

# **TTIC 31230, Fundamentals of Deep Learning**

David McAllester, Autumn 2022

Diffusion Model Basics

# Denoising Diffusion Probabilistic Models

Ho, Jain and Abbeel, June 2020



## Progressive VAEs

Diffusion models are a special case of progressive VAEs.

A progressive VAE has layers of latent variables  $z_1, \dots, z_L$ .

The encoder defines  $P_{\text{enc}_0}(z_1|y)$  and  $P_{\text{enc}_\ell}(z_{\ell+1}|z_\ell)$  (defining a Markov chain).

We have a prior  $P_{\text{pri}}(z_L)$  and a decoder defined by  $P_{\text{dec}_\ell}(z_\ell|z_{\ell+1})$  and  $P_{\text{dec}_0}(y|z_1)$ .

Ho et al. take  $L = 1000$ .

## The Ho et al. Diffusion Model

The encoder is not trained — The encoder just adds noise.

The encoder is designed so that  $Z_L$  is distributed as  $\mathcal{N}(0, I)$ .  
The prior is not trained.

The decoder is a single network applied to all layers but taking the layer as an argument.

## A Gaussian Noise Encoder

We assume  $y \in R^d$  and  $z_\ell \in R^d$ .

In the case of images every  $z_\ell$  is an image with the same dimension as  $y$ .

For notational convenience we define  $z_0 = y$ .

The index  $\ell$  will always range from 1 to  $L$  where  $z_{\ell-1}$  might be  $z_0$ .

The encoder is fixed (not trained) and is defined by a sequence of noise levels  $\sigma_1, \dots, \sigma_L$  where for  $1 \leq \ell \leq L$  we have

$$z_\ell = \sqrt{1 - \sigma_\ell^2} z_{\ell-1} + \sigma_\ell \epsilon \quad \epsilon \sim \mathcal{N}(0, I)$$

## A Gaussian Noise Encoder

$$z_\ell = \sqrt{1 - \sigma_\ell^2} z_{\ell-1} + \sigma_\ell \epsilon \quad \epsilon \sim \mathcal{N}(0, I)$$

This is designed so that if  $z_\ell$  has unit variance in each dimension then  $z_{\ell+1}$  also has unit variance in each dimension.

$z_0 = y$  is scaled so that each coordinate is in the interval  $[0, 1]$ .

The constant variance of  $z_\ell$  is important because the same decoder is being used at all  $\ell$  and we want the scale of the decoder input to be independent of  $\ell$ .

## A Gaussian Noise Encoder

$$z_\ell = \sqrt{1 - \sigma_\ell^2} z_{\ell-1} + \sigma_\ell \epsilon \quad \epsilon \sim \mathcal{N}(0, I)$$

Note that the mean of  $z_\ell$  equals  $\sqrt{1 - \sigma_\ell^2}$  times the mean of  $z_{\ell-1}$ .

This repeated reduction drives the mean of  $z_L$  to a negligible level.

We then get that  $z_L$  is distributed as  $\mathcal{N}(0, I)$ .

## A Gaussian Noise Encoder

$$z_\ell = \sqrt{1 - \sigma_\ell^2} z_{\ell-1} + \sigma_\ell \epsilon \quad \epsilon \sim \mathcal{N}(0, I)$$

They increase the noise at higher levels. After some experimentation they use

$$\sigma_\ell^2 = 10^{-4} + .02 \frac{\ell}{L} \quad (\text{with } L = 1000)$$



## Direct Sampling of $z_\ell$

We can sample from  $P_{\text{enc}}(z_\ell|z_0)$  directly

$$\text{define } \alpha_\ell = \prod_{\ell=1}^{\ell} \sqrt{1 - \sigma_\ell^2}$$

$$z_\ell = \alpha_\ell z_0 + \sqrt{1 - \alpha_\ell^2} \epsilon \quad \epsilon \sim \mathcal{N}(0, I)$$

The variance of the noise term can be derived by solving a recurrence relation or by observing that  $z_\ell$  must have unit variance for unit variance  $z_0$ .

## A Natural but Problematic Loss Function

$$\text{dec}^* = \underset{\text{dec}}{\operatorname{argmin}} \quad E_{z_0, \ell, z_{\ell-1}, z_\ell} \quad ||z_{\ell-1} - \text{dec}(z_\ell, \ell)||^2$$

Here  $z_{\ell-1}$  is sampled directly from  $z_0$  and  $z_\ell$  is sampled from  $z_{\ell-1}$ .

## Reducing the Decoder's Dependence on $\ell$ .

We have already reduced the decoder's dependence on  $\ell$  by making  $z_\ell$  have unit variance for all  $\ell$ .

But predicting  $z_{\ell-1}$  from  $z_\ell$  behaves differently for different  $\ell$ .

For small  $\ell$ , where  $\sigma_\ell$  is small,  $z_{\ell-1}$  is near  $z_\ell$ . For large  $\ell$  we have that  $z_{\ell-1}$  is far from  $z_\ell$ . Since the same network is used at all  $\ell$  we want to reduce the dependence on  $\ell$ .

This can be accomplished by using an  $\epsilon$ -decoder.

## An $\epsilon$ -Decoder

First we solve for  $z_{\ell-1}$  in terms of  $z_\ell$  and  $\epsilon$ .

$$z_\ell = \sqrt{1 - \sigma_\ell^2} z_{\ell-1} + \sigma_\ell \epsilon \quad \epsilon \sim \mathcal{N}(0, I)$$

$$z_{\ell-1} = \frac{1}{\sqrt{1 - \sigma_\ell^2}} (z_\ell - \sigma_\ell \epsilon)$$

$$\text{dec}(z_\ell, \ell) = \frac{1}{\sqrt{1 - \sigma_\ell^2}} (z_\ell - \sigma_\ell \epsilon_\Phi(z_\ell, \ell)) + \sigma_\ell \delta \quad \delta \sim \mathcal{N}(0, I)$$

Here  $\epsilon_\Phi(z_\ell, \ell)$  is a trained network whose target value has the same behavior at all levels of  $\ell$ .

## An $\epsilon$ -Decoder

$$\text{dec}(z_\ell, \ell) = \frac{1}{\sqrt{1 - \sigma_\ell^2}} (z_\ell - \sigma_\ell \epsilon_\Phi(z_\ell, \ell))$$

However, SGD on the loss  $\|z_{\ell-1} - \text{dec}(z_\ell, \ell)\|^2$  now scales the SGD gradients on  $\epsilon$  differently for different  $\ell$ .

We effectively have different learning rates for different  $\ell$ .

## Training the $\epsilon$ -Decoder

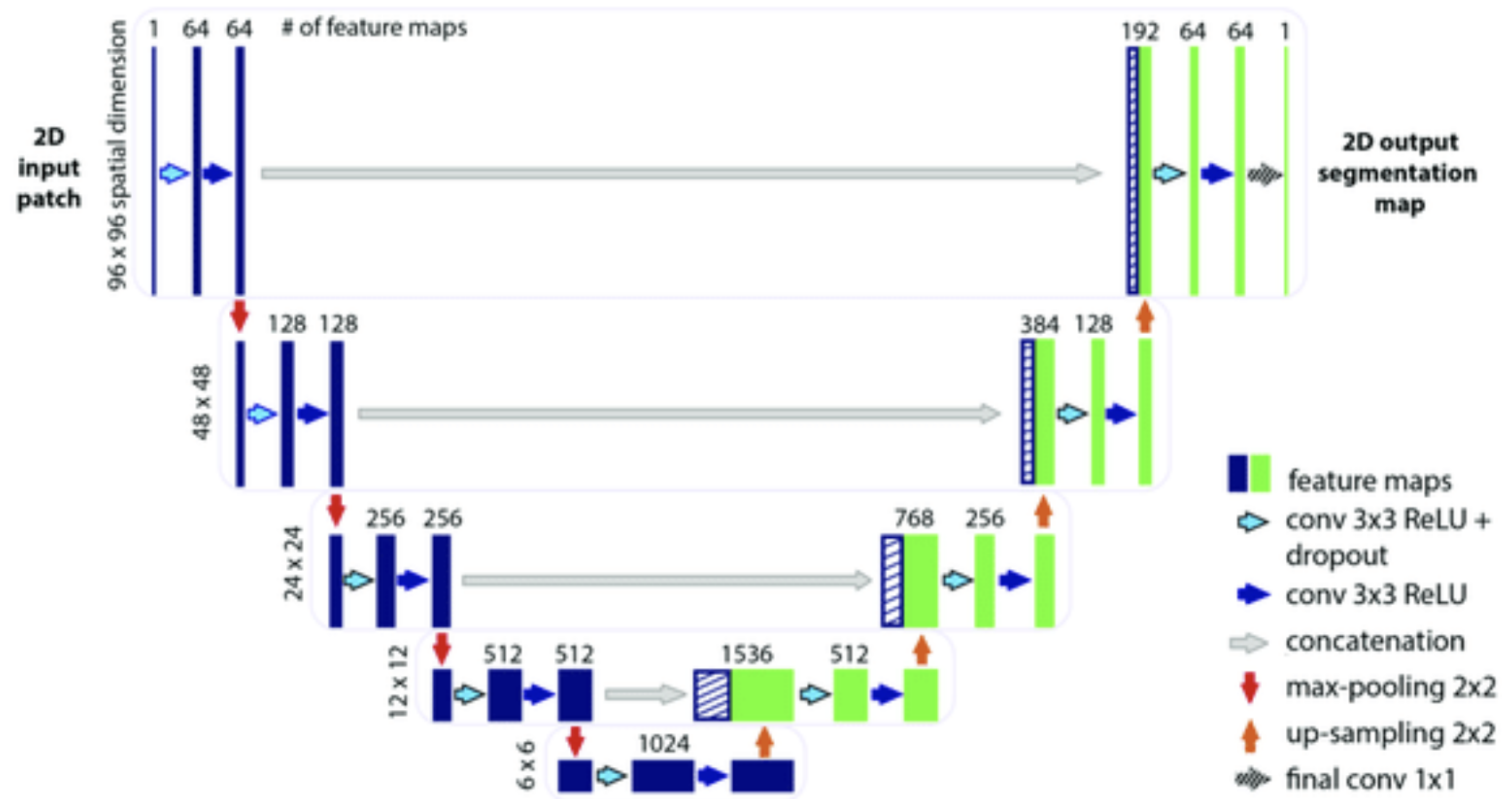
To make the scale of the SGD gradients independent of  $\ell$  we use the following loss.

$$\epsilon^* = \underset{\epsilon}{\operatorname{argmin}} \left\{ \begin{array}{l} E_{z_0, \ell, z_{\ell-1}, \epsilon \sim \mathcal{N}(0, I)} \\ ||\epsilon - \epsilon_{\Phi}(z_{\ell}(z_{\ell-1}, \epsilon), \ell)||^2 \end{array} \right.$$

We now repeatedly sample  $z_0$ ,  $\ell$ ,  $z_{\ell-1}$  and  $\epsilon$  and do gradient updates on  $\epsilon_{\Phi}$ .

# $\epsilon$ -Decoder Architecture

The  $\epsilon$ -decoder is a U-Net.



# Voila





## Progressive VAEs

Diffusion models are a special case of progressive VAEs.

A progressive VAE has layers of latent variables  $z_1, \dots, z_L$ .

It is not clear whether the diffusion model is the best way to do this.

It could be that the important idea is simply having large  $L$  with  $I(z_0, z_\ell)$  going from  $H(z_0)$  to 0 as  $\ell$  goes from 0 to  $L$ .

## Progressive VAEs

A progressive VAE has layers of latent variables  $z_1, \dots, z_L$ .

It seems more natural for image modeling to have spacial resolution decline with increasing  $\ell$ .

This should eliminate the need for U-Nets and reduce computational cost.

**END**