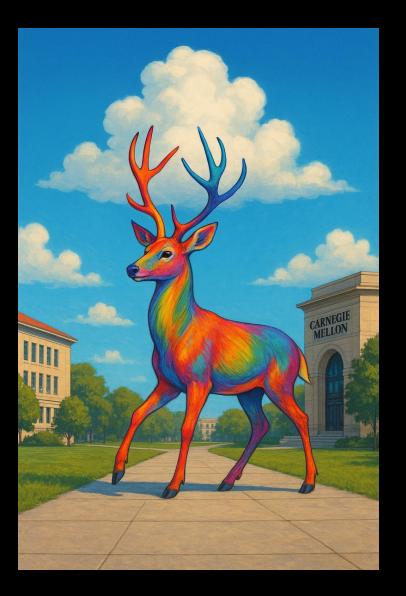
Diffusion Model

Lei Li

Carnegie Mellon University
Language Technologies Institute

Image and Video Generation

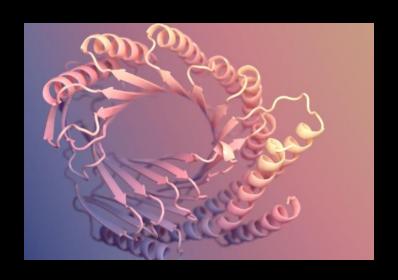
generate an image of a colorful deer riding on the walking into the sky at cmu campus



Diffusion Models are state-of-the-art models for (Continuous) Data Generation

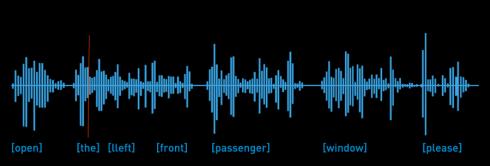
- Time series data
- Audio/speech
- 3D Objects (coordinates)
 - Molecule structures

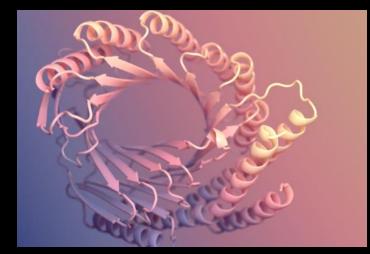




Data Representation

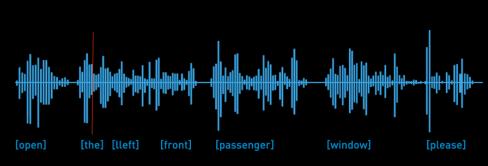


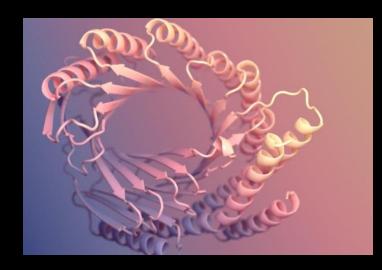




Probabilistic Model for Data





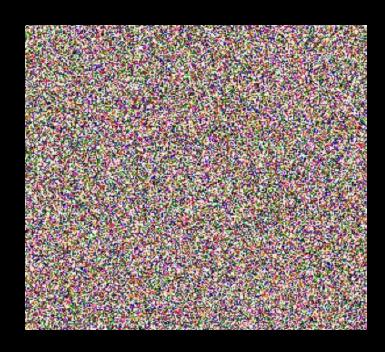


Probabilistic Generative Process

Start from initial distribution $p_T(x_T)$ e.g. Gaussian N(0, I)

Generate slightly improved data $x_{t-1} \sim p_{t-1}(x_{t-1}|x_t)$

final data $x_0 \sim p_0(x_0|x_1)$







Learning the generative model = learning the parameters for each $p_{t-1}(x_{t-1}|x_t)$

It is difficult to directly construct a series of probability distributions

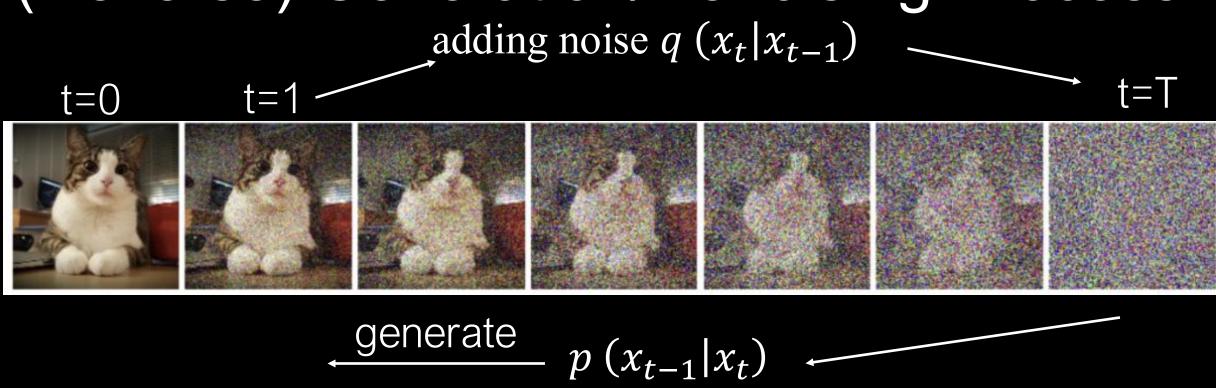
(Forward) Noising Process

• adding Gaussian noise at each time step t $x_t \sim q(x_t|x_{t-1}) = N(\sqrt{\alpha_t}x_{t-1},\beta_t I)$, $\alpha_t = 1 - \beta_t$

equivalently
$$\mathbf{x}_t = \sqrt{\alpha_t} x_{t-1} + \sqrt{\beta_t} \epsilon, \epsilon \sim N(0, I)$$
 $t=0$ $t=1$ $t=T$

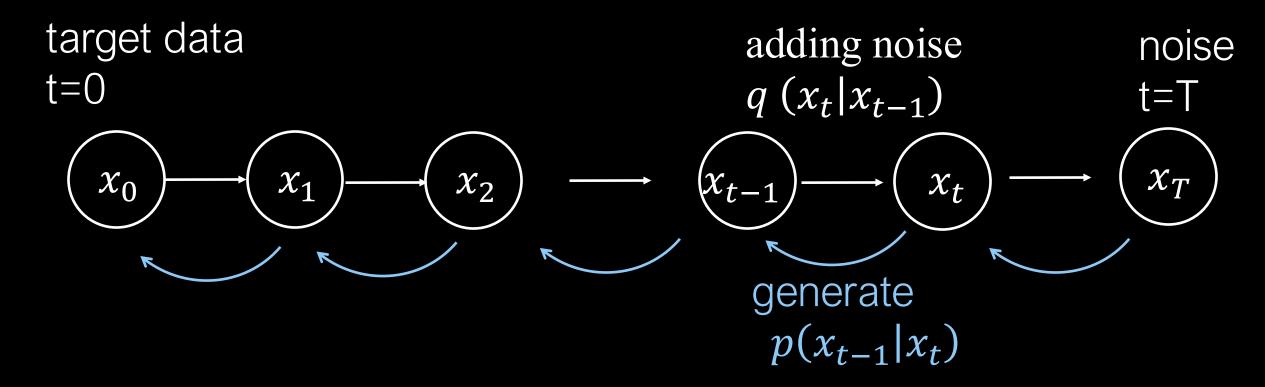
Consequence: with sufficient large T, the result will be Gaussian noise

(Reverse) Generation/Denoising Process



Learning problem: how to find parameters of p to match the sequence of noised images?

Diffusion Model



Learning problem: how to find parameters of p to match the sequence of noised images?

Denoising Diffusion Probabilistic Models (DDPM)

Reverse Process (Generation), aka denoising

$$p_{\theta} = \prod_{t=1}^{T} p_{\theta}(x_{t-1}|x_t)$$

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

Forward Process (Diffusion), aka noising

$$q(x_{1:T}|x_0) = \prod_{t=1}^{T} q(x_t|x_{t-1})$$

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

noising schedule: β_t in [0, 1], and follows a predefined schedule

DDPM Training

- For each data sample x_0 (e.g. image), random pick step t
 - o sample a noise from standard Gaussian N(0,I) => ϵ
 - o noise image at step t: $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 \bar{\alpha}_t} \epsilon$, $\bar{\alpha}_t = \prod_{s=1}^t (1 \beta_s)$
 - o estimate noise using Neural network: $\hat{\epsilon}_{\theta} = NN_{forward}(x_t, t)$
 - o compute MSE loss: $L = \|\epsilon \hat{\epsilon}_{\theta}\|^2$
 - o compute gradient via backward through NN
 - o update model parameters using optimization alg (e.g. Adam)

DDPM training is essentially predicting the noise based on image

DDPM Inference

- randomly sample from Gaussian(0, I) $\rightarrow x_T$
- For step t = T to 1
 - o estimate noise using Neural network: $\hat{\epsilon}_{\theta} = NN_{forward}(x_t, t)$
 - o generate noise from N(0, I) $\rightarrow \eta$
 - o estimate data mean: $\tilde{\mu}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t \frac{\beta_t}{\sqrt{1-\overline{\alpha}_t}} \hat{\epsilon}_{\theta} \right)$
 - o update data: $x_{t-1} = \tilde{\mu}_{t-1} + \sqrt{\beta_t} \eta$

DDPM inference is essentially iteratively removing predicted noise

Intuition and Theory

Training objective:

$$L = -\log p_{\theta}(x_0)$$

$$\leq -E \left[\log \frac{p_{\theta}(x_{0:T})}{q(x_{1:T}|x_0)} \right] = L_{elbo}$$
 elbo loss

$$\begin{split} & L_{elbo} \\ &= \sum_{t=1}^{T} \text{KL} \big(q(x_{t-1} | x_t, x_0) \parallel p_{\theta}(x_{t-1} | x_t) \big) + \text{KL} \big(q(x_T | x_0) \parallel p_{\theta}(x_T) \big) \\ & - \text{E} \log p_{\theta}(x_0 | x_1) \end{split}$$

Hint: Jensen's inequality, posterior for linear Gaussian distribution

$$q(x_t|x_0) = N(x_t; \sqrt{\overline{\alpha}_t}x_0, (1 - \overline{\alpha}_t)I)$$

$$\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$$

$$q(x_{t-1}|x_t, x_0) = N(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}_t I)$$

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}$$

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

just set

$$p_{\theta}(x_{t-1}|x_t) = N(\mu_{\theta}(x_t, t), \sigma_t^2 I)$$

$$KL(q(x_{t-1}|x_t,x_0) || p_{\theta}(x_{t-1}|x_t))$$

becomes

$$KL\left(N(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}_t I) \parallel N(\mu_{\theta}(x_t, t), \sigma_t^2 I)\right)$$

let $\sigma_t^2 = \tilde{\beta}_t$, and let

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \hat{\epsilon}_{\theta}(x_t, t) \right)$$
$$\min \|\tilde{\mu}_t - \mu_{\theta}(x_t, t)\|^2$$

is equivalent to

$$\min \frac{\beta_t^2}{\alpha_t (1 - \bar{\alpha}_t)} \|\epsilon - \hat{\epsilon}_{\theta}(x_t, t)\|^2$$

Essentially, DDPM is using NN to predict noise given noised data

Code Example (see colab)

Summary

Diffusion Model

- o is a probabilistic generative model by iteratively removing noise
- o construct a sequence of corrupted data by applying Gaussian noises
- use a neural network to estimate the noise given the corrupted data (without knowing the original clean data)
- o once estimate the noise, use posterior estimation to remove noise
- o can use any suitable NN (e.g. UNet)
- o much better than VAE, but much slower