

DDIM 2021

Monday, July 7, 2025 4:37 PM

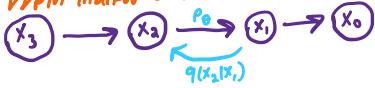
Denoising Diffusion Implicit Models

- More efficient/fast compared to DDPM (10 to 50x faster)
- Generates superior samples

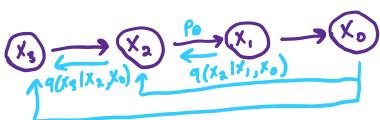
Introduction

- Introduces non-Markovian inference models

DDPM Markov Chain



DDIM non-Markovian inference Model



DDPM Goal: learn model distribution $p_0(x_0)$ that approximates $q(x_0)$ and easy to sample from

Variational Inference for Non-Markovian Forward Processes

- Need to find faster inference process to reduce # of iterations in forward process
- observation: The DDPM objective, L_θ , only depends on $q(x_t|x_0)$, but not $q(x_{1:T}|x_0)$

Forward Process

- Consider family q of inference distributions, indexed by the vector $\sigma \in \mathbb{R}_{\geq 0}^T$

$$\text{Formula 1: } q_\sigma(x_{1:T}|x_0) := q_\sigma(x_T|x_0) \prod_{t=2}^T q_\sigma(x_{t-1}|x_t, x_0)$$

distribution of the final noisy sample, x_T
distribution of the previous-step latent, x_{t-1} , given current noisy sample, x_t , and the original clean image x_0 , under a noise schedule, σ

Meaning: The full forward noising process, expressed as $q_\sigma(x_{1:T}|x_0)$ (where σ denotes we parameterize by the noise), is defined as the joint distribution of all the noisy variables x_1, x_2, \dots, x_T given the original data x_0 .

Purpose: Makes the reverse process deterministic (no neural network like with the DDPM), so the DDIM just learns to predict the noise.

$$\text{Formula 2: } q_\sigma(x_t|x_{t-1}, x_0) := N\left(\hat{x}_t, \sigma_t^2 I\right)$$

distribution of less noisy image given it's more noisy image, and the fully denoised image
 $\hat{x}_t = \text{original image should remain at step } t-1$
 $\sigma_t = \text{the SD of the noise added at step } t$
 $\sigma_t^2 = \text{scale factor added to } \hat{x}_t \text{ to get } x_t$
 $\sigma_t^2 = \text{standard variance (in all directions)}$

the mean is the sum of a clean signal (from x_0) and a residual noise signal (from x_0), carefully scaled so the resulting sample x_{t-1} stays consistent with the diffusion schedule.

Meaning: Given how noisy the image is now, and what the original image looked like, how should I compute the next denoised image?

Purpose: Defines how to take one step backwards in the diffusion process.

$$\text{Bayes' Theorem: } P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

Formula 3: $q_\theta(x_t|x_{t-1}, x_0) = \frac{\text{the "reverse" step (formula 2)}}{\text{overall distribution of noisy image } x_t \text{ given } x_0}$

(derived from Bayes')

$q_\theta(x_{t-1}|x_t, x_0) q_\theta(x_t|x_0)$

forward transition probability from $x_{t-1} \rightarrow x_t$ given the original image

Meaning: the forward transition probability from $x_{t-1} \rightarrow x_t$ given the clean image

Purpose: defines a forward step using reverse-style steps, so the model can be reversed deterministically.

Reverse Process

- Generative process: $p_\theta(x_0:T)$
- each $p_\theta^{(t)}(x_{t-1}|x_t)$ leverages knowledge of $q_\theta(x_{t-1}|x_t, x_0)$ (defined already)

Formula 4: $f_\theta^{(t)}(x_t) := \frac{x_t - \sqrt{1-\alpha_t} \cdot e_\theta^{(t)}(x_t)}{\sqrt{\alpha_t}}$

noisy image

predicted noise

estimate of a clean image at time step

Meaning: The model's estimate of a clean image is the noisy image minus the predicted noise (scaled).

Purpose: tells us how to predict the clean image estimate of clean image and standard unit variance

Formula 5: $p_\theta^{(t)}(x_{t-1}|x_t) := \begin{cases} N(f_\theta^{(t)}(x_t), \sigma_t^2 I) & \text{if } t \neq 1 \\ q_\theta(x_{t-1}|x_t, f_\theta^{(t)}(x_t)) & \text{otherwise} \end{cases}$

Forward transition probability given the predicted clean image (similar to Formula 3)

Meaning: The predicted next image given the current more noisy image is either:

- Add a small Gaussian centered at the prediction of the clean image if we're at the end ($t=1$)
- Use the clean image estimate to deterministically guide the reverse process in all other cases

Purpose: tells us how to sample next less noisy image

Training Objective

- Training Loss function: $J_\theta(\epsilon_\theta)$ distribution over noisy images x_1, x_2, \dots, x_T generated from clean image x_0 in reverse, from pure noise (x_T) to a clean image (x_0)

Formula 6: $J_\theta(\epsilon_\theta) := \mathbb{E}_{x_0:T} q_\theta(x_0:T) [\log q_\theta(x_{1:T}|x_0) - \log p_\theta(x_{0:T})]$

... function

$\log(a) - \log(b) = \log\left(\frac{a}{b}\right) \Rightarrow$ comparing two stories
 $a: \text{starting from a real image } x_0, \text{ if I add noise in a known way, I know how images should look at each step}$
 $b: \text{... from random noise } x_T, \text{ I use my neural network to try to ... I might be wrong or uncertain.}$

↓
the expectation/weighted
average over the full
trajectory (x_0, x_1, \dots, x_T)
that were sampled from the
known forward process q_0

↑
starting from x_0 ,
reverse-engineer the path back to x_0 - but in
 $\log \left(\frac{\text{true/noising process}}{\text{model's reverse process}} \right)$

$= \mathbb{E}_{x_0:T \sim q_0(x_0:T)} \left[\log(q_0(x_T|x_0)) + \sum_{t=1}^T \log q_0(x_{t-1}|x_t, x_0) \right]$

sum of known conditional probabilities
describing how each prior noisy image x_{t-1}
looks given the next one, x_t , and the
original, x_0

the prior/initial belief of the final
noised sample, x_T . (usually just
a standard normal: $p_0(x_T) = N(x_T; 0, I)$)

$- \sum_{t=1}^T \log p_0^{(t)}(x_{t-1}|x_t) - \log p_0(x_T)$

models predicted probabilities for
reversing noise - from x_t to x_{t-1}
at every step.

Meaning: The objective function quantifies the expected difference between the true forward noising process and the model's learned reverse denoising process, using KL-divergence

Purpose: we take the derivative of the loss to update parameters (using autograd)

Sampling

L_1 loss = mean absolute error

L_2 loss = mean squared error

Denoising Diffusion Implicit Models

We generate a sample $x_{t-1} \sim p_0(x_{t-1})$ as follows:

$$\text{Formula 7: } x_{t-1} := \sqrt{\alpha_t} \left(\frac{x_t - \bar{\alpha}_t \cdot E_0^{(t)}(x_t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot E_0^{(t)}(x_t) + \sigma_t e_t$$

$\alpha_0 = 1$

image sample
predicted clean image
direction pointing to x_t , the more noisy image
random noise
 $e_t = N(0, I)$

Meaning: given the current (more noisy) image, we predict/generate the next (slightly denoised) image by adding the predicted clean image, a vector pointing in the direction of the more noisy image, and random noise (to make it more stochastic) - all with scaling factors.

Purpose: tells us how to generate the next denoised sample, x_{t-1} , given x_t during sampling.

Accelerated Generation

- key to the performance speed-up:
- the DDPM takes 100's to 1000's of steps to turn noise into an image
- the DDIM speeds up generation by changing how we sample

- From DDPM, we know $q(x_t|x_0) = N(\sqrt{\alpha_t}x_0, (1-\alpha_t)I)$
 - we then know you can go directly from a clean image to a noisy version in one step
- Predict x_0 from any x_t using the model
 - The trained model $E_0(x_t, t)$ predicts the noise in x_t
 - We rearrange the forward equation: $\hat{x}_0 = \frac{x_t - \sqrt{1-\bar{\alpha}_t} \cdot E_0(x_t, t)}{\sqrt{\bar{\alpha}_t}}$
 - gives a guess of the clean image from any noisy image

- gives --

3) Define deterministic reverse rule
- using predicted \hat{x}_0 , we define a new reverse rule to go from $x_t \rightarrow x_{t-1}$, without sampling

$$x_{t-1} = \sqrt{\hat{\alpha}_{t-1}} \cdot \hat{x}_0 + \sqrt{1 - \hat{\alpha}_{t-1}} \cdot E_\theta(x_t, t)$$

- resembles the forward process - using model's predicted noise

→ fully deterministic

→ works with any schedule of steps

- e.g., $t=1000 \rightarrow 950 \rightarrow 900$

- trade off: \hat{x}_0 isn't perfect so balance results with desired sampling time

Summary

DDPM vs. DDIM

Feature	DDPM (original)	DDIM (improved)
Sampling type	stochastic (adds noise at each step)	Deterministic (no noise added)
Reverse Step Formula	$x_{t-1} \sim N(\mu_\theta, \sigma_t^2)$	$x_{t-1} = \sqrt{\hat{\alpha}_{t-1}} \hat{x}_0 + \sqrt{1 - \hat{\alpha}_{t-1}} \cdot E_\theta$
Steps Required	~ 1000 for quality images	10-50 (speed vs performance)
Randomness During Inference	Yes - samples at each step	No - completely deterministic
Can Skip Steps	No - must go step-by-step	Yes - jumps between custom steps
key trick	Learns to predict noise	Uses predicted noise to recover \hat{x}_0 , then rewinds deterministically using known marginals
why it works	Reverse process modeled as a markov chain of Gaussians	reverse process derived from forward marginals $q(x_t x_0)$ and the model's noise predictions
Noise Prediction Used	Yes - to define mean of each reverse Gaussian	Yes - to reconstruct \hat{x}_0 , which then determines the next image in the reverse process.