Diffusion Models

By Tejomay Kishor Padole (PhD Student, IIT Bombay)

Generative Modeling

• Learning to generate new samples given a set of data points.

Estimate: $p(x_i) \ \forall x_i \in \mathcal{D}$

Generate: $x_{\mathrm{new}} \sim p(x_i)$

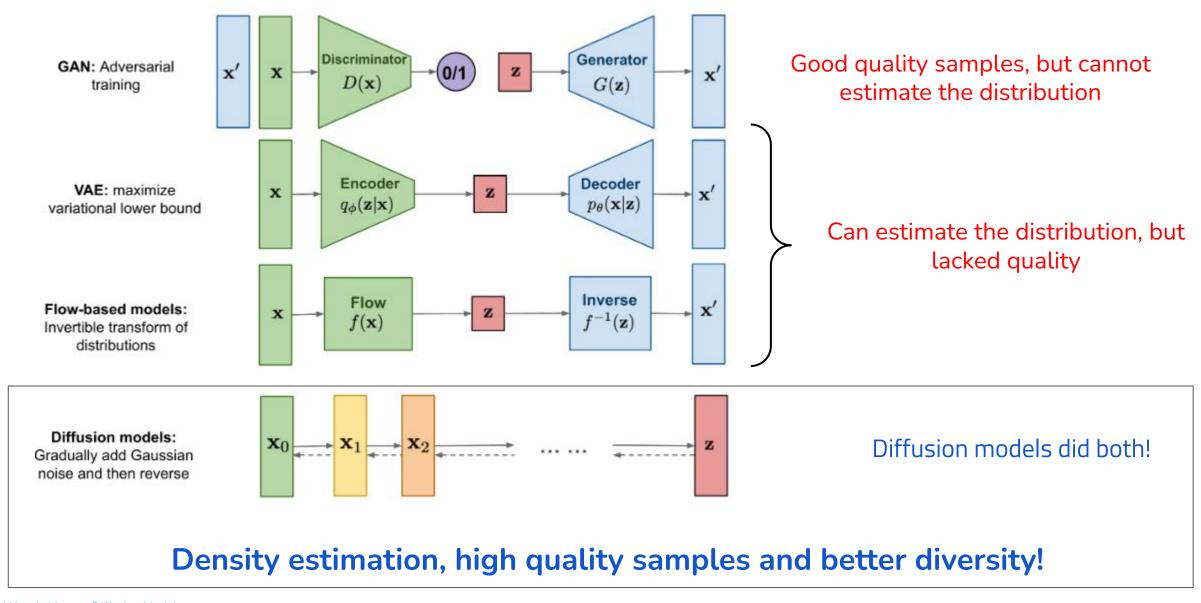
Autoregressive LLMs, Variational Autoencoders, Normalizing flows

Explicitly model the data distribution

Generative Adversarial
Networks

Modeling the generation process

Generative Models



Why Diffusion Models are Trending?

Diffusion models enabled high resolution image generation!



Diffusion Models

Inspired from non-equilibrium statistical physics.

• **Forward process:** slowly and iteratively destroys the structure of the input until it gets completely transformed to random noise.

• Reverse process: learns to iteratively construct the data from random noise using a Neural Network.

• The noising process can be seen as similar to a drop of ink diffusing into clear water.



Generative Process as a form of Iterative Denoising

image

Noising Process

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I)$$
 (Markovian Gaussian transitions)



 $p_{data}(x)$



 $\mathcal{N}(0,I)$

Scaling factor for variance preservation

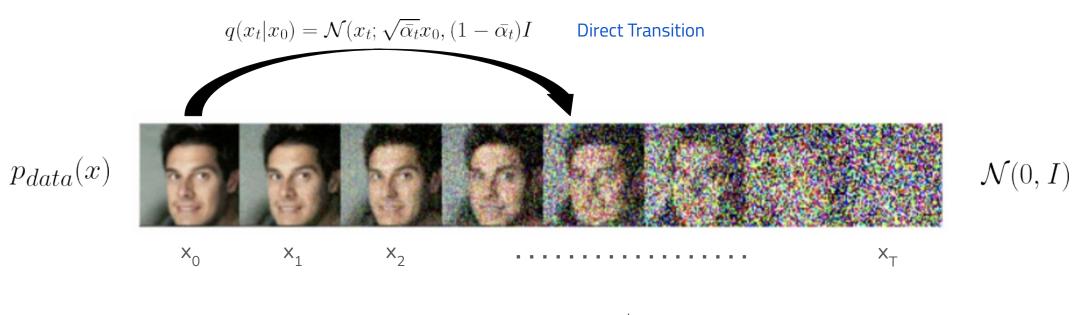
$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I) \xrightarrow{}$$

$$x_t = \sqrt{1-\beta_t}x_{t-1} + \sqrt{\beta_t}\epsilon \text{ where } \epsilon \sim \mathcal{N}(0, I)$$
 Downscaled Noise previous timestep component

Variance of the forward process, monotonic increasing function of t

Generative Process as a form of Iterative Denoising

Noising Process



We define
$$\alpha_t = 1 - \beta_t$$
 and $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1-\bar{\alpha}_t)I)$$
 Part of the original mage Part of random gaussian noise

Generative Process as a form of Iterative Denoising

Noising Process

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I) \qquad \text{(Markovian Gaussian transitions)}$$

$$p_{data}(x)$$

$$x_0 \qquad x_1 \qquad x_2 \qquad \dots \qquad x_T$$

$$p(x_T) = \mathcal{N}(x_T; 0, I) \quad \text{and} \quad p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t))$$
 Parameterized with a Neural Network

How to learn this?

Learning the Generative Process



Generative Process

Initialize: $p(x_T) = \mathcal{N}(x_T; 0, I)$ and

Iteratively sample from: $p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t))$

$$-\log p_{\theta}(x_{0}) \leq -\log p_{\theta}(x_{0}) + D_{\mathbf{KL}}(q(x_{1:T}|x_{0}) \mid\mid p_{\theta}(x_{1:T}|x_{0}))$$

$$= -\log p_{\theta}(x_{0}) + \mathbb{E}_{x_{1:T} \sim q(x_{1:T}|x_{0})} \left[\log \frac{q(x_{1:T}|x_{0})}{p_{\theta}(x_{1:T}|x_{0})} \right]$$

$$= -\log p_{\theta}(x_{0}) + \mathbb{E}_{x_{1:T} \sim q(x_{1:T}|x_{0})} \left[\log \frac{q(x_{1:T}|x_{0})}{\frac{p_{\theta}(x_{0:T})}{p_{\theta}(x_{0})}} \right]$$

$$= \mathbb{E}_{x_{1:T} \sim q(x_{1:T}|x_{0})} \left[\log \frac{q(x_{1:T}|x_{0})}{p_{\theta}(x_{0:T})} \right]$$

(Variational Lower Bound on log-likelihood of data)

Expanding KL-divergence

Re-writing the denominator

Got rid of the log-likelihood term!

Learning the Generative Process



Generative Process

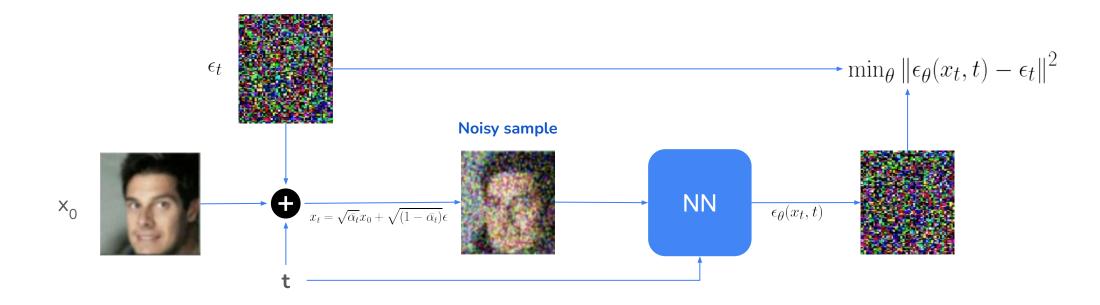
Final Loss Function

$$L_{VLB} = \mathbb{E}_q \left[\underbrace{D_{\mathrm{KL}}(q(x_T, x_0) \mid\mid p_{\theta}(x_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\mathrm{KL}}(q(x_{t-1} \mid x_t, x_0) \mid\mid p_{\theta}(x_{t-1} \mid x_t))}_{L_{t-1}} - \underbrace{\log p_{\theta}(x_0 \mid x_1)}_{L_0} \right]$$

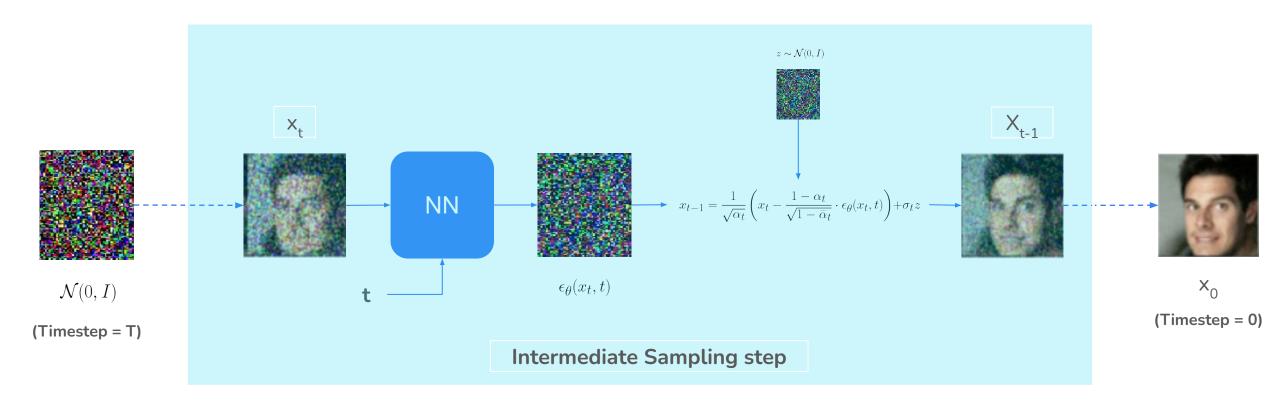
$$= \mathbb{E}_{x_0} \left[\frac{1}{2 \mid\mid \Sigma_{\theta} \mid\mid_2^2} \mid\mid \tilde{\mu}(x_t, x_0) - \mu_{\theta}(x_t, t) \mid\mid_2^2 \right] \qquad \text{Closed-form solution for KLD between two gaussians (and variance ignored)}$$

$$= \mathbb{E}_{x_t, \epsilon_t} \mid\mid \epsilon_t - \epsilon_{\theta}(x_t, t) \mid\mid^2 \qquad \qquad \text{Reparameterized as:} \quad \mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(x_t, t))$$

Training Algorithm



Sampling (Generation) Algorithm



The Score Function and Diffusion Models

The Score Function

$$x_i \sim p(x)$$

Goal: Estimate the distribution with a NN

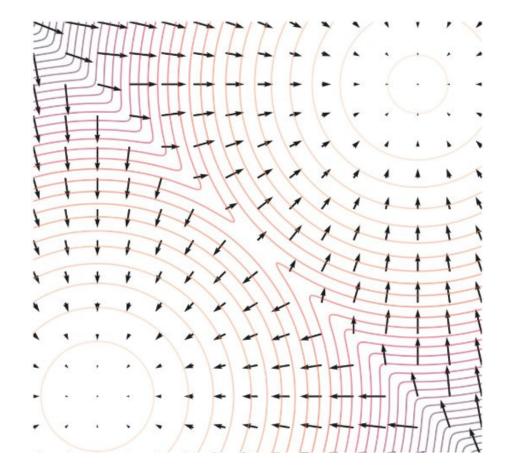


$$p_{\theta}(x_i) = \frac{\exp\left(-f_{\theta}(x_i)\right)}{\boxed{Z_{\theta}}} \underbrace{\begin{array}{c} \text{Intractable!!} \\ \\ Z_{\theta} \end{array}}_{} Z_{\theta} = \int_{x} \exp(-f_{\theta}(x)) dx$$

$$\ln p_{\theta}(x_i) = -f_{\theta}(x_i) - \ln Z_{\theta}$$

$$\nabla x_i \ln p_{\theta}(x_i) = -\nabla x_i f_{\theta}(x_i) - \nabla x_i \mathcal{Z}_{\theta}$$

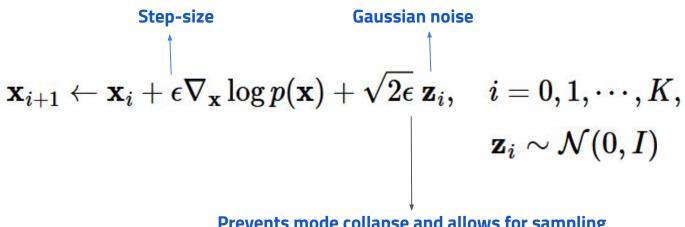
Score function



Score function interpretation on a mixture of two gaussians in 2D.

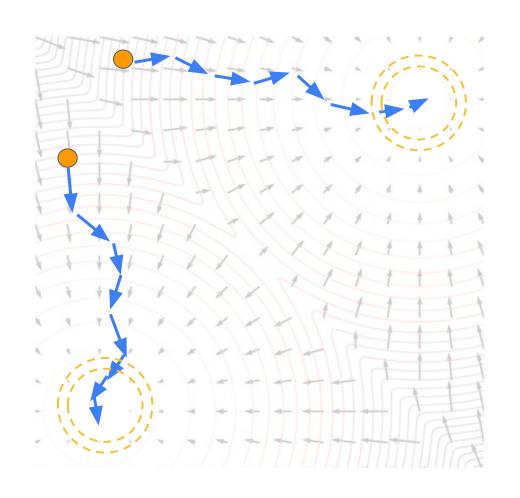
Interpretation of the Score Function

 Learning the score function allows sampling from the data distribution.



Prevents mode collapse and allows for sampling from regions with low-data density

Langevin Dynamics



Score-based Diffusion Models

Noising Process

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}$$





 $\mathcal{N}(0,I)$

score function

$$\mathrm{d}\mathbf{x} = \left[\mathbf{f}(\mathbf{x},t) - g^2(t) \nabla_{\mathbf{x}} \log p_t(\mathbf{x})\right] \mathrm{d}t + g(t) \mathrm{d}ar{\mathbf{w}}$$

Generative Process

$$\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x},t)$$

(Time-conditioned score network)

Training objective:
$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \mathbb{E}_t \Big\{ \lambda(t) \mathbb{E}_{\mathbf{x}(0)} \mathbb{E}_{\mathbf{x}(t)|\mathbf{x}(0)} \big[\| \mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_{0t}(\mathbf{x}(t) \mid \mathbf{x}(0)) \|_2^2 \big] \Big\}$$

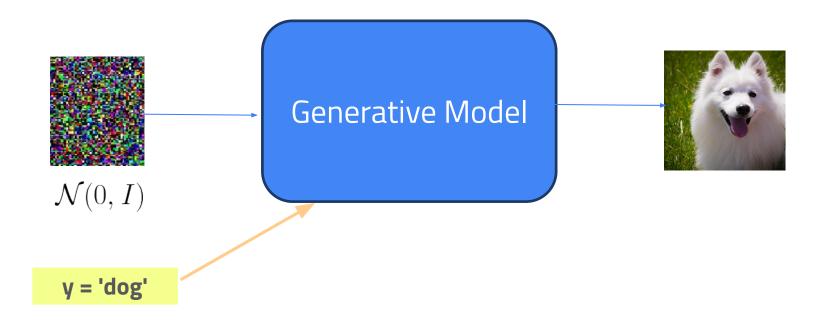
Conditioning of Diffusion Models

Class Conditional Diffusion Models

Generate: $x_{\mathrm{new}} \sim p(x_i)$ (Unconditional)

What if we want to generate: $x_{new} \sim p(x_i|y)$ (where 'y' is a class label)

General strategy: Train the generative model with class awareness.



Classifier Guided Diffusion

We want: $x_{\text{new}} \sim p(x|y)$ (where 'y' is a class label)

model)

$$p(x \mid y) = rac{p(y \mid x) \cdot p(x)}{p(y)}$$
 $\log p(x \mid y) = \log p(y \mid x) + \log p(x) - \log p(x)$
Label probabilities (can be obtained from a pre-trained Unconditional data probability (modeled using a diffusion

Dhariwal and Nichol, 2021 propose an additional guidance scale:

$$\log p_{\gamma}(x|y) \propto \log p(x) + \log p(y|x)$$
 Guidance scale

image classifier

Classifier Guided Diffusion

Varying the guidance scale allows us to have a trade-off between diversity and fidelity.

$$abla_x \log p_\gamma(x \mid y) =
abla_x \log p(x) + \gamma
abla_x \log p(y \mid x).$$

$$\gamma = 1.0$$
 $\gamma = 10.0$





Generated images from the class "Pembroke Welsh Corgi" with different guidance scales.

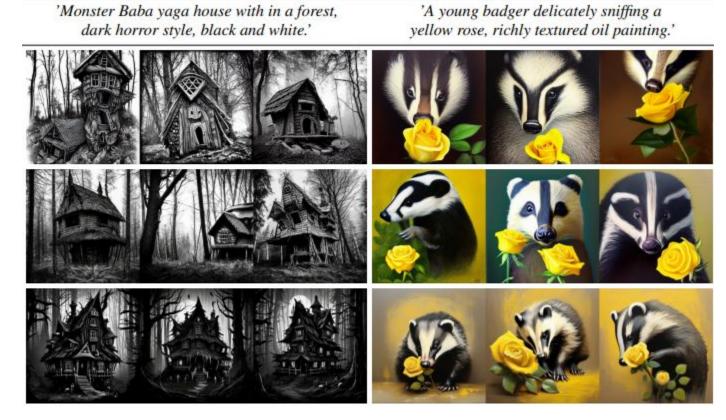
Arbitrary Conditioning of Diffusion Models

Classifier guidance doesn't allow for arbitrary conditioning mechanisms.

What if we want to generate:

$$x_{new} \sim p(x_i|y)$$

(where 'y' can be natural language text, image layout, etc.)



Classifier-free Guidance

What if we want to generate:

$$x_{new} \sim p(x_i|y)$$

(where 'y' can be natural language text, image layout, etc.)

Re-formulate:

$$p(y \mid x) = rac{p(x \mid y) \cdot p(y)}{p(x)}$$

$$\log p(y \mid x) = \log p(x \mid y) + \log p(y) - \log p(x)$$

From classifier guidance:

$$\log p_{\gamma}(x|y) \propto \log p(x) + \gamma \log p(y|x)$$

$$\log p_{\gamma}(x|y) \propto \log p(x) + \gamma(\log p(x|y) + \log p(y) - \log p(x))$$

$$\log p_{\gamma}(x|y) \propto (1-\gamma) \log p(x) + \gamma \log p(x|y)$$

Train a conditional and an unconditional model in tandem.

Classifier-free Guidance

$$\log p_{\gamma}(x|y) \propto (1-\gamma) \log p(x) + \gamma \log p(x|y)$$

 $\gamma = 1.0$ $\gamma = 3.0$



Samples from OpenAl GLIDE text-to-image model (*Nichol et al.*, 2021) with classifier-free guidance.

Prompt: A stained glass window of a panda eating bamboo.

Diffusion Models For Text

Why Diffusion Models for Text?



Autoregressive model: suffers from sampling drifts, major cause of hallucination

Humans do not necessarily process text sequentially.

Eg: The food service of the restaurant was brilliant. The plate was cold and so was the food.

Sequential generation does not take advantage of GPU parallelism (which can boost generation in long-contexts)

What are Diffusion Models doing Differently?

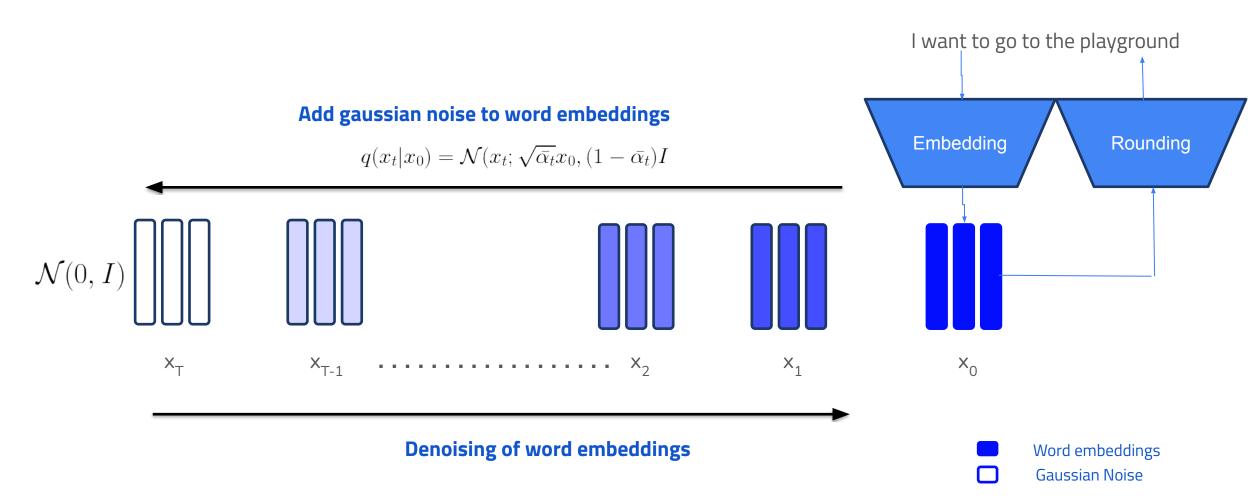
Non-autoregressive Generation: Generate all tokens in parallel Iterative Denoising: allows reconsideration of already generated tokens

Exhibit better control

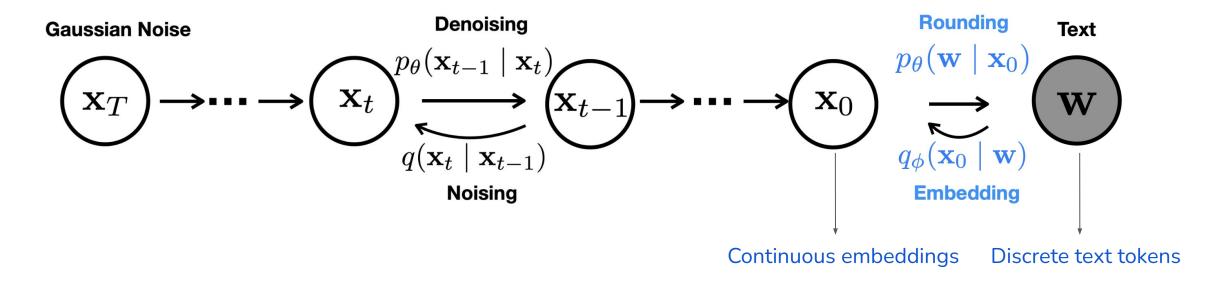
Continuous Diffusion Models For Text

Continuous Diffusion for Text

• Gaussian diffusion processes defined on word embeddings.



Continuous Diffusion for Text



Embedding step: Maps discrete tokens to an embedding space.

$$q_{\phi}(x_0|w) = \mathcal{N}(\mathcal{E}(w), \sigma_0 I)$$

• **Rounding step:** Finds out the most likely token given an embedding (reverse of what embedding step does).

$$p_{\theta}(w|x_0) = \prod_{i=1}^n p_{\theta}(w_i|x_i)$$
 (each component represents a position in the sentence modeled with a softmax function)

The rest of the forward process is defined in the same way as shown earlier.

Training Continuous Diffusion for Text

Loss =
$$L_{VLB}$$
 + $E_{q_{\phi}(x_{0}|w)}$ $\log \frac{q_{\phi}(x_{0}|w)}{p_{\theta}(w|x_{0})}$
= $E_{x_{t} \sim q(x_{t}|x_{0}), x_{0} \sim q_{\phi}(x_{0}|w)}$ $\left[\|\epsilon_{\theta}(x_{t}, t) - \epsilon_{t}\|^{2} + \log q_{\phi}(x_{0}|w) - \log p_{\theta}(w|x_{0}) \right]$
= $E_{x_{t} \sim q(x_{t}|x_{0}), x_{0} \sim q_{\phi}(x_{0}|w)}$ $\left[\|\epsilon_{\theta}(x_{t}, t) - \epsilon_{t}\|^{2} + \|\mathcal{E}(w) - \mu_{\theta}(x_{1}, 1)\|^{2} - \log p_{\theta}(w|x_{0}) \right]$

The first term can be reparameterized by letting the neural network directly predict $\mathbf{x_0}$ instead of predicting **noise**. Li et al., 2022 found that doing so reduced rounding errors.

$$= \underset{x_t \sim q(x_t|x_0), \ x_0 \sim q_{\phi}(x_0|w)}{\mathbb{E}} \left[\|h_{\theta}(x_t, t) - x_0\|^2 + \|\mathcal{E}(w) - \mu_{\theta}(x_1, 1)\|^2 - \log p_{\theta}(w|x_0) \right]$$

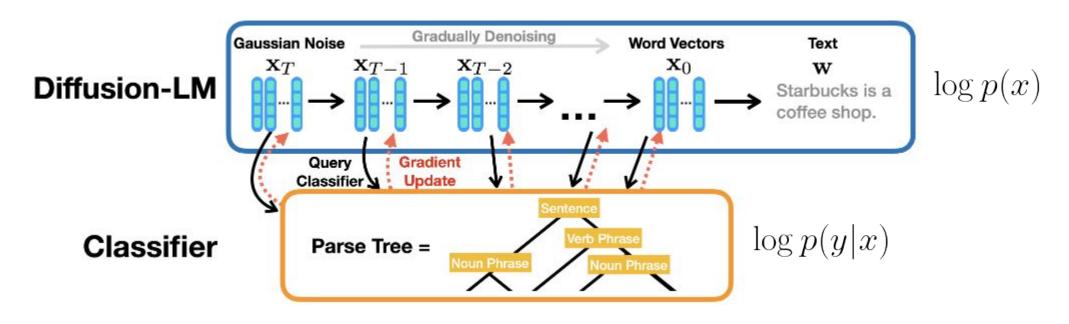
$$\downarrow \qquad \qquad \qquad \downarrow$$
predicted x₀ at timestep t

Cross-entropy loss

Controllable Continuous Diffusion for Text

We can develop controllable diffusion LMs with classifier guidance!

$$x_{new} \sim p(x_i|y)$$
 ('y' is a control signal)



$$\log p_{\gamma}(x|y) \propto \log p(x) + \gamma \log p(y|x)$$

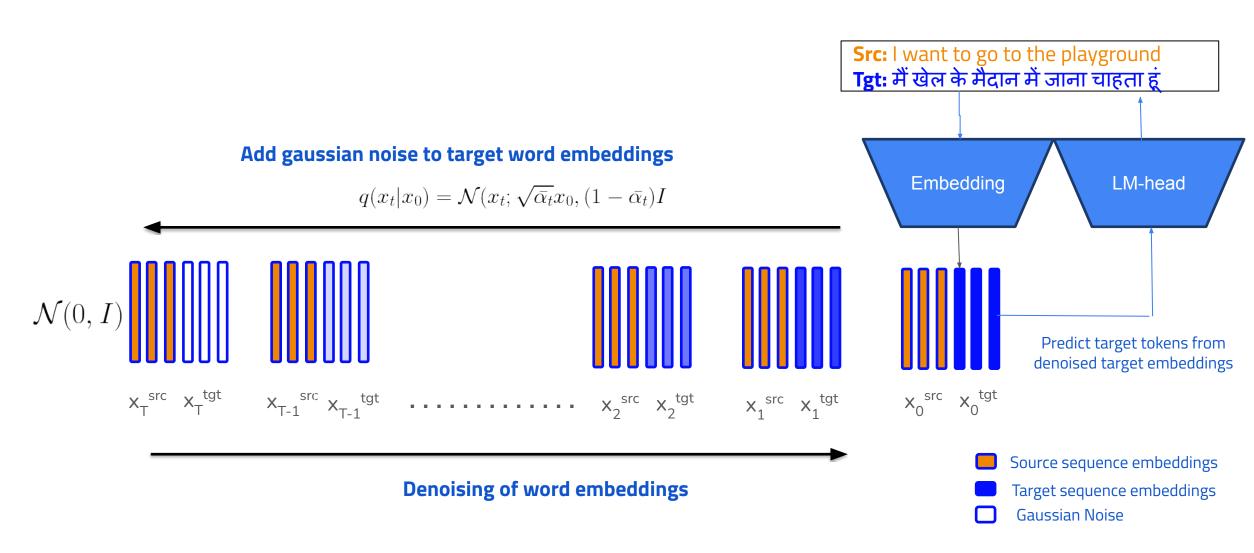
Controllable Continuous Diffusion for Text

We can develop controllable diffusion LMs with classifier guidance!

target POS	PROPN AUX DET ADJ NOUN NOUN VERB ADP DET NOUN ADP DET NOUN PUNCT
FUDGE	Aromi is a non family - friendly fast food coffee shop in the riverside area with a low Customer Rating.
Diffusion-LM FT	Fitzbillies is a cheap coffee shop located on the outskirts of the city. Aromi is a fast food pub located at the centre of the city.
target POS	PROPN AUX DET NOUN VERB NOUN ADJ NOUN PUNCT PRON NOUN NOUN AUX ADJ
FUDGE	Cocum is a family - friendly coffee shop, that has a low price range and a low customer rating
Diffusion-LM FT	Zizzi is a pub providing restaurant Chinese food. Its customer rating is low Zizzi is a pub providing kids friendly services. Its customer rating is average

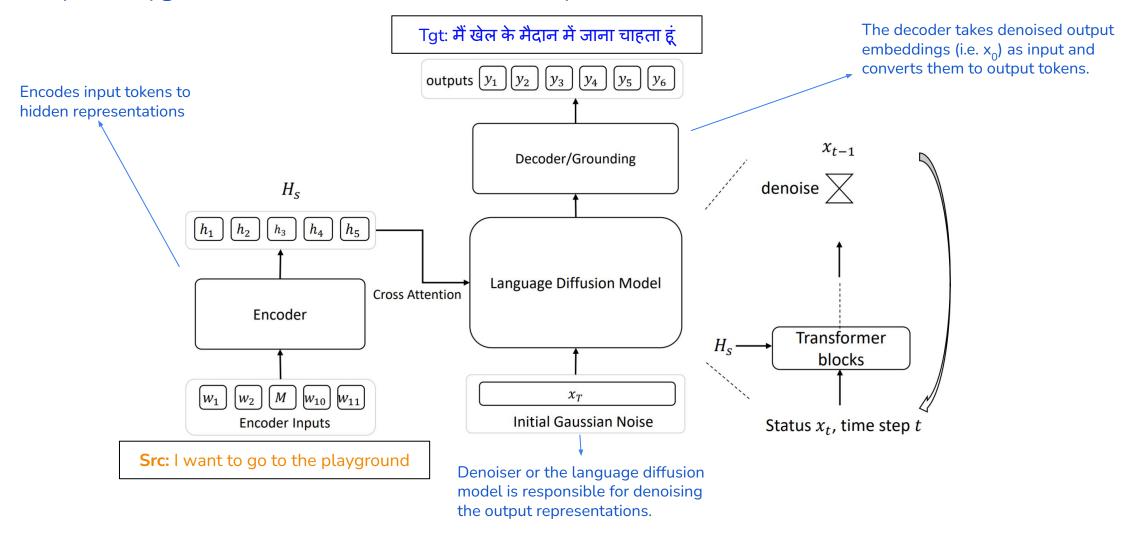
DiffuSeq (Gong et al., 2022)

Gaussian diffusion processes on word embeddings corresponding to the target sequence.



GENIE(Lin et al., 2023)

Seq-to-seq generation with encoder-decoder style architecture.



Discrete Diffusion Language Models

Discrete Diffusion

- Can we directly denoise in the token space (which is discrete)?
- No need to define extra embedding and rounding steps which are harder to train jointly with the diffusion model.

Continuous Diffusion	Discrete diffusion
Add Gaussian noise (continuous distribution)	Add discrete noise (what is discrete noise)?
Diffusion process converges to a gaussian	Diffusion process converges to a categorical distribution (which categorical distribution)?

Masked Diffusion Language Models

Defines the diffusion processes as iterative masking/unmasking of tokens.

 $q(x_t|x_{t-1})$

Noising Process

What do you do in the morning? I go [M] a morning walk with my dog every day.

What do you do in the morning? I go [M] a morning walk [M] my dog every day.

What do you do in the morning? I go [M] [M] morning walk [M] my [M] every day.

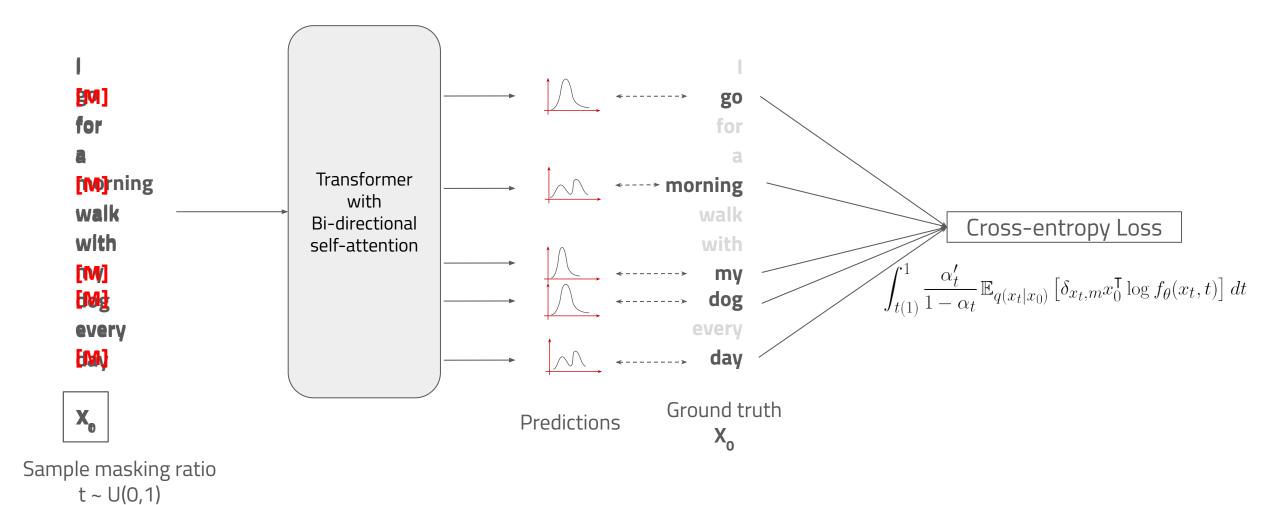
What do you do in the morning? I go [M] [M] morning walk [M] my [M] every day.

What do you do in the morning? [M] go [M] [M] [M] [M] walk [M] my [M] [M] [M].

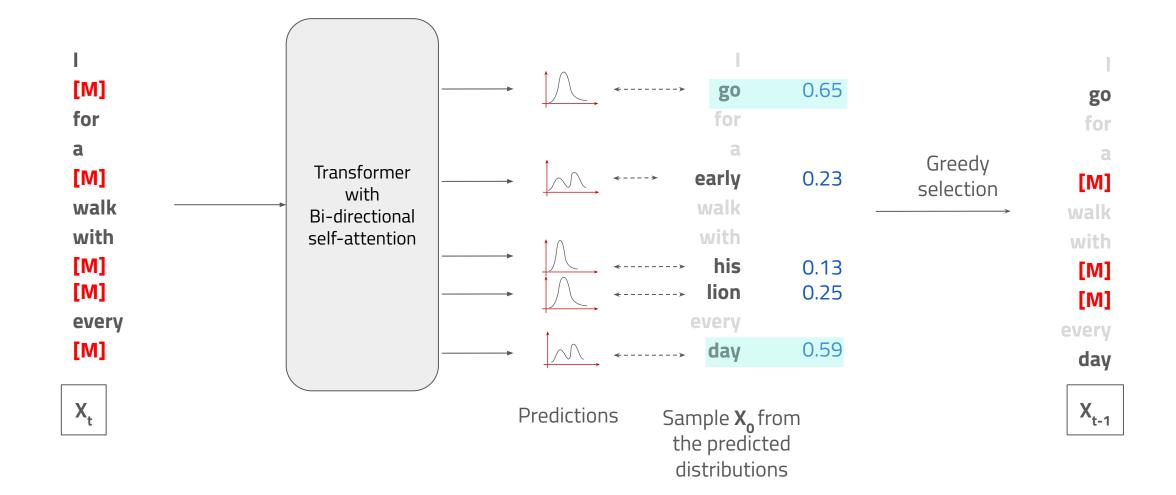
 $p_{\theta}(x_{t-1}|x_t)$

Generative Process

Training Masked Diffusion Language Models



Greedy Sampling Step in Masked Diffusion Language Models

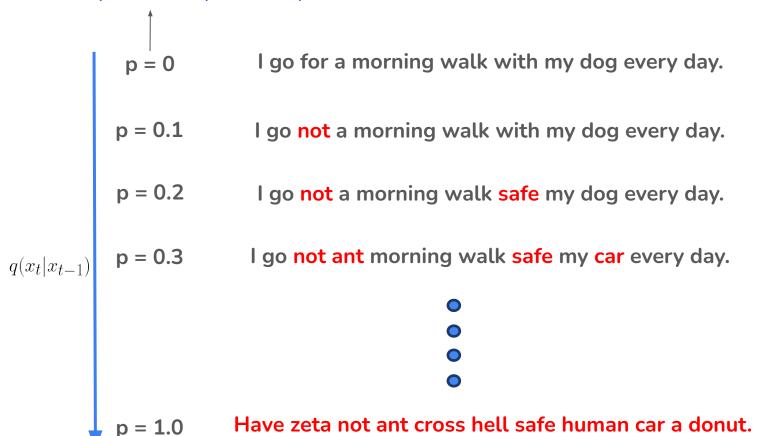


Uniform Diffusion Language Models

Vocabulary

A, ability, able, ... z, zeta

Random replacement probability



Every position in this sentence is simply a uniform randomly sampled word from the vocabulary

 $p_{\theta}(x_{t-1}|x_t)$

UDLM: Formal Definition

Let
$$x = \langle x^{(1)}, x^{(2)}, \dots, x^{(n)} \rangle$$
 (sequence of tokens)

Let x_o be a one-hot vector of length equal to vocabulary size V

$$q(x_t^i|x_0^i) = \alpha_t x_0^i + (1-\alpha_t)\pi, \text{ where } \pi = \frac{\mathbb{I}}{|V|} \qquad \text{Each value in vector } \pi \text{ is 1/|V|} \qquad \text{Forward process}$$

$$q(x_{t-1}|x_t,x_0) - f_{\theta}(x_t,t) = \mathbf{NN}_{\theta}(x_t,t)$$
 Estimate of \mathbf{x}_0
$$p_{\theta}(x_{t-1}|x_t) = q(x_{t-1}|x_t,f_{\theta}(x_t,t))$$

$$\text{Loss Function: } \mathcal{L}^{\infty} = \int_{t=0}^{t=1} \mathbb{E}_{q} \left[\frac{\alpha_{t}'}{N\alpha_{t}} \left[\frac{N}{\bar{\mathbf{x}}_{i}} - \frac{N}{(\bar{\mathbf{x}}_{\theta})_{i}} - \sum_{j \text{ s.t. } (\mathbf{z}_{t})_{j} = 0} \left(\frac{\bar{\mathbf{x}}_{j}}{\bar{\mathbf{x}}_{i}} \right) \log \left[\left(\frac{(\bar{\mathbf{x}}_{\theta})_{i} \cdot \bar{\mathbf{x}}_{j}}{(\bar{\mathbf{x}}_{\theta})_{j} \cdot \bar{\mathbf{x}}_{i}} \right) \right] \right] \right] \mathrm{d}t$$



Diffusion vs. Autoregression

Diffusion vs. Autoregression

Autoregressive LMs	Diffusion LMs
Sequential generation	Non-autoregressive generation with iterative denoising
Harder to control, exhibit sampling drifts.	Easier to control
Lot of innovation (Eg: KV-caching, group-query attention, etc.)	Nascent stage of research, not completely practical yet

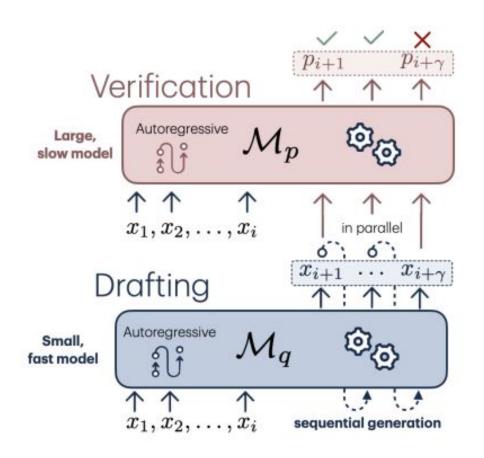
Can we combine them?

Can we Combine Autoregressive LLMs and Diffusion LMs?

- Why combine? To try to make use of the qualities of both.
- **Challenge:** Completely different training paradigms!
- Ways to combine:
 - Diffusion model as a helper for LLM
 - Adapting an LLM to a diffusion model

Speculative Decoding

Employs a small LM to help parallelize token generation from a larger model.

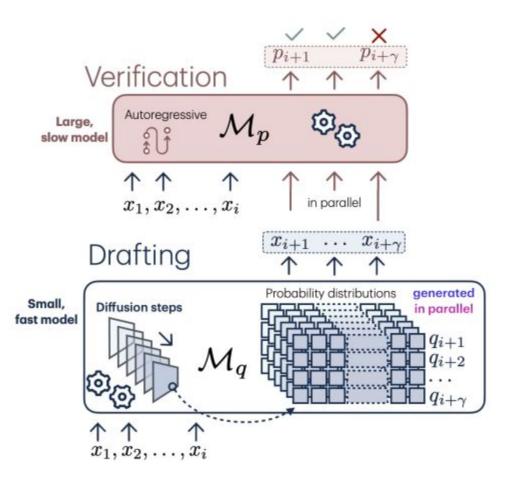


Verify candidate sequences in parallel with a larger model

Generate candidate sequences sequentially from a smaller but faster model

Speculative Diffusion Decoding

Use a small diffusion LM to generate potential candidates.



Quality vs speed trade-off can be done by tuning the diffusion steps

Generate candidate sequences non-autoregressively with a small diffusion model

Adapting LLMs to Diffusion LMs

- Given that we have a lot of pre-trained autoregressive LLMs, can we take advantage of their pre-training and try adapting them to diffusion LMs?
- Naive fine-tuning won't work due to difference in training paradigms.
- How to bring these training paradigms closer?
- Gong et al., 2024 proposes masked diffusion LMs as a solution!

Adapting LLMs to Masked Diffusion LMs

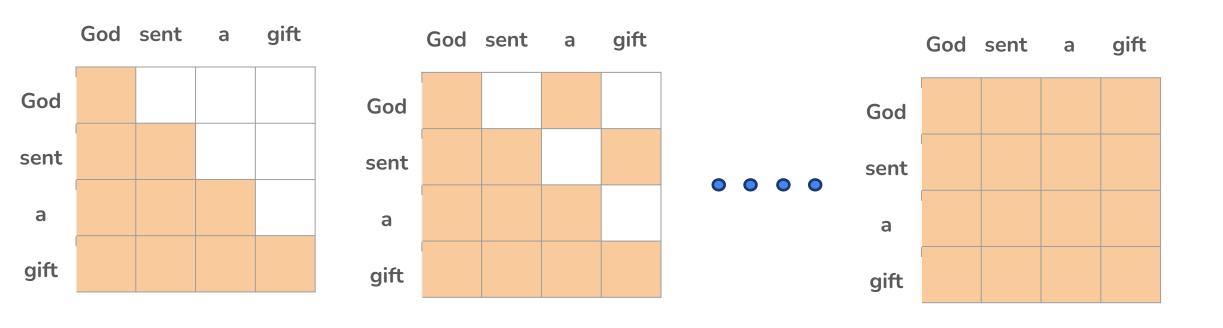
- Previously, we observed that masked diffusion LM loss can be simplified to cross-entropy.
- Cross-entropy loss is also used to train pre-trained LMs.

(Masked Diffusion Loss)
$$\mathcal{L}_T = \int_0^1 \frac{\alpha_t'}{1 - \alpha_t} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} [\delta_{\boldsymbol{x}_t, \boldsymbol{m}} \boldsymbol{x}_0^\top \log f_{\theta}(\boldsymbol{x}_t)] \, dt.$$

(LLM Loss)
$$\mathcal{L}_{AR}^{1:N} = -\sum_{n=1}^{N} (\mathbf{x}_0^n)^ op \log f_{ heta}(\mathbf{x}_0^{1:n-1})_{n-1}$$

- How do we adapt?
 - **Issue no .1:** LLM uses a causal attention mask while diffusion uses bi-drectional attention.
 - **Issue no. 2:** LLMs predicts for the next position while diffusion predicts on the same one.

Causal Attention to Bi-directional Attention



Fine-tuning progress

Training Algorithm

Algorithm 1 Adaptation Training

```
1: Input: network f_{\theta} initialized by existing models,
   training corpus p_{data}(x_0^{1:N}), mask token m.
```

- 2: Output: model parameters θ .
- 3: repeat

4: Draw $x_0^{1:N} \sim p_{data}$ and set labels $\leftarrow x_0^{1:N}$

- 5: Sample $t \in Uniform(0, 1)$ Forward diffusion step \downarrow
 - 6: Sample $x_t^{1:N} \sim q(x_t|x_0)$
 - Anneal the attention mask attn_mask
 - 8: Forward $logits \leftarrow f_{\theta}(\boldsymbol{x}_{t}^{1:N})$ with $attn_mask$
 - Right shift logits by one position

for masked diffusion

- Masked diffusion loss: 10: $\mathcal{L}_t = \frac{1}{t} \delta_{x_t,m} \text{CE}(logits, labels) \triangleright \text{Eq.} 7$
 - Backprop with \mathcal{L}_t and update θ 11:
 - 12: **until** end training

Supplementary Material

- DDPM blog: https://lilianweng.github.io/posts/2021-07-11-diffusion-models/
- Score SDE blog: https://fanpu.io/blog/2023/score-based-diffusion-models/
- Flow matching and Diffusion Models lectures: https://youtube.com/playlist?list=PL57nT7tSGAAUDnli1LhTOoCxlEPGS19vH&si=L8T3fjOFbhooLj9T

Thank You