edu](mailto:akg67@cs.cornell.edu)

**Edgar Marroquin**  
Cornell Tech, NYC, USA.  
[emm392@cornell.edu](mailto:emm392@cornell.edu)

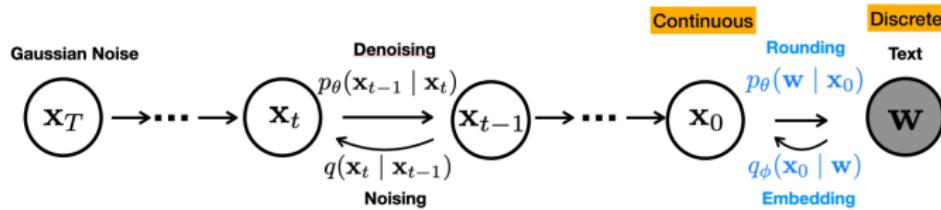
**Justin T Chiu**  
Cornell Tech, NYC, USA.  
[jtc257@cornell.edu](mailto:jtc257@cornell.edu)

**Alexander Rush**  
Cornell Tech, NYC, USA.  
[ar459@cornell.edu](mailto:ar459@cornell.edu)

**Volodymyr Kuleshov**  
Cornell Tech, NYC, USA.  
[kuleshov@cornell.edu](mailto:kuleshov@cornell.edu)

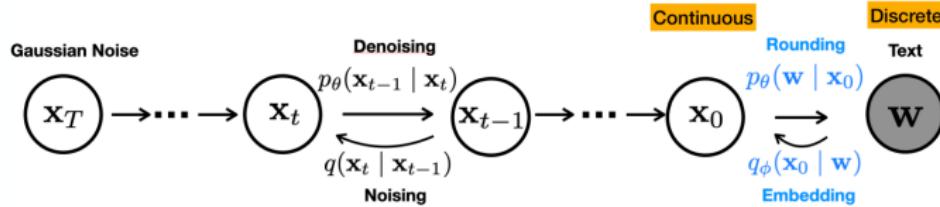
## Diffusion LM

Encode discrete texts into continuous space and add Gaussian noise



- ▶ Requires embedding and rounding steps
  - ▶ Lifting to continuous space allows one to compute gradients
  - ▶ Ultimately helps with controllable text generation

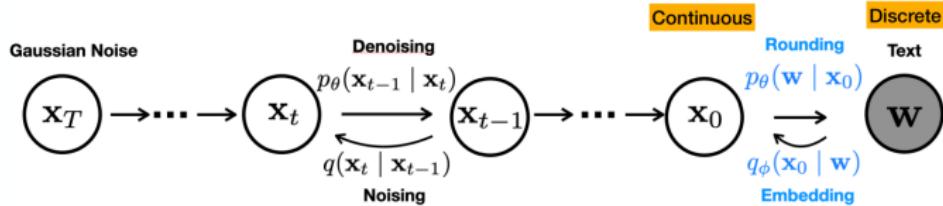
# Diffusion LM



- ▶ Embedding step? End-to-end learning (e2e), parameterize

$$q_\phi(x_0 | w) = \mathcal{N}(\text{Emb}_\phi(w), \sigma^2 \cdot I)$$

# Diffusion LM



- ▶ Embedding step? End-to-end learning (e2e), parameterize

$$q_\phi(x_0 | w) = \mathcal{N}(\text{Emb}_\phi(w), \sigma^2 \cdot I)$$

- ▶ Modified objective

$$\mathcal{L}_{\text{e2e}}(w) = \mathbb{E}_{x_0 \sim q_\phi} [\mathcal{L}_{\text{simple}}(x_0) - \log p_\theta(w | x_0)]$$

where  $\mathcal{L}_{\text{simple}}(x_0)$  is the standard objective for continuous time diffusion models

$$\mathcal{L}_{\text{simple}}(x_0) = \sum_{t=1}^T \mathbb{E}_{q(x_t | x_0)} [\|\hat{\mu}_\theta(x_t, t) - \mu_{t-1}(x_t, x_0)\|^2].$$

# Diffusion LM

How to specify rounding step?

- ▶ Ideally, we would like to take most likely word from softmax output

$$\hat{w}_i = \arg \max_{w_i} p_\theta(w_i | x_{0:i})$$

# Diffusion LM

How to specify rounding step?

- ▶ Ideally, we would like to take most likely word from softmax output

$$\hat{w}_i = \arg \max_{w_i} p_\theta(w_i | x_{0:i})$$

- ▶ Issue:  $x_0$  does not ‘commit’ to a single word embedding

# Diffusion LM

How to specify rounding step?

- ▶ Ideally, we would like to take most likely word from softmax output

$$\hat{w}_i = \arg \max_{w_i} p_\theta(w_i | x_{0:i})$$

- ▶ Issue:  $x_0$  does not ‘commit’ to a single word embedding
- ▶ Solution, modify the standard objective further

$$\mathcal{L}_{\text{simple}}(x_0) \leftarrow \mathcal{L}_{\text{round}}(x_0) = \sum_{t=1}^T \mathbb{E}_{q(x_t|x_0)} [\|\hat{f}_\theta(x_t, t) - x_0\|^2].$$

# Diffusion LM

How to specify rounding step?

- ▶ Ideally, we would like to take most likely word from softmax output

$$\hat{w}_i = \arg \max_{w_i} p_\theta(w_i | x_{0,i})$$

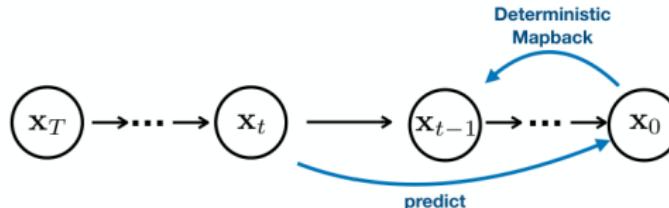
- ▶ Issue:  $x_0$  does not ‘commit’ to a single word embedding
- ▶ Solution, modify the standard objective further

$$\mathcal{L}_{\text{simple}}(x_0) \leftarrow \mathcal{L}_{\text{round}}(x_0) = \sum_{t=1}^T \mathbb{E}_{q(x_t|x_0)} [\|\hat{f}_\theta(x_t, t) - x_0\|^2].$$

- ▶ At decoding, use *clamping*

$$x_{t-1} = \sqrt{\bar{\alpha}} \cdot \text{Clamp}(f_\theta(x_t, t)) + \sqrt{1 - \bar{\alpha}} \cdot \epsilon.$$

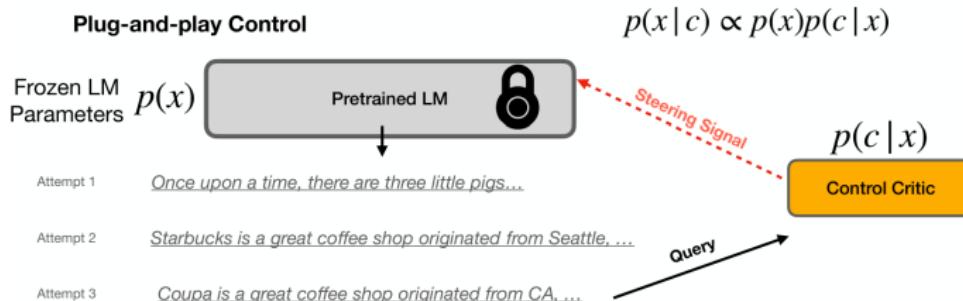
where  $\epsilon \sim \mathcal{N}(0, I)$ .



# Diffusion LM

Controllable text generation via plug-and-play

Goal: generate positive reviews about a coffee shop called Coupa?



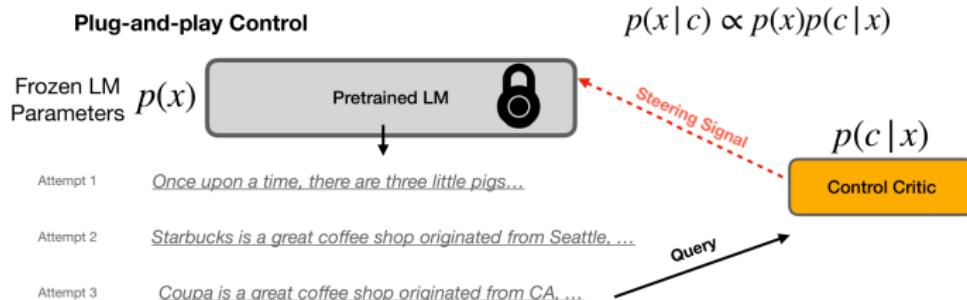
- ▶ Adds control in a post-hoc fashion

$$e_i^h \leftarrow e_i^h + \gamma \nabla_{e_i^h} \log p(c | e_i^h)$$

## Diffusion LM

# Controllable text generation via plug-and-play

## Goal: generate positive reviews about a coffee shop called Coupa?

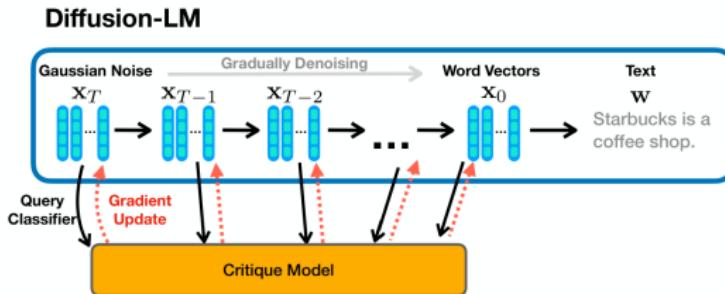


- ▶ Adds control in a post-hoc fashion

$$e_i^h \leftarrow e_i^h + \gamma \nabla_{e_i^h} \log p(c|e_i^h)$$

- ▶ Can we integrate control directly into the sampling process?

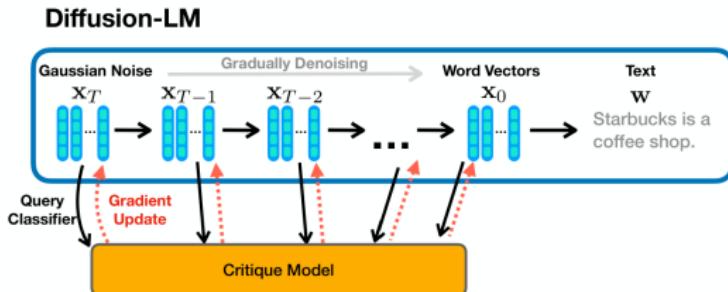
# Diffusion LM



- ▶ Goal: learn to sample from

$$p(x_{0:T}|c) = \prod_{t=1}^T p(x_{t-1}|x_t, c) \propto \prod_{t=1}^T p(x_{t-1}|x_t) \cdot p(c|x_{t-1}, x_t)$$

# Diffusion LM



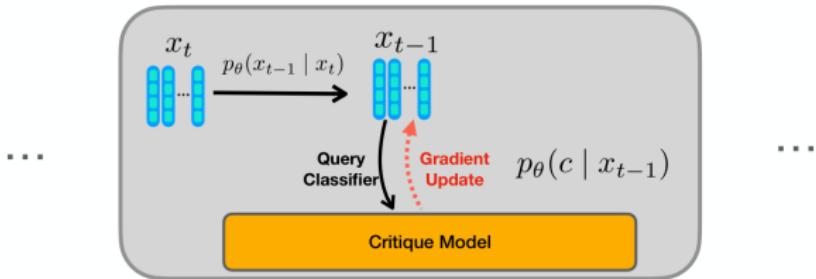
- ▶ Goal: learn to sample from

$$p(x_{0:T}|c) = \prod_{t=1}^T p(x_{t-1}|x_t, c) \propto \prod_{t=1}^T p(x_{t-1}|x_t) \cdot p(c|x_{t-1}, x_t)$$

- ▶ Motivates the following gradient update

$$\nabla_{x_{t-1}} \log p(x_{t-1}|x_t, c) = \nabla_{x_{t-1}} \log p(x_{t-1}|x_t) + \nabla_{x_{t-1}} \log p(c|x_{t-1})$$

# Diffusion LM



Update  $x_{t-1}$  in the following gradient direction:

$$\nabla_{x_{t-1}} \log p(\mathbf{x}_{t-1} | \mathbf{x}_t) + \nabla_{x_{t-1}} \log p(\mathbf{c} | \mathbf{x}_{t-1}),$$

Diffusion-LM  
transition                      Classifier  
                                  score

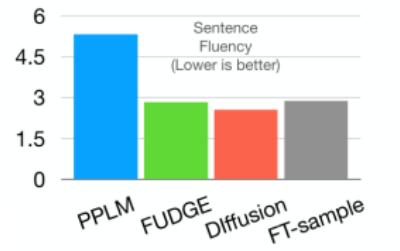
- ▶ use *fluency regularization*

$$\lambda \nabla_{x_{t-1}} \log p(x_{t-1} | x_t) + \nabla_{x_{t-1}} \log p(c | x_{t-1})$$

balances maximizing and sampling akin to nucleus sampling and temperature scaling

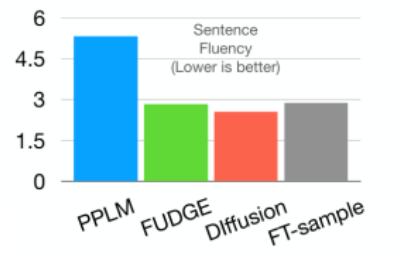
# Diffusion LM

- ▶ Semantic Content task: generate a sensible sentence with a given field and value
- ▶ Food = Japanese → “Browns Cambridge is good for Japanese food and also children friendly.”



# Diffusion LM

- ▶ Semantic Content task: generate a sensible sentence with a given field and value
- ▶ Food = Japanese → “Browns Cambridge is good for Japanese food and also children friendly.”



80M parameters, sequence length  $n = 64$ , diffusion steps  $T = 2000$ , square-root noise schedule; teacher LM is *GPT-2 model*

# Diffusion LM

input (Semantic Content)	food : Japanese
output text	Browns Cambridge is good for Japanese food and also children friendly near The Sorrento .
input (Parts-of-speech)	PROPN AUX DET ADJ NOUN NOUN VERB ADP DET NOUN ADP DET NOUN PUNCT
output text	Zizzi is a local coffee shop located on the outskirts of the city .
input (Syntax Tree)	(TOP (S (NP (#) (#) (#)) (VP (#) (NP (NP (#) (#))))))
output text	The Twenty Two has great food
input (Syntax Spans)	(7, 10, VP)
output text	Wildwood pub serves multicultural dishes and is ranked 3 stars
input (Length)	14
output text	Browns Cambridge offers Japanese food located near The Sorrento in the city centre .
input (left context)	My dog loved tennis balls.
input (right context)	My dog had stolen every one and put it under there.
output text	One day, I found all of my lost tennis balls underneath the bed.

- ▶ for infilling control task, learn a conditional score

$$\nabla_{x_{t-1}} \log p(x_{t-1} | x_t, c)$$

where  $c = (O_L, O_R)$  is left and right context

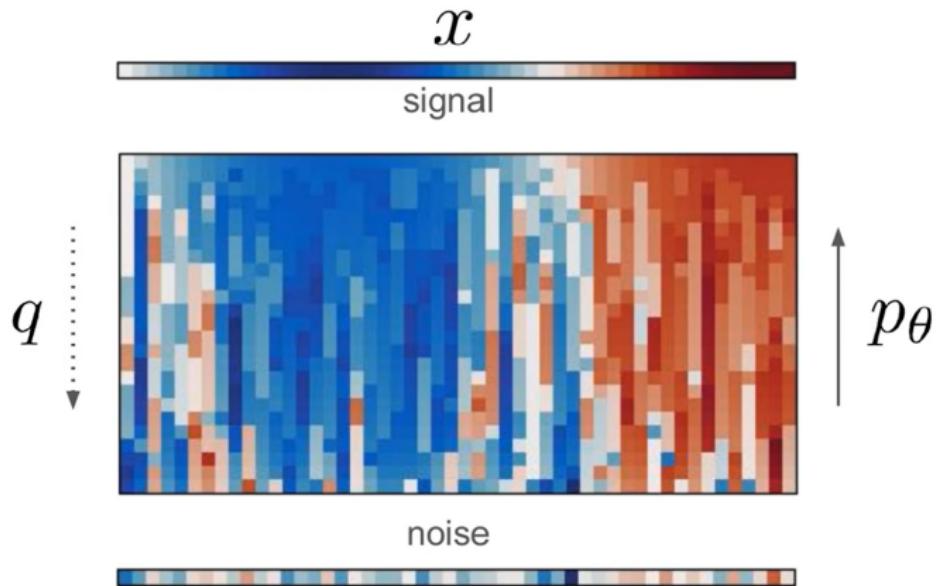
- ▶ this is more similar to question/response text generation in LLMs

# Diffusion LM

	Semantic Content		Parts-of-speech		Syntax Tree		Syntax Spans		Length	
	ctrl ↑	lm ↓	ctrl ↑	lm ↓	ctrl ↑	lm ↓	ctrl ↑	lm ↓	ctrl ↑	lm ↓
PPLM	9.9	5.32	-	-	-	-	-	-	-	-
FUDGE	69.9	2.83	27.0	7.96	17.9	<b>3.39</b>	54.2	4.03	46.9	3.11
Diffusion-LM	<b>81.2</b>	<b>2.55</b>	<b>90.0</b>	<b>5.16</b>	<b>86.0</b>	3.71	<b>93.8</b>	<b>2.53</b>	<b>99.9</b>	<b>2.16</b>
FT-sample	72.5	2.87	89.5	4.72	64.8	5.72	26.3	2.88	98.1	3.84
FT-search	89.9	1.78	93.0	3.31	76.4	3.24	54.4	2.19	100.0	1.83

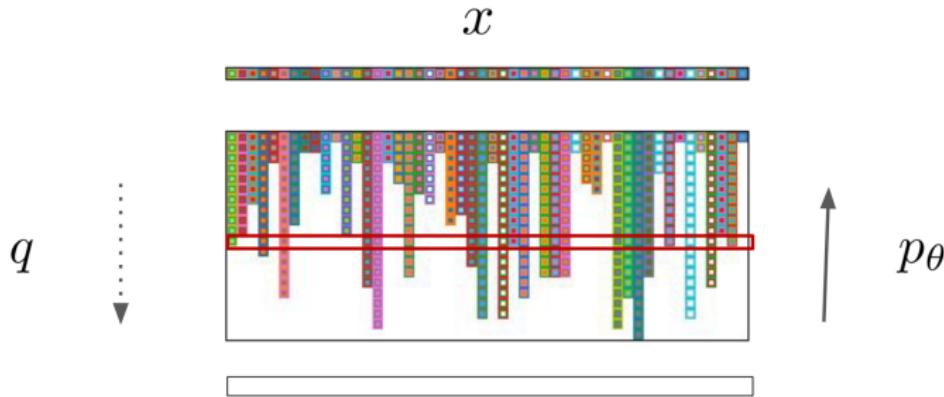
**Figure:** Diffusion-LM achieves high success rate (ctrl) and good fluency (lm) across all 5 control tasks, outperforming the PPLM and FUDGE baselines. Our method even outperforms the fine-tuning oracle (FT) on controlling syntactic parse trees and spans.

# Masked Diffusion LM



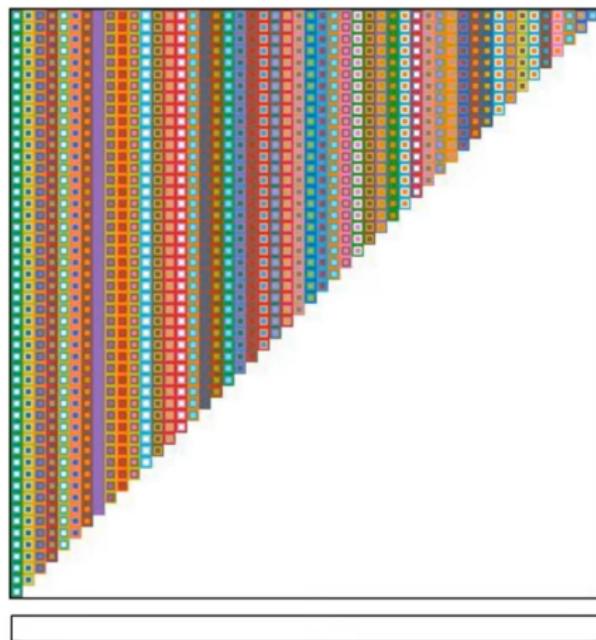
## Masked Diffusion LM

A demonstration of the discrete forward and backward process in the model

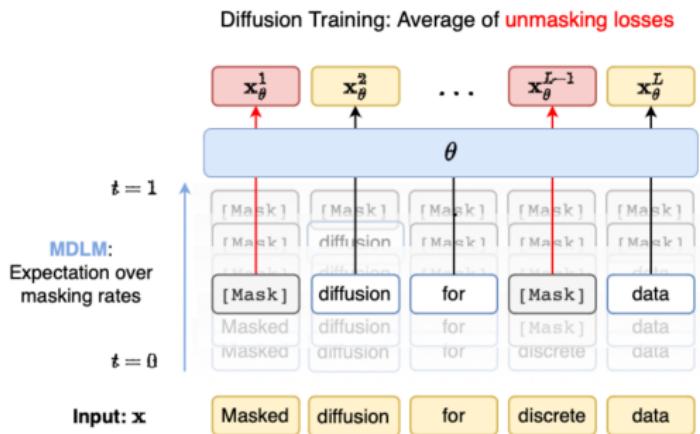


- ▶  $q$  is the measure along the forward process
  - ▶ Limit is “masked”— $[m]$ , a special symbol
  - ▶  $p_\theta$  is the measure along the backward process

## Autoregressive Analogy



# Training Procedure



The forward Markov chain  $\{Z_t\}_{t \in [0,1]}$ ,  $Z_0 = X$ ,  $Z_t = \alpha_t X + (1 - \alpha_t)m$ :

$$P(Z_t | X) = Q_t \dots Q_{t_1} X$$

- ▶ Mixture of data  $X$  and noise  $m$ ; Mixing time control  $\alpha_t$
- ▶ Simulate by a sequence of Markov transition kernels  $Q_{t_i}$ ,  $i = 1, \dots, T$

# Training Objective Function

Information from the training procedure: discrete score function or **A denoising predictor  $x_\theta(Z_t, t)$** . The loss function (continuous version) is

$$\mathcal{L}(\theta) = -\mathbb{E}_{Z_t} \int_0^1 \beta_t \log \langle x_\theta(Z_t, t), X \rangle dt$$

The prediction model is set such that  $\langle x_\theta(Z_t, t), m \rangle = 0$ . Empirical version:

$$\mathcal{L}(\theta) = -\mathbb{E}_{t \sim \text{Unif}[0,1], X \sim p_{\text{data}}} \beta_t \log \langle x_\theta(\alpha_t X + (1 - \alpha_t)m, t), X \rangle dt$$

# Training Objective Function

Information from the training procedure: discrete score function or **A denoising predictor  $x_\theta(Z_t, t)$** . The loss function (continuous version) is

$$\mathcal{L}(\theta) = -\mathbb{E}_{Z_t} \int_0^1 \beta_t \log \langle x_\theta(Z_t, t), X \rangle dt$$

The prediction model is set such that  $\langle x_\theta(Z_t, t), m \rangle = 0$ . Empirical version:

$$\mathcal{L}(\theta) = -\mathbb{E}_{t \sim \text{Unif}[0,1], X \sim p_{\text{data}}} \beta_t \log \langle x_\theta(\alpha_t X + (1 - \alpha_t)m, t), X \rangle dt$$

Implementation of denosing predictor

- ▶ Tokenization is critical to performance, small vocabulary ↓
- ▶ Architectural matters: diffusion transformer (encoder-only, rotary positional embeddings for time conditioning)
- ▶ Choose correlated  $t_i$ , low-discrepancy sequence

# Sampling Procedure

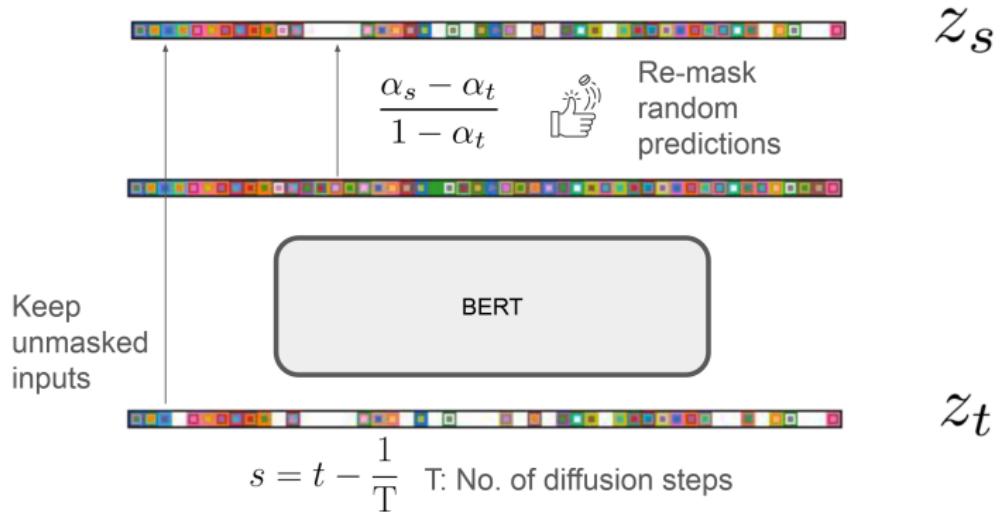
Sampling procedure is a reversed version of forward Markov chain, for any  $s < t$

$$Z_s | Z_t = \begin{cases} Z_t^i & \text{if } Z_t^i \neq m \\ \frac{1}{1-\alpha_t}((1-\alpha_s)x_\theta(Z_t, t)^i + (\alpha_s - \alpha_t)m) & \text{if } Z_t^i = m \end{cases}$$

Entrywise denoising process based on the rule:

1. Begin with pure noise, all tokens are masked
2. When a token is unmasked ( $\neq m$ ), it remains unchanged during reverse diffusion
3. For masked token, w.p.  $(\alpha_s - \alpha_t)/(1 - \alpha_t)$  unmask it and apply denoiser  $x_\theta$

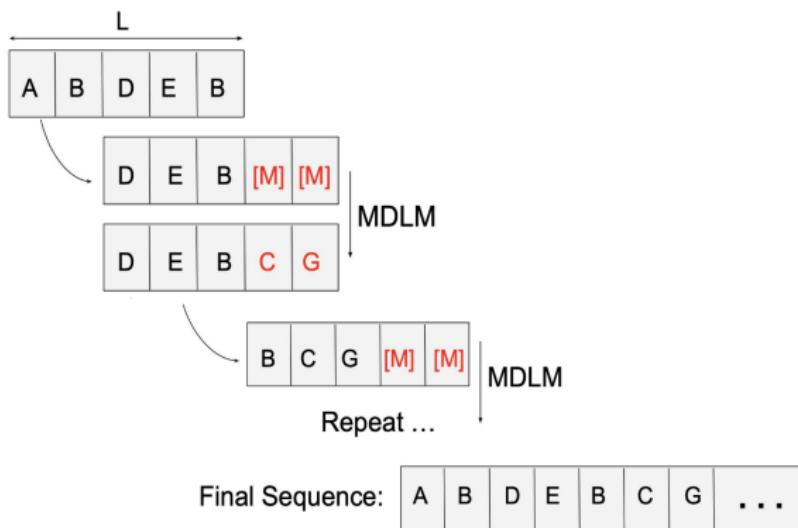
## One Step of Sampling Procedure



## Extension: Semi-autoregressive Method

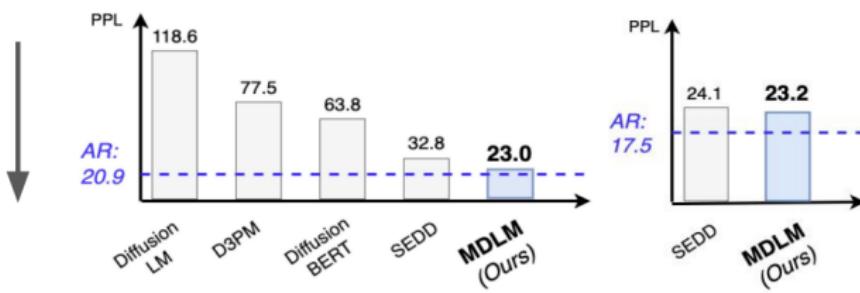
Instead of generating the whole sequence at once, generate part of it and continue to generate sequence of arbitrary length

- Add "masked" somewhere in the sequence and run the backward process again



# Experiments

- ▶ Tested on standard One Billion Words (LM1B) and OpenWebText (OWT) datasets; evaluate perplexity
- ▶ ~ 110M parameters; trained on ~ 300B tokens; optional time conditioning
- ▶ Beat other diffusion-based models, 10% ~ 25% worse than AR models
- ▶ faster sampling compared with other diffusion-based models, slower than AR models (same perplexity)

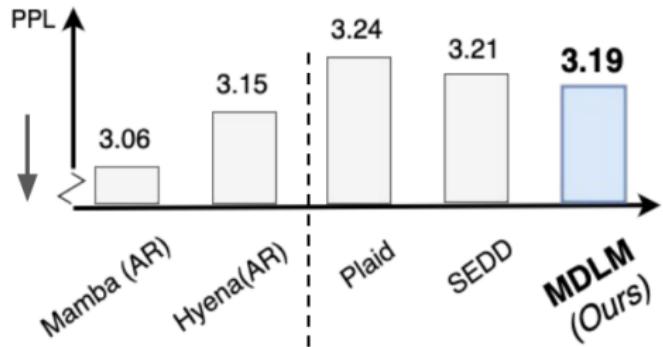


LM1B

OpenWebText

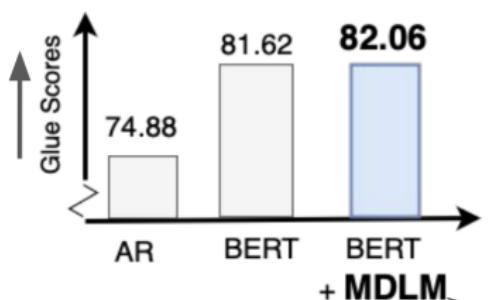
## Experiments

- ▶ Apply to the generative modeling of DNA sequences; finetune Caduceus DNA model with both AR and diffusion LMs
  - ▶ ~ 500K parameters, trained on the HG38 human genome dataset
  - ▶ Evaluate on the downstream bio tasks — weak correlation between performance and perplexity



## Experiments

- ▶ Tested on General Language Understanding Evaluation (GLUE): sentiment analysis, linguistic acceptability...
  - ▶ Score: Improve BERT a little bit, better than AR models
  - ▶ Perplexity: Diffusion LM fine tuning reduce perplexity ( $78 \rightarrow 35$ ); still higher than AR models ( $\sim 23$ )



Trained BERT on C4  
Finetuned with MDLM  
Evaluated it on GLUE

Generative BERT

# Summary of the Diffusion LM Frameworks

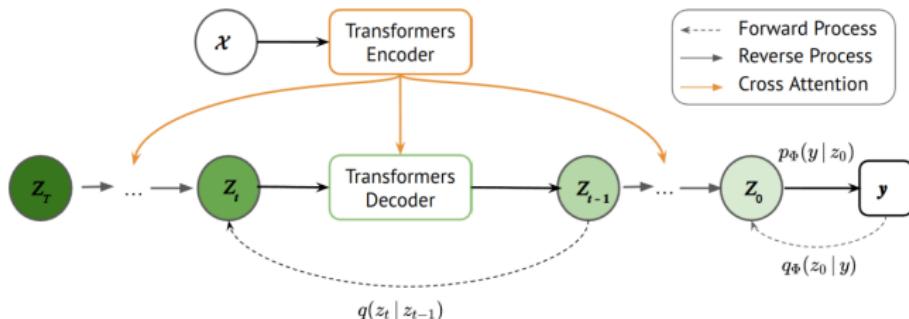
- ▶ **Parallel Generation:** sequentially conditioned on preceding tokens → parallel generation of all output tokens. Enhance the speed and efficiency of some real world applications
- ▶ **Text Interpolation:** generate intermediate sentences between two given parts, smooth transitions and coherent outputs.
- ▶ **Token-level control:** conditional diffusion process
- ▶ **Robust to corruption:** denoising mechanism aids in mitigating errors and noise present in the input sequence
- ▶ **Different Paradigm...**

# Summary of the Diffusion LM Frameworks

- ▶ **Parallel Generation**: sequentially conditioned on preceding tokens → parallel generation of all output tokens. Enhance the speed and efficiency of some real world applications
- ▶ **Text Interpolation**: generate intermediate sentences between two given parts, smooth transitions and coherent outputs.
- ▶ **Token-level control**: conditional diffusion process
- ▶ **Robust to corruption**: denoising mechanism aids in mitigating errors and noise present in the input sequence
- ▶ **Different Paradigm...**

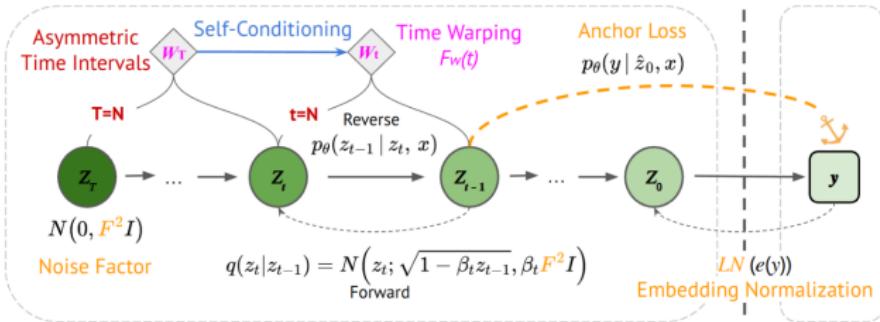
Models	Advantages				Disdvantages	
	Parallel Generation	Text Interpolation	Token-level Controls	Robustness to Input Corruption	Training Complexity	Model Interpretability
AR Models	✗	✗	✗	✗	$\mathcal{O}(n)$	😊
Diffusion Models	✓	✓	✓	✓	$\mathcal{O}(n^T)$	😊

# Summary of the Diffusion LM Frameworks



- ▶ transformers use bi-directional attention
- ▶ encoding used for contextual embeddings
- ▶ has spurred *lots* of algorithms & techniques for handling discrete data

## Summary of the Diffusion LM Frameworks



**Anchor Loss, Noise Factor, Embedding Normalization:** *Difformer*

### Asymmetric Time Intervals: Bit Diffusion

### **Self-Conditioning: Bit Diffusion, SED, CDCCD, SeqDiffuseSeq, DiNojSeq**

**Time Warning:** CPDCD, SeqDiffuSeq

- ▶ Broadly fall into two categories (i) adapting models to discrete variables, and (ii) improved sampling procedures

# Summary of the Diffusion LM Frameworks

Model	Tasks	Schedule	Sampling
<b>Discrete Diffusion Models</b>			
Multinomial Diffusion (Hoogeboom et al., 2021)	unconditional text generation, unsupervised spell-checking	Transition matrices	—
D3PM (Discrete Denoising Diffusion Probabilistic Models) (Austin et al., 2021b)	char-level text and image generation	Uniform Transition Matrices	—
Zero-shot Diffusion (Nachmani and Dovrat, 2021)	machine translation	Partial Noising	Classifier-free conditional denoising
SUNDAE (Step-unrolled Denoising Autoencoders) (Savinov et al., 2021)	machine translation and unconditional text generation	Uniform Transition Matrices	Low-temperature sampling, Argmax-unrolled decoding, fewer token update
DiffusionBERT (He et al., 2022)	unconditional text generation	Spindle	$\times$ 0-parameterization
SSD-LM (Semi-autoregressive Simplex-based Diffusion) (Han et al., 2022)	unconditional and controlled text generation	Logits generation	Greedy projection, Sampling, Multi-hot
Bit Diffusion (Generating Discrete Data using Diffusion Models with Self-Conditioning) (Chen et al., 2023b)	categorical image generation and image captioning	—	Self-Conditioning, Asymmetric Time Intervals
DiffusER (Discrete Diffusion via Edit- and-Refine) (Zou et al., 2023)	machine translation, summarization	Edit-based Corruption	Beam Search, 2D Beam

For more architectures, see Zou et al. "A survey of diffusion models in natural language processing."

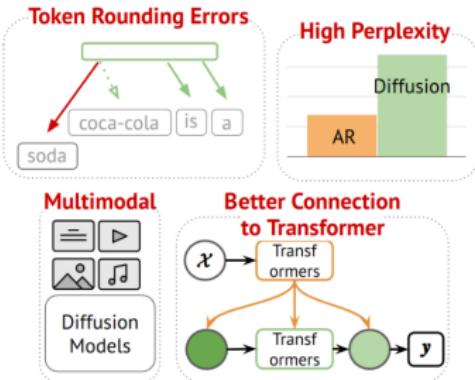
# Summary of the Diffusion LM Frameworks

Size and computational cost of models?

Model	# Parameters	GPU Hours
Stable Diffusion 2.0	1B	≈ 200k A100 GPU hours
Diffusion-LM (variant)	500M	≈ 50k A100 GPU hours
GPT-3	175B	≈ 1,000k A100 GPU hours

There's a gap in the computational aspect, so text-based diffusion models have potential to perform well

## Summary of the Diffusion LM Frameworks



- ▶ connecting spatiality of transformers and temporality of diffusion models, e.g., including time step  $t$  into transformer
  - ▶ diffusion models naturally excel in image generation, how to create more semantically accurate visual outputs?

## References I

- 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.
- Subham Sahoo, Marianne Arriola, Yair Schiff, Aaron Gokaslan, Edgar Marroquin, Justin Chiu, Alexander Rush, and Volodymyr Kuleshov. Simple and effective masked diffusion language models. *Advances in Neural Information Processing Systems*, 37:130136–130184, 2025.
- Hao Zou, Zae Myung Kim, and Dongyeop Kang. A survey of diffusion models in natural language processing. *arXiv preprint arXiv:2305.14671*, 2023.