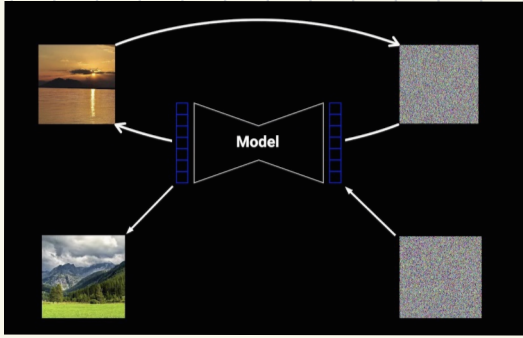Diffusion model: Systematically destroy data distribution through an iterative forward diffusion process, then learn a reverse diffusion process that restores structure in data
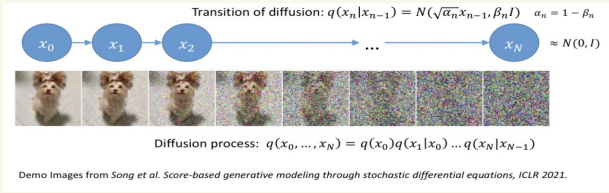


: From YouTube "Outlier"

- Diffusion process gradually injects noise to data.
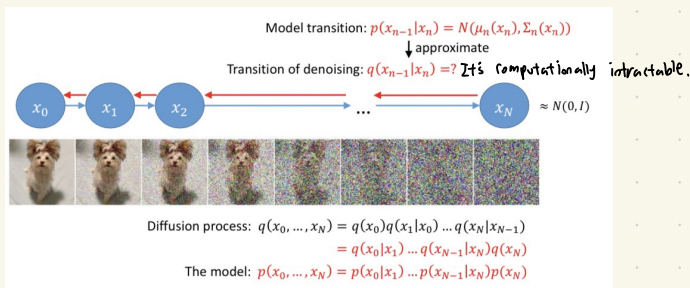  ↳ Described by a Markov Chain: $q(X_0,...,X_N) = q(X_0) q(X_1|X_0) \cdots q(X_N|X_{N-1})$
  
  where $q(X_t|X_{t-1}) = \mathcal{N}(X_t|\sqrt{1-\beta_t}\, X_{t-1}, \beta_t I)$
  
  Then $X_t = \sqrt{1-\beta_t}\, X_{t-1} + \sqrt{\beta_t}\, \varepsilon$, where $\varepsilon \sim \mathcal{N}(0,I)$



Transition of diffusion: $q(x_n|x_{n-1}) = N(\sqrt{\alpha_n}x_{n-1}, \beta_n I)$    $\alpha_n = 1 - \beta_n$

$x_0$ → $x_1$ → $x_2$ → ... → $x_N$ ≈ $N(0, I)$

Diffusion process: $q(x_0, ..., x_N) = q(x_0)q(x_1|x_0)...q(x_N|x_{N-1})$

Demo Images from Song et al. Score-based generative modeling through stochastic differential equations, ICLR 2021.

- Diffusion process in the reverse direction : Denoising process
  ↳ Reverse factorization : $q(X_0,...,X_N) = q(X_0|X_1) \cdots q(X_{N-1}|X_N)\, q(X_N)$



Model transition: $p(x_{n-1}|x_n) = N(\mu_n(x_n), \Sigma_n(x_n))$
↓ approximate
Transition of denoising: $q(x_{n-1}|x_n) = ?$ It's computationally intractable.

$x_0$ ← $x_1$ ← $x_2$ ← ... ← $x_N$ ≈ $N(0, I)$

Diffusion process: $q(x_0, ..., x_N) = q(x_0)q(x_1|x_0)...q(x_N|x_{N-1})$
$= q(x_0|x_1)...q(x_{N-1}|x_N)q(x_N)$
The model: $p(x_0, ..., x_N) = p(x_0|x_1)...p(x_{N-1}|x_N)p(x_N)$

Notation

$X_0$ = Original image, $X_T$ = Isotropic Gaussian
$q(X_t|X_{t-1})$ : Forward Process
$p(X_{t-1}|X_t)$ : Reverse Process
  ↳ $p$ takes in an image $X_t$ and produces a sample $X_{t-1}$ using the neural network.

Forward Process: $q(X_t | X_{t-1})$

$q(X_t | X_{t-1}) = \mathcal{N}(X_t, \sqrt{1-\beta_t} X_{t-1}, \beta_t I)$

    ↑     ↑     ↑

  output  mean  variance, and $\beta$ refers to the Schedule, and $\beta = [0,1]$

Linear Schedule: $\beta_{start} = 0.0001$, $\beta_{end} = 0.02$ Grow linearly up to $0.02$

Forward Process has a closed form Solution.

① let $\alpha_t = 1 - \beta_t$

② $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$

③ By applying "reparameterization trick" ($\mathcal{N}(\mu, \sigma^2) = \mu + \sigma \cdot \varepsilon$), $q(X_t|X_{t-1}) = \mathcal{N}(X_t, \sqrt{1-\beta_t} X_{t-1}, \beta_t I) = \sqrt{1-\beta_t} X_{t-1} + \sqrt{\beta_t} \varepsilon$, where $\varepsilon \sim \mathcal{N}(0,1)$

④ Replace $1-\beta_t = \alpha_t$, then : $\sqrt{\alpha_t} X_{t-1} + \sqrt{1-\alpha_t} \varepsilon$

$\qquad\qquad\qquad = \sqrt{\alpha_t \alpha_{t-1}} X_{t-2} + \sqrt{1-\alpha_t \alpha_{t-1}} \varepsilon$

$\qquad\qquad = \cdots = \sqrt{\alpha_t \alpha_{t-1} \cdots \alpha_1 \alpha_0} X_0 + \sqrt{1-\alpha_t \alpha_{t-1} \cdots \alpha_1 \alpha_0} \varepsilon$

$\qquad\qquad\qquad\qquad = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1-\bar{\alpha}_t} \varepsilon = q(X_t|X_0) = \mathcal{N}(X_t : \sqrt{\bar{\alpha}_t} X_0, (1-\bar{\alpha}_t)I)$

For Sampling : $X_t = \sqrt{\bar{\alpha}_t} X_0 + \sqrt{(1-\bar{\alpha}_t)} \varepsilon$     where $\varepsilon \sim \mathcal{N}(0,1)$
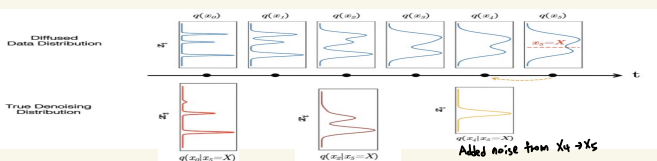
Reverse diffusion Process

- The goal of a diffusion model is to learn the reverse denoising process to iteratively undo the forward process

- Since $q(X_t|X_{t-1})$ is unknown (intractable), we try to approximate $q(X_t|X_{t-1})$ using $P_\theta$ function which is basically what neural network does

Why intractable?

$q(X_t | X_{t-1}) = q(X_t | X_{t-1}) \cdot \dfrac{q(X_{t-1})}{q(X_t)}$, where $q(X_t) = \int q(X_t | X_{t-1}) q(X_{t-1}) dx$

  ↳ Computing this is computationally intractable



Diffused Data Distribution / True Denoising Distribution

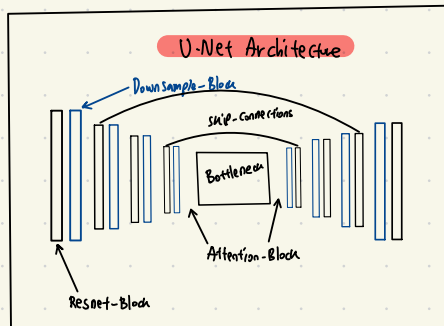- The reverse process step $q(X_{t-1}|X_t)$ can be estimated as a Gaussian distribution too. Hence, we can parameterize the learned reverse process as

$P_\theta (X_{t-1} | X_t) = \mathcal{N}(X_{t-1}, \mu_\theta(X_t, t), \sigma_t^2 I)$

                      ⌣

            Trainable network
            (U-Net, MLP, Denoising Autoencoder)

and, $P(X_T) = \mathcal{N}(X_T, 0, I)$

         ↑  ↑  ↑

        Output  $\mu$  $\sigma$

U-Net Architecture



Down Sample-Block
Skip-connections
Bottleneck
Attention-Block
Resnet-Block

$$: P_\theta(X_0|X_1) \, P_\theta(X_1|X_2) \cdots P_\theta(X_{T-1}|X_T) = \frac{P_\theta(X_0,X_1,\ldots,X_T)}{P_\theta(X_1,\ldots,X_{T-1},X_T)} \cdot \frac{P_\theta(X_1,\ldots,X_{T-1},X_T)}{P_\theta(X_2,\ldots,X_{T-1},X_T)} \cdot \cdots \cdot \frac{P_\theta(X_{T-1},X_T)}{P_\theta(X_T)}$$

$\underbrace{\phantom{P_\theta(X_0|X_1) \, P_\theta(X_1|X_2) \cdots P_\theta(X_{T-1}|X_T)}}$
Markov chain

$$\Rightarrow P_\theta(X_{0:T}) = P(X_T) \cdot \prod_{t=1}^{T} P_\theta(X_{t-1}|X_t)$$

Now, how do we train diffusion model?

: We need to maximize the likelihood that an image we generate looks like it comes from the original data distribution.
We can bound the likelihood using ELBO (Variational Lower bound) just like VAE.

$$L_{VLB} = L_T + L_{T-1} + \cdots + L_0$$

where $L_T = D_{KL}(q(X_T|X_0) \parallel P_\theta(X_T))$

$\qquad\quad L_t = D_{KL}(q(X_t|X_{t+1},X_0) \parallel P_\theta(X_t|X_{t+1}))$ for $1 \le t \le T-1$

$\qquad\quad L_0 = -\log P_\theta(X_0|X_1)$

$$q(X_{t-1}|X_t,X_0) = \mathcal{N}(X_{t-1}, \, \tilde{\mu}_t(X_t,X_0), \, \tilde{\beta}_t I),$$

where $\tilde{\mu}_t(X_t,X_0) := \dfrac{\sqrt{\bar{\alpha}_{t-1}}\,\beta_t}{1-\bar{\alpha}_t}X_0 + \dfrac{\sqrt{\bar{\alpha}_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}X_t$ and $\tilde{\beta}_t = \dfrac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$

And after doing some algebra, we get

$$L_{t-1} = \mathbb{E}_{X_0,\varepsilon}\left[\frac{1}{2\|\Sigma_\theta\|_2^2}\|\tilde{\mu}(X_t,X_0) - \mu_\theta(X_t,t)\|_2^2\right]$$

$$= \mathbb{E}_{X_0,\varepsilon}\left[\frac{1}{2\|\Sigma_\theta\|_2^2}\left\|\frac{1}{\sqrt{\bar{\alpha}_t}}\left(X_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon\right) - \mu_\theta(X_t,t)\right\|_2^2\right]$$

Instead of predicting $\mu$, we should predict epsilon instead.

$$\tilde{\mu}_t(X_t,X_0) = \frac{1}{\sqrt{\bar{\alpha}_t}}\left(X_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon\right) \longrightarrow \mu_\theta(X_t,t) = \frac{1}{\sqrt{\bar{\alpha}_t}}\left(X_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_\theta(X_t,t)\right)$$

Therefore, our loss becomes

$$L_{t-1} = \mathbb{E}_{X_0,\varepsilon}\left[\frac{\beta_t^2}{2\alpha_t(1-\bar{\alpha}_t)\|\Sigma_\theta\|_2^2}\left\|\varepsilon - \varepsilon_\theta\left(\sqrt{\bar{\alpha}_t}X_0 + \sqrt{1-\bar{\alpha}_t}\,\varepsilon, t\right)\right\|_2^2\right]$$

And the authors of DDPM says it's okay to drop all that baggage in front.

$$L_{t-1} = \mathbb{E}_{X_0,\varepsilon}\left[\left\|\varepsilon - \varepsilon_\theta\left(\sqrt{\bar{\alpha}_t}X_0 + \sqrt{1-\bar{\alpha}_t}\,\varepsilon, t\right)\right\|_2^2\right]$$

In Summary: Algorithm 1: Training

1: repeat

2: $X_0 \sim q(X_0)$

3: $t \sim \text{Uniform}(\{1,\ldots,T\})$

4: $\varepsilon \sim \mathcal{N}(0,I)$

5: Take gradient descent step on
$$\nabla_\theta \left\|\varepsilon - \varepsilon_\theta\left(\sqrt{\bar{\alpha}_t}X_0 + \sqrt{1-\bar{\alpha}_t}\,\varepsilon, t\right)\right\|^2$$

6: until converged.

Algorithm 2 : Sampling (Inference)

1: $X_T \sim N(0, I)$

2: for $t = T, \ldots, 1$ do

3:    $z \sim N(0, I)$ if $t > 1$, else $z = 0$

4:    $X_{t-1} = \frac{1}{\sqrt{a_t}} \left( X_t - \frac{1 - a_t}{\sqrt{1 - \bar{a}_t}} \varepsilon_\theta(X_t, t) \right) + \sigma_t z$

5: end for

6: return $X_0$

$\star$ $\sigma_t = \frac{1 - \bar{a}_{t-1}}{1 - \bar{a}_t} \beta_t$

Model will always be designed for each time step and the way of telling which time step we are is done using the sinosidal embeddings