

# Nichol and Dhariwal on DDPM (2021)

Sunday, November 6, 2022 2:53 PM

- This note is about two papers published in 2021 by [Alex Nichol](#) and [Prafulla Dhariwal](#).
  - [Improved Denosing Diffusion Probabilistic Models](#) (ICML 2021)
  - [Diffusion Models Beat GANs on Image Synthesis](#) (NeurIPS 2021)

## Improved DDPM Paper

- DDPM
  - Reintroduced and popularized by Ho et al. in 2020 [\[LINK\]](#).
  - Shown to produce high quality images.
    - Measured by the FID score.
    - Competitive to GANs and autoregressive models.
  - However, does not do so well with the log-likelihood metrics.
    - Not competitive with autoregressive models yet.

### - Log-likelihood metrics

\* We have a data distribution  $p_{\text{data}}$  that generates the data,

\* We also have a generative model that, given a data item  $x$ , can compute the probability  $p_{\theta}(x)$  of generating  $x$

\* The average log-likelihood metric is simply

$$\mathbb{E}_{x \sim p_{\text{data}}} [-\log p_{\theta}(x)]$$

\* Note that it's actually **negative** log-likelihood

\* Commonly abbreviated as NLL

\* The common unit is "bits per dimension".

- Seems like people like to take logarithm with base 2.

- Also, the unit implies we must divide the log-likelihood with the number of dimensions of the data item.

\* For a DDPM, the log-likelihood cannot be computed directly.

\* However, we know that the log-likelihood is greater than its variational lower bound derived from Jensen's inequality.

$$\mathbb{E}_{x^{0:T} \sim q} [-\log p_{\theta}(x^{0:T})]$$

$$\leq \mathbb{E}_{x^{0:T} \sim q} [\log q(x^{1:T} | x^{(0)}) - \log p_{\theta}(x^{0:T})]$$

} ①

Here,  $q$  is the probability of the forward process

↳ so  $q(x^{(0)}) = p_{\text{data}}(x^{(0)})$

$p_{\theta}$  is the probability of the backward process

$T$  is the number of time steps.

See the [PDF write-up on DDPM](#) for other definitions and notations

\* Now, ① does not really tell us how to compute it.

We need another rewrite:

$$\begin{aligned} E_{x^{(0:T)} \sim q} & \left[ \log q(x^{(1:T)} | x^{(0)}) - \log p_{\theta}(x^{(0:T)}) \right] \\ &= E_{x^{(0:T)} \sim q} \left[ D_{\text{KL}}(q(x^{(T)} | x^{(0)}) \parallel p_{\theta}(x^{(T)})) \right. \\ & \quad - \log p_{\theta}(x^{(0)} | x^{(1)}) \\ & \quad \left. + \sum_{t=2}^T D_{\text{KL}}(q(x^{(t-1)} | x^{(0)}, x^{(t)}) \parallel p_{\theta}(x^{(t-1)} | x^{(t)})) \right] \end{aligned}$$

② {

We know how to compute these.

\* So, to evaluate the variational lower bound of data sample  $x^{(0)}$

① Run the forward process to get  $x^{(1)}, x^{(2)}, \dots, x^{(T)}$

② Use the data items in the forward process to compute the above 3-term expression.

- The paper proposes ways to improve log-likelihood of DDPMs

① Use learned noise schedules rather than fixed ones in [Ho et al. 2020].

② Propose new algorithm to sample time step to optimize to reduce gradient noise.

③ Propose a new noise schedule that does not destroy info too quickly.

- Background

\* The variational lower bound ② is denoted by  $L_{\text{vlb}}$  and the terms are numbered.

$$L_{\text{vlb}} := L_0 + L_1 + L_2 + \dots + L_T$$

$$L_0 := -\log p_{\theta}(x^{(0)} | x^{(1)})$$

$$L_{t-1} := D_{\text{KL}}(q(x^{(t-1)} | x^{(t)}, x^{(0)}) || p(x^{(t-1)} | x^{(t)}))$$

$$L_T := D_{\text{KL}}(q(x^{(T)} | x^{(0)}) || p_{\theta}(x^{(T)}))$$

\* Evaluating  $L_0$  needs care

- Pixel values are discrete, not continuous.
- The paper divide the real line into 256 bins
- It then compute the probability of  $x^{(0)}$  landing into the right bin given  $x^{(1)}$ .
- $p(x^{(0)} | x^{(1)})$  is a Gaussian, so  $\int_a^b p(x^{(0)} | x^{(1)}) dx^{(0)}$  can be computed with Gaussian CDF.

\* For  $L_T$

- It does not depend on  $\theta$ , but we need to compute it for NLL computation.
- Will be close to 0 if the forward process sufficiently destroys the information in the input.

\* Recall that the Hø et. al. paper uses the backward process of the following form:

$$p(x^{(t-1)} | x^{(t)}) = N(x^{(t-1)}; \mu_{\theta}(x^{(t)}, t), \Sigma_{\theta}(x^{(t)}, t)).$$

Here,  $\mu_{\theta}(x^{(t)}, t)$  is a neural network.

Actually,  $\Sigma_{\theta}(x^{(t)}, t)$  can be a neural network too, but Hø et. al. use predetermined constants.

\* Hø et al. experiment with two choices of  $\Sigma_{\theta}(x^{(t)}, t)$ .

$$\textcircled{1} \quad \Sigma_{\theta}(x^{(t)}, t) = \beta_t I \quad (\text{the noise schedule of the forward process})$$

$$\textcircled{2} \quad \Sigma_{\theta}(x^{(t)}, t) = \tilde{\beta}_t I = \left( \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \right) I$$

Here,  $\tilde{\beta}_t$  is the variance that show up in  $q(x^{(t-1)} | x^{(t)}, x^{(0)})$ .

\* Ho et al. mention that:

①  $\beta_+$  is optimal when  $x^{(0)} \sim N(0, I)$

②  $\tilde{\beta}_+$  is optimal when  $x^{(0)}$  is a single sample.

They cite the 2015 paper by Sohl-Dickstein et al. as the source of the two statements above [\[LINK\]](#). However, I could not find the proofs in the paper.

\* They found that either choice did not yield difference in FID score. So, they just set  $\Sigma_\theta(x^{(t)}, t)$  to  $\beta_+$  for simplicity.

\* Instead of predicting the mean of  $p(x^{(t-1)} | x^{(t)})$  directly, the Ho et al. paper trains a network  $\xi_\theta(x^{(t)}, t)$  such that it would predict the noise  $\xi$  that is used to turn  $x^{(0)}$  to  $x^{(t)}$  according to

$$x^{(t)} = \sqrt{\alpha_t} x^{(0)} + \sqrt{1 - \alpha_t} \xi$$

In this way,

$$\mu_\theta(x^{(t)}, t) = \frac{1}{\sqrt{\alpha_t}} \left( x^{(t)} - \frac{\beta_+}{\sqrt{1 - \alpha_t}} \xi_\theta(x^{(t)}, t) \right).$$

\* To train  $\xi_\theta(\cdot, \cdot)$ , the Ho et al. paper does not optimize  $L_{vll}$  directly.

Instead, it optimizes

$$L_{\text{simple}} = E_{t \sim \{0, 1, \dots, T\}, x^{(0)} \sim q, \xi \sim N(0, I)} \left[ \|\xi - \xi_\theta(x^{(t)}, t)\|^2 \right]$$

where  $x^{(t)} = \sqrt{\alpha_t} x^{(0)} + \sqrt{1 - \alpha_t} \xi$

\* Notice that  $L_{\text{simple}}$  does not contain any element that can be used to adjust the variance of  $p(x^{(t-1)} | x^{(t)})$ .

- Improving the log-likelihood ①: Learning  $\Sigma_\theta(x^{(t)}, t)$

\* The setup of Ho et al. achieves good FID, but poor NLL.

\* Nichol and Prafulla found that NLL can be improved if  $T$  is increased from 1000 to 4000.

\* They make the following observations.

-  $\beta_+$  and  $\tilde{\beta}_+$  are extreme values that one can set for the variance  
0 (t-1) (t).

- $\beta_t$  and  $\tilde{\beta}_t$  are extreme values that one can set for the variance of  $p(x^{(t+1)} | x^{(t)})$

↳ I actually don't believe this.

- In the Ho et al paper, the choices between  $\beta_t$  and  $\tilde{\beta}_t$  do not yield differences in FID score because  $\beta_t$  is close to  $\tilde{\beta}_t$  except for when  $t$  is near 0.
- They also observe that the first few terms of  $L_0, L_1, \dots, L_T$  contribute the most to NLL.

- Hence, while the choice of variance of  $p(x^{(t+1)} | x^{(t)})$  might not matter for FID, it might matter for NLL

↳ So, we might be able to improve NLL by optimizing the variance of  $p(x^{(t+1)} | x^{(t)})$  when  $t$  is near 0.

\* Nichol and Dhariwal propose learning  $\Sigma_\theta(x^{(t)}, t)$  instead of using fixed values.

- Their  $\Sigma_\theta(x^{(t)}, t)$  is a diagonal matrix.

Each dimension of the data item has its own variance.

- They parameterize  $\Sigma_\theta(x^{(t)}, t)$  so that it interpolates between  $\beta_t$  and  $\tilde{\beta}_t$  in log space

$$\Sigma_\theta(x^{(t)}, t) = \exp(v \log \beta_t + (1-v) \log \tilde{\beta}_t)$$

- They modify the noise network  $\xi_\theta(x^{(t)}, t)$  so that it has another head that outputs a  $v$  for each dimension.
- They did not constrain  $v$  to be in the range  $[0, 1]$  as they found they did not need to do so in practice.
- To train  $\Sigma_\theta(x^{(t)}, t)$ , they use the loss

$$L_{\text{hybrid}} = L_{\text{simple}} + 0.001 L_v$$

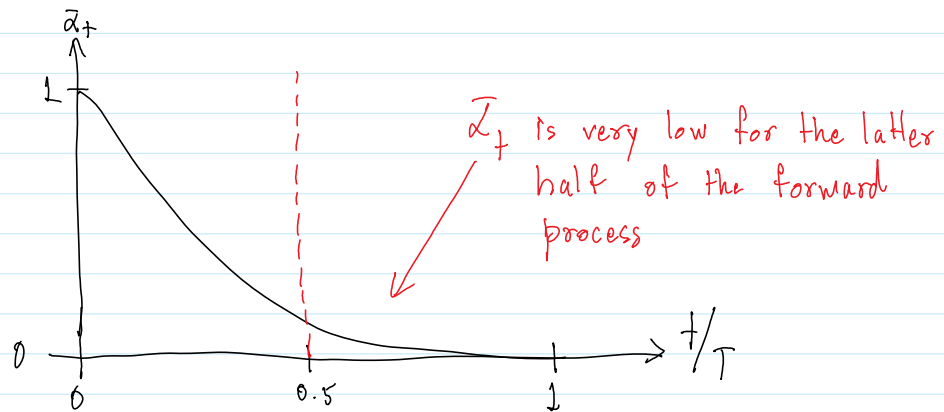
$$L_{\text{hybrid}} = L_{\text{simple}} + 0.001 L_{\text{vlb}}$$

- When evaluating  $L_{\text{vlb}}$ , they freeze the gradient on  $\xi_{\theta}(x^{(t)}, t)$  so that  $L_{\text{vlb}}$  does not affect the noise calculation.

- In other words,  $\xi_{\theta}(x^{(t)}, t)$  is mainly guided by  $L_{\text{simple}}$   
 $\Sigma_{\theta}(x^{(t)}, t)$  is exclusively guided by  $L_{\text{vlb}}$ .

- Improving log-likelihood (2): Better noise schedule

\* If you graph  $\bar{\alpha}_t$  as a function of  $t$ , you will get



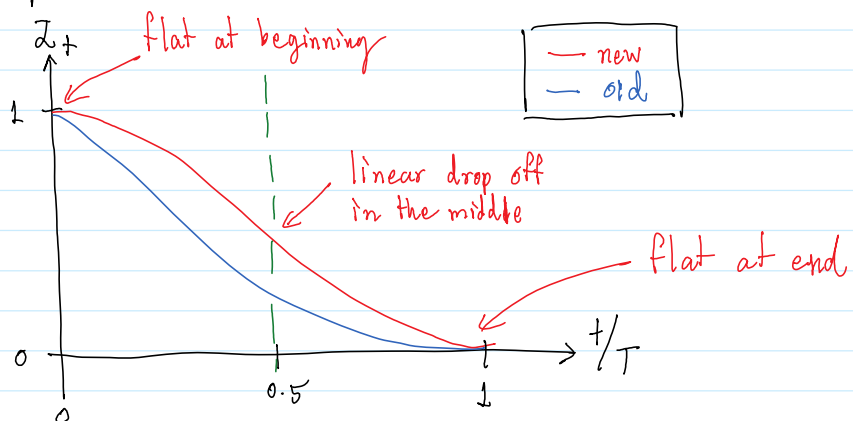
In other words, the noise schedule in the Ho et al. paper seems to destroy information too quickly. As we reach the latter half of the forward process, the distribution already looks so much like  $N(0, I)$ .

\* Nichol and Dhariwal propose using

$$\bar{\alpha}_t = \frac{f(t)}{f(0)} \quad \text{where} \quad f(t) = \cos\left(\frac{t/T + s}{1+s} \cdot \frac{\pi}{2}\right)^2$$

$s = 0.008$

The graph of  $\bar{\alpha}_t$  now looks like this:



\* Rationale for  $\cos^2$

- Does not destroy info too fast.
- Linear drop off near middle  $\rightarrow$  should be where most work is done
- Flat at beginning and end  $\rightarrow$  easier to match distros there
- $\cos^2$  has the right shape, but other functions can also be used.

\* Rationale for the  $s$  constant

- Prevent  $\beta_t$  from being too small near  $t=0$

$\hookrightarrow$  Recall that  $x^{(t)} = \sqrt{\alpha_t} x^{(0)} + \sqrt{1-\alpha_t} \xi$ .

$\hookrightarrow$  So, if  $\beta_t$  is too small near  $t=0$ ,  $\sqrt{1-\alpha_t} \xi$  will be small too.

$\hookrightarrow$  As a result, it becomes hard to predict  $\xi$

- $s$  is chosen so that  $\sqrt{\beta_0}$  is slightly smaller than the bin size  $\frac{1}{127.5}$ .

seems like pixel is normalized to  $[-1, 1]$  range. Pixel value is 0 to 255, and so 255 bins



\* Moreover, the paper also clip  $\beta_t = 1 - \frac{\alpha_t}{\alpha_{t-1}}$  to be no larger than 0.999 so that there would be no singularity near  $t=T$ .

$\hookrightarrow$  Wait, isn't  $\beta_t$  supposed to be low so that each step would be easily reversible?

- Improving log-likelihood ③: Reducing gradient noise

\* During training, if we graph  $L_{\text{simple}}$  and  $L_{\text{vb}}$  against number of iting iterations, the graph would look jaggy.

\* Nichol and Dhariwal hypothesized that sampling  $t$  uniformly causes the loss values to be noisy, especially  $L_{\text{vb}}$ .

\* They propose importance sampling based on  $E[L_t^2]$ .

So,  $L_{\text{vb}} := E_{t \sim p_t} [L_t / p_t]$  where  $p_t \propto \sqrt{E[L_t^2]}$  and  $\sum p_t = 1$ .

\*  $E[L_t^2]$  is estimated by maintaining a history of the last 10 evaluations of  $L_t$ .

$\hookrightarrow$  Initially,  $t$  is sampled uniformly until we have 10 samples for each

$T$ .  
↳ Initially,  $t$  is sampled uniformly until we have 10 samples for each  $L_t$ .

\* After the new sampling strategy is introduced, the graph of  $L_{v/b}$  becomes very smooth.

↳ It's almost too good to be true. Did they cheat?

- Ablation

\*  $\cos^2$  noise schedule improves both FID and NLL.

\* Training with  $L_{\text{simple}} \rightarrow$  good FID, bad NLL

Training with  $L_{v/b} \rightarrow$  bad FID, good NLL

Training with  $L_{\text{hybrid}} \rightarrow$  good FID, good NLL

- Additional benefit: Faster sampling

\* The paper's DDPM is trained with  $T = 4000$ , so it is 4 times slower than that in Ho et al.'s paper.

\* However, Nichol and Dhariwal also experimented with sampling using subsequences of  $(1, 2, 3, \dots, T)$ .

\* Say, the subsequence is  $(s_1, s_2, \dots, s_K)$ . We need to modify the values that  $\Sigma_\theta(\cdot, \cdot)$  interpolate to

$$\beta_{s_k} = 1 - \frac{\bar{\alpha}_{s_k}}{\bar{\alpha}_{s_K}}, \quad \tilde{\beta}_{s_k} = \frac{1 - \bar{\alpha}_{s_k}}{1 - \bar{\alpha}_K} \beta_{s_K}$$

\* The paper experimented with using subsequences of length  $K = 25, 50, 100, 200, 400, 1000$

\* They found that DDPMs trained with  $L_{\text{simple}}$  suffers huge loss in FID scores when  $K$  becomes smaller.

However, for models trained with  $L_{\text{hybrid}}$ , the FID scores did not deteriorate so dramatically.

\* They also compared their results to DDIM [\[LINK\]](#)

- When  $K < 50$ , DDIM achieved better FID scores.

- When  $K \geq 50$ , their DDPM achieved better FID scores.

- Additional result: impact of model size on performance



\* The paper changes model size by changing the number of channels in the first layer.

↳ 64, 96, 128, 192

\* They found that FID and NLL improve as model size increases.

\* FID improvement follow a power law (i.e.  $\log \text{FID}$  is asymptotically linear).

\* NLL improvement does not follow a power law.

---

DDPMs beats GANs paper

- What the paper does.

① Improve DDPM architecture.

② Experiment on classifier guidance.

This allow DDPM to trade diversity for fidelity, like GANs can.

- Result: Their DDPMs achieved better metric values than GANs on multiple datasets.

- What has been used so far.

\* Ho et al. 2020

- UNet architecture

- Global attention layer at  $16 \times 16$  resolution (single head)

- Timestep embedding is projected and added into each residual block.

\* Nichol and Dhariwal 2020 (previous section)

- Learned  $\Sigma_{\theta}(x^{ct}, t)$  and its parameterization

-  $\cos^2$  noise schedule

- Hybrid loss function

- Sampling of  $t$  based on  $E[L_t^2]$  during training.

\* Song et al. 2020 (the DDIM one)

- Deterministic sampling procedure.

- Metrics used to evaluate generative models.

- Metrics used to evaluate generative models.

\* FID is used to assess image quality (diversity + fidelity)

\* IS is used to assess fidelity

\* Improved precision and recall (Kynkääniemi et al. 2019 [LINK](#))

- Precision is used to measure fidelity.

- Recall is used to measure diversity.

- Architecture improvements

\* The changes

- Increasing depth versus width, holding model size relatively constant.

- Increasing number of attention heads

- Using attention at  $32 \times 32$ ,  $16 \times 16$ , and  $8 \times 8$  instead of just  $16 \times 16$

- Using BigGAN residual block for upsampling and downsampling activations.

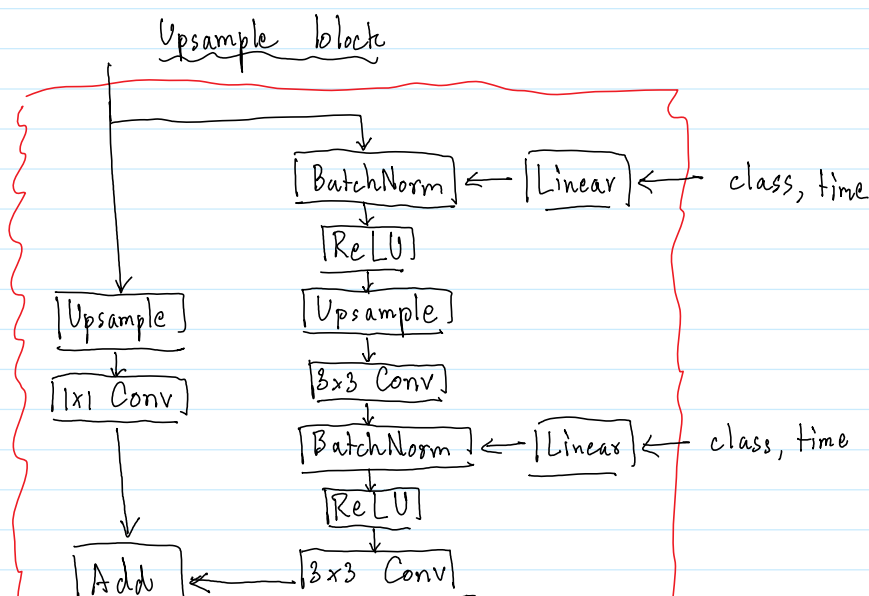
- Rescaling residual connections with  $1/\sqrt{2}$ .

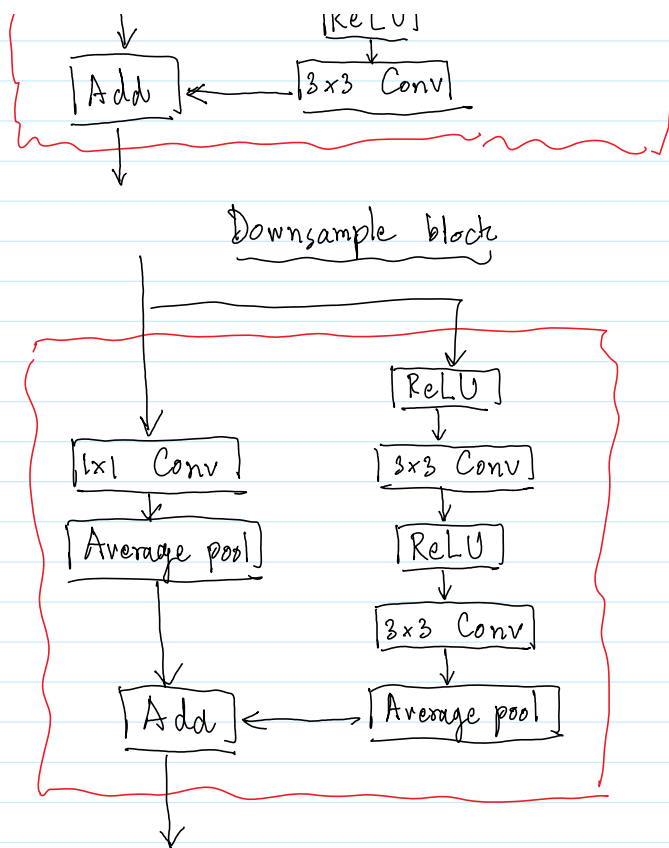
- Adaptive Group Normalization layer (AdaGN)

\* BigGAN residual block

- See the BigGAN paper [LINK](#) for more details.

- Two types of blocks: upsample and downsample blocks





- The BigGAN paper doesn't seem to specify the normal residual block. However, I think it's just removing the Upsample operator from the upsampling block above.

\* Scaling residual connection by  $1/\sqrt{2}$

- Used in StyleGAN and StyleGAN 2 papers
- Residual block typically computes  $x \mapsto x + h(x)$ .
- StyleGANs computes  $x \mapsto \frac{x + h(x)}{\sqrt{2}}$

\* Adaptive Group Normalization (AdaGN)

- BigGAN residual block uses batch norm.
- However, the DDPM in this paper uses AdaGN

$$\text{AdaGN}(x, y_s, y_b) = y_s \text{GroupNorm}(x) + y_b$$

↑ ↑  
linear projections of class and time embedding.

\* Best configuration.

- Base num channels = 128

- 2 residual blocks per resolution
- multi-resolution attention ( $32 \times 32$ ,  $16 \times 16$ ,  $8 \times 8$ )
- 64 channels per attention head (# heads determined by # channels?)
- Classifier guidance
  - \* Used to force a DDPM to generate images from a certain class.
  - \* Train a classifier  $p_\theta(y | x^{(t)}, t)$  on noisy image  $x^{(t)}$ .
  - \* Use gradient  $\nabla \log p_\theta(y | x^{(t)}, t)$  during the sampling process.
    - ↳ DDPM training procedure does not change.
    - ↳ Classifier can be trained separately from DDPM.

#### \* Derivation

- To condition on class  $y$ , it is enough to sample  $x^{(t-1)}$  according to

$$p_{\theta, \theta}(x^{(t-1)} | x^{(t)}, y) = Z p_\theta(x^{(t-1)} | x^{(t)}) p_\theta(y | x^{(t-1)})$$

where  $Z$  is a normalizing constant.

- In general, this is intractable, but it can be approximated.
- First, recall the form of  $p_\theta(x^{(t-1)} | x^{(t)})$ .

$$p_\theta(x^{(t-1)} | x^{(t)}) = N(x^{(t-1)}; \mu, \Sigma)$$

$$\log p_\theta(x^{(t-1)} | x^{(t)}) = -\frac{1}{2} (x^{(t-1)} - \mu)^T \Sigma^{-1} (x^{(t-1)} - \mu) + C$$

- We may assume that  $\log p_\theta(y | x^{(t-1)})$  is not as peaky as

$$\log p_\theta(x^{(t-1)} | x^{(t)})$$

↳ This assumption is reasonable because as  $T \rightarrow \infty$ ,  $\Sigma$  will become smaller and smaller,

- Given the assumption, we can approximate  $\log p_\theta(y | x^{(t-1)})$  with a first order Taylor series expansion around  $\mu$ .

$$\begin{aligned} \log p_\theta(y | x^{(t-1)}) &\approx \log p_\theta(y | \mu) + (x^{(t-1)} - \mu)^T \nabla \log p_\theta(y | x^{(t-1)}) \Big|_{x^{(t-1)} = \mu} \\ &= (x^{(t-1)} - \mu)^T g + C_1 \end{aligned}$$

- This gives

- This gives

$$\log(p_\theta(x^{(t-1)} | x^{(t)}, y) p_\theta(y | x^{(t-1)}))$$

$$\approx -\frac{1}{2} (x^{(t-1)} - \mu)^T \Sigma^{-1} (x^{(t-1)} - \mu) + (x^{(t-1)} - \mu)^T g + C_2$$

Since  $x^{(t-1)}$  is hard to write, let us abbreviate it with just  $x$  for now.

We have that

$$-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) + (x - \mu)^T g + C_2$$

$$= -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) + \frac{1}{2} x^T \Sigma^{-1} \Sigma g - \frac{1}{2} \mu^T \Sigma^{-1} \Sigma g - \frac{1}{2} g^T \Sigma \Sigma^{-1} \mu + \frac{1}{2} g^T \Sigma \Sigma^{-1} x - \frac{1}{2} g^T \Sigma^{-1} g + C_2$$

constant

$$= -\frac{1}{2} (x - \mu + \Sigma g)^T \Sigma^{-1} (x - \mu + \Sigma g) + C_3$$

$$\approx \log N(x; \mu + \Sigma g, \Sigma) + C_4$$

- So, conditional transition can be approximated by shifting the mean

$$\mu_\theta(x^{(t)}, t) \text{ by } \Sigma \nabla \log p_\theta(y | \mu_\theta(x^{(t)}, t))$$

- To make the process more customizable, we will scale  $\Sigma \nabla \log p_\theta(y | \mu)$  by a scaling factor  $s$  before adding to  $\mu$ .
- The conditional sampling algorithm is as follows.

$$x^{(T)} \leftarrow \text{sample from } N(0, I)$$

for  $t \leftarrow T, T-1, T-2, \dots, 1$  do

$$\mu \leftarrow \mu_\theta(x^{(t)}, t)$$

$$\Sigma \leftarrow \Sigma_\theta(x^{(t)}, t)$$

$$x^{(t-1)} \leftarrow \text{sample from } N(\mu + s \Sigma \nabla \log p_\theta(y | \mu), \Sigma)$$

end for

return  $x^{(0)}$

### \* Sampling algorithm for DDIM

- The above conditional algorithm only works for vanilla DDPM. It doesn't work for DDIM.
- In general, for any DDPM, we train a noise prediction network  $\hat{\epsilon}_\theta(x^{(t)}, t)$ . It has the property that

in general, we only have the noise prediction network  $\xi_\theta(x^{(t)}, t)$ . It has the property that

$-\frac{\xi_\theta(x^{(t)}, t)}{\sqrt{1-\alpha_t}}$  is a good approximation of  $\nabla \log p_\theta(x^{(t)})$ .

$$\begin{aligned} - \text{So, } \nabla \log(p_\theta(x^{(t)}) p_\phi(y|x^{(t)})) \\ = \nabla \log p_\theta(x^{(t)}) + \nabla \log p_\phi(y|x^{(t)}) \\ = -\frac{\xi_\theta(x^{(t)}, t)}{\sqrt{1-\alpha_t}} + \nabla \log p_\phi(y|x^{(t)}) \end{aligned}$$

- The new noise prediction  $\hat{\xi}_\theta(x^{(t)}, y, t)$  is given by

$$\begin{aligned} \hat{\xi}_\theta(x^{(t)}, y) &= -\sqrt{1-\alpha_t} \nabla \log(p_\theta(x^{(t)}) p_\phi(y|x^{(t)})) \\ &= -\sqrt{1-\alpha_t} \nabla \log p_\theta(x^{(t)}) - \sqrt{1-\alpha_t} \nabla \log p_\phi(y|x^{(t)}) \\ &= \xi_\theta(x^{(t)}, t) - \sqrt{1-\alpha_t} \nabla \log p_\phi(y|x^{(t)}) \end{aligned}$$

We can then use this new noise prediction to compute  $x^{(t-1)}$  from  $x^{(t)}$ .

### \* Scaling classifier guidance

- Dhariwal and Nichol found that that had to scale the gradient by a factor of more than 1 (that is  $s > 1$ ).

↳ When  $s = 1$ , the classifier assigns reasonable probability (around 50%) to the final sample  $x^{(0)}$ . However, the sample's class does not match the intended class upon visual inspection.

↳ When  $s > 1$ , the class probability increased to nearly 100%.

- We can understand the effect of  $s$  by noting that

$$s \cdot \nabla \log p_\phi(y|x) = \nabla \log \frac{(p_\phi(y|x))^s}{Z}$$

where  $Z$  is the normalization constant.

When  $s > 1$ , the distribution  $(p_\phi(y|x))^s$  becomes peakier by exponentiation. So,  $(p_\phi(y|x))^s$  focuses more on the modes of

the distribution. This leads to higher fidelity in the generated samples.

#### \* Summary of results

- Guidance trades recall for precision.
- Increasing  $s$  improves precision at the cost of recall
- Comparison to GANs and other generative models.
- \* Dhariwal's and Nichol's DDPM generally yielded higher FID and sFID scores than other models, even without guidance.
- \* With guidance, FID and precision improves further.
- \* However, their DDPMs consistently lose to BigGAN-deep in precision and DCTransformers in recall