

We start with initial images  $x_0 \sim g(x_0)$ , where  $g$  is a distribution that we want to approximate. Then, we assume Gaussian noise in the forward process in the form  $g(x_t | x_{t-1}) = N(x_t; \sqrt{2\alpha_t} x_{t-1} (1-\alpha_t)^{\frac{1}{2}})$ .

Due to the "nice" properties of Gaussians we can get a closed form solution of  $g(x_t | x_0)$ . We can look inductively:

$$y_t | x_t \sim g(x_t | x_{t-1})$$

$$x_t = \sqrt{\alpha_t} x_{t-1} + \sqrt{1-\alpha_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}} x_{t-2} + \sqrt{1-\alpha_{t-1}} \varepsilon'_{t-2} \right)$$

$$+ \sqrt{1-\alpha_t} \varepsilon_{t-1}$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \overbrace{\sqrt{\alpha_t - \alpha_{t-1}} \varepsilon'_{t-2} + \sqrt{1-\alpha_t} \varepsilon_{t-1}}$$

$$= \overbrace{\sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + \sqrt{\alpha_t - \alpha_{t-1}} \varepsilon'_{t-2}} + \overbrace{\sqrt{1-\alpha_t} \varepsilon_{t-1}}$$

$$= \int 2x_{t+1} X_{t+2} + \int 1 - 2x_{t+1} \epsilon_{t+2}$$

Where I used the fact that the sum of two Gaussian random variables is still a Gaussian random variable with mean and variance as a sum of the means and variances of the two rand. variables.

$$\text{We see that } g(X_t | X_0) = \mathcal{N}(X_t | \sqrt{\sum_{i=1}^t \alpha_i} X_0,$$

$1 - \sum_{i=1}^t \alpha_i \mathbb{I})$ . We will use the notation from the DDPM paper,  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

From there we can calculate the posterior

$$g(X_{t-1} | X_t, X_0) = \frac{g(X_t | X_{t-1}, X_0) g(X_{t-1} | X_0)}{g(X_t | X_0)}$$

$$= \frac{\mathcal{N}(X_t; \sqrt{\alpha_t} X_{t-1}, (1 - \alpha_t) \mathbb{I}) \mathcal{N}(X_{t-1}; \sqrt{\bar{\alpha}_{t-1}}(1 - \bar{\alpha}_t)^{-\frac{1}{2}} X_0)}{\mathcal{N}(X_t; \sqrt{\bar{\alpha}_t}, (1 - \bar{\alpha}_t) \mathbb{I})}$$

which is a lengthy calculation by writing the exponentials. See "Understanding Diffusion

models : A unified perspective " (or a full derivation-

$$\propto \mathcal{N}(x_{t+1}; \frac{\sqrt{\alpha_t(1-\bar{\alpha}_{t-1})}x_t + \sqrt{\bar{\alpha}_{t-1}(1-\alpha_t)}x_0}{1-\bar{\alpha}_t}, \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} I).$$

The ELBO in this case is  $\log p_\theta(x) \geq E_{q(x_{1:T}|x_0)}$

$$\left[ \log \frac{p(x_0:T)}{q(x_{1:T}|x_0)} \right],$$
 where  $p_\theta$  is the reverse process dependent on the neural network's parameters

where  $\log p(x_0) = \log \int p(x_0:T) dx_{1:T}$ .

T being the final step in the diffusion process.

$$so, E_{q(x_{1:T}|x_0)} \left[ \log \frac{p(x_0:T)}{q(x_{1:T}|x_0)} \right]$$

In the rest of the calculation I will stop writing the expectation

$$\Rightarrow \log \frac{p_\theta(x_0:T)}{q(x_{1:T}|x_0)} = \log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{\prod_{t=1}^T q(x_t|x_{t-1})}$$

$$= \log \frac{p_\theta(x_T) p_\theta(x_0|x_1) \prod_{t=2}^T p_\theta(x_{t-1}|x_t)}{g(x_1|x_0) \prod_{t=2}^T g(x_t|x_{t-1})}$$

$$= \log \frac{p_\theta(x_T) p_\theta(x_0|x_1)}{g(x_1|x_0)}$$

$$+ \log \frac{\prod_{t=2}^T p_\theta(x_{t-1}|x_t)}{\prod_{t=2}^T g(x_t|x_{t-1})}$$

$$= \log \frac{p_\theta(x_T) p_\theta(x_0|x_1)}{g(x_1|x_0)} + \log \frac{\prod_{t=2}^T p_\theta(x_{t-1}|x_t)}{g(x_t|x_{t-1})}$$

$$= \log \frac{p_\theta(x_T) p_\theta(x_0|x_1)}{g(x_1|x_0)}$$

$$+ \log \frac{\prod_{t=2}^T \frac{p_\theta(x_{t-1}|x_t)}{g(x_t|x_{t-1}, x_0) g(x_t|x_0)}}{g(x_T|x_{t-1}, x_0)}$$

$$\begin{aligned}
 &= \log \frac{p_\theta(x_T) p_\theta(x_0|x_1)}{q(x_1|x_0)} \\
 &\quad + \log \frac{q(x_1|x_0)}{q(x_T|x_0)} + \log \prod_{t=2}^T \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)} \\
 &= \log \frac{p_\theta(x_T) p_\theta(x_0|x_1)}{q(x_T|x_0)} + \sum_{t=2}^T \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1}, x_0)}
 \end{aligned}$$

$$\begin{aligned}
 &= \left[ E_{g(x_1|x_0)} \log p_\theta(x_0|x_1) \right] - D_{KL}(q(x_T|x_0) || p(x_T)) \\
 &\quad - \sum_{t=2}^T \left[ E_{g(x_t|x_0)} [D_{KL}(q(x_{t-1}|x_t, x_0) || p_\theta(x_{t-1}|x_t))] \right]
 \end{aligned}$$

After plugging back the expectations.  
 All the terms in the above expression  
 will receive a minus sign because of mini-  
 mizing the negative log-likelihood, see  
 the sign in the ELBO changes. The KL-di-

Divergence terms can be calculated exactly since  $g(x_{t-1}|x_t, x_0)$  and  $p_\theta(x_{t-1}|x_t)$  are Gaussians.  $D_{KL}(g(x_T|x_0) \parallel p(x_T))$  is zero because of assuming  $g(x_T|x_0)$  to be a standard Gaussian as well as  $p(x_T)$ . In the calculation above, we also used the Markovian assumption  $p_\theta(x_t|x_{t-1}, x_{t-2}, \dots) = p_\theta(x_t|x_{t-1})$ .

For the KL's,  $D_{KL}(g(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t))$

$$= \frac{1}{2} \left\| \mu_\theta - \frac{\sqrt{\alpha_t(1-\alpha_{t-1})}x_t + \sqrt{\alpha_{t-1}(1-\alpha_t)}x_0}{1-\bar{\alpha}_t} \right\|_2^2$$

$\nearrow$

the variance of  $g(x_{t-1}|x_t, x_0)$

Where  $p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t), \Sigma_g^{-1} I)$

By a reparametrization of  $\mu_\theta$  as

$$\mu_\theta = \frac{\sqrt{\alpha_t(1-\alpha_{t-1})}x_t + \sqrt{\alpha_{t-1}(1-\alpha_t)}\hat{x}_g(x_t)}{1-\bar{\alpha}_t}$$

$$\text{we get : } \frac{1}{2} \frac{1}{\nabla^2 g} \cdot \frac{\bar{x}_{t-1} (1-\bar{x}_t)^2}{(1-\bar{x}_t)^2} \left\| \hat{x}_g(x_t, t) - x_0 \right\|^2$$

Furthermore, to improve training efficiency we can uniformly sample from  $[0, T]$  instead of taking every  $t \in [0, T]$ . Indeed, this provides an approximation for the actual value. There is a deeper question regarding the meaning of using a NN to predict  $x_0$  given the noisy image.

Does it imply that when we sample from  $\mathcal{N}(0, I)$ , the resulting image generated is the prediction of the NN? Well, in theory yes. If the network is trained so that it captures the data characteristics, then it will produce a high quality image. In almost all scenarios, the dataset is very complex, and one might give a high quality prediction

do the NN might give a noisy result.  
For less noisy steps, but for steps closer to  
T the result might still be very noisy.  
The key word here is sample quality. You  
might want to think that we want to minimize  
ELBO and model this reverse process,  
and to do that we need the prediction at  
step  $t$  be close to the initial image.

Afterwards, we use the reverse process  
to sample. So, we have this construction  
in order to model the data, otherwise try  
using the NN above and compare it to see  
if there are major differences (there are).

Even so, it was tested and predicting the  
noise instead of  $\hat{x}_\theta(x_t, t)$  achieves a better  
sampling quality.

$$\text{Let's start again: } x_t = \sqrt{2}_t x_0 + \sqrt{1-\alpha_t} \epsilon$$
$$\Rightarrow x_0 = \frac{x_t - \sqrt{1-\alpha_t} \epsilon}{\sqrt{2}_t}$$

$$\Rightarrow \mu_g(x_t, x_0) = \frac{\sqrt{2t}(1-\bar{2}_{t-1})x_t + \sqrt{2_{t-1}}(1-2_t)x_0}{1-\bar{2}_t}$$

After substituting  $x_0$  in  $\mu_g$  we get

$$\mu_g = \frac{1}{\sqrt{2t}}x_t - \frac{1-2_t}{\sqrt{1-2_t}\sqrt{2t}}\epsilon$$

Thus, we can reparametrize  $\mu_0(x_t, t) = \frac{1}{\sqrt{2t}}x_t - \frac{1-2_t}{\sqrt{1-2_t}\sqrt{2t}}\hat{\epsilon}_\theta(x_t, t)$ , predicting the noise.

Empirically there is an improved sampling quality. The question about the fact we know  $x_0$  and  $\epsilon$  during training, then why we use neural networks for them is a very good one. This is something usually omitted, but of crucial significance. During training, even though we know  $x_0$  and  $\epsilon$  because we use them to make the training dataset, the NN doesn't know about them (and that's what it has to learn).

why is learning). The NN at step  $t$  takes  
an input  $(x_t, t)$ , where  $x_t = x_0 + \text{noise}_t$ .  
Therefore, it needs to decouple the original  
image and the noise from  $x_t$ .

Furthermore, the term  $L_0$  can be ignored  
because the reconstructed image will be  
very close to the original at the first  
step where the added noise is substantially  
smaller. Also, the gradient from  $L_0$  is very  
small after the model trained for big  
steps.