

# DSCI 565: GENERATIVE DEEP MODELING

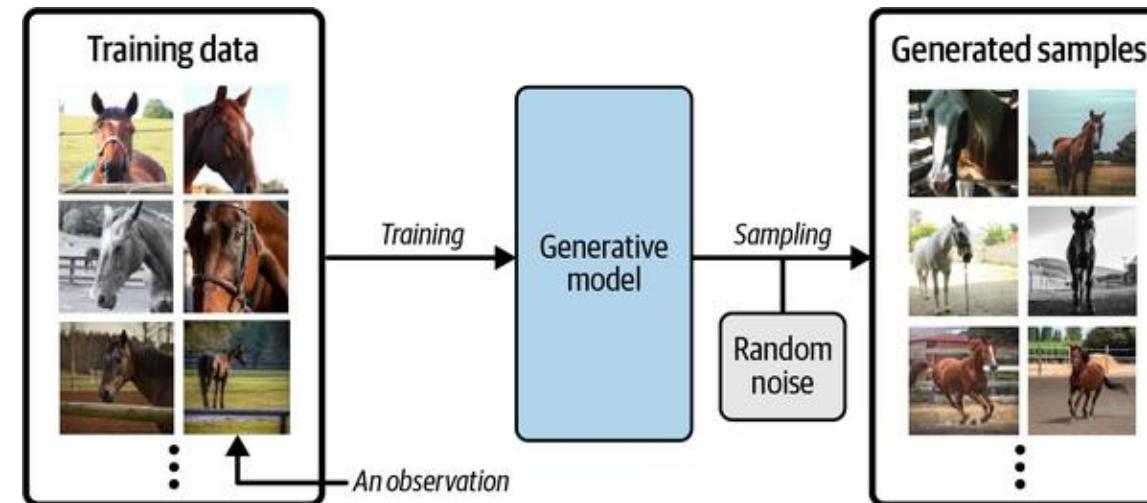
*This content is protected and may not  
be shared, uploaded, or distributed.*

Ke-Thia Yao  
Lecture 3: 2025 November 17

# Generative Modeling

2

- Given a training set, learn a **probabilistic model** that we can sample to generate **synthetic data** that resemble the training set distribution
- For example, given a set of photo-realistic images, learning a model that can generate synthetic photo-realistic images
- Approaches include
  - **Generative Adversarial Networks (GAN)**
  - **Diffusion models**



# References

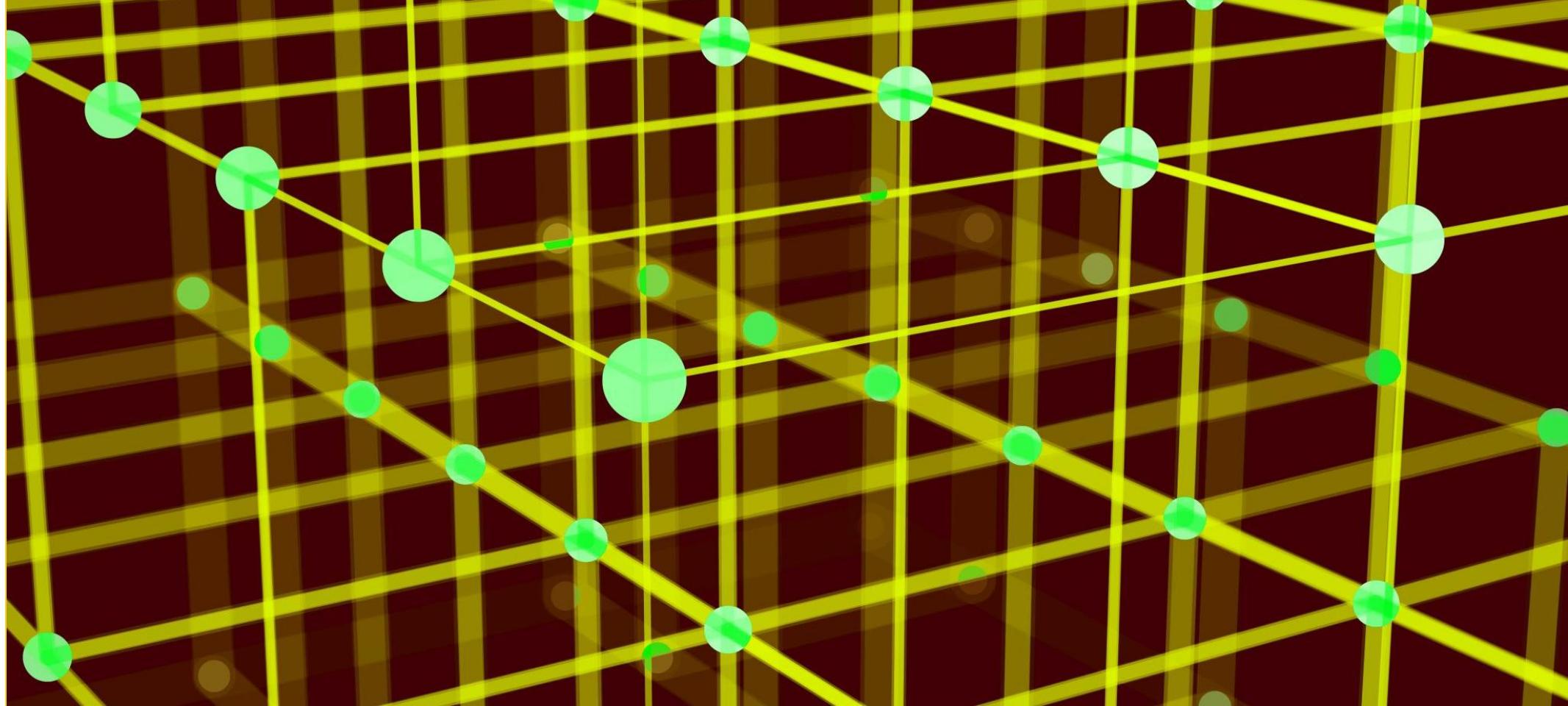
3

## □ References

- Zhang, Aston, Zachary C. Lipton, Mu Li, and Alexander J. Smola. 2023. *Dive into Deep Learning*. 1st edition. Cambridge University Press.
- Generative Deep Learning: Teaching Machines to Paint, Write, Compose, and Play, 2<sup>nd</sup> Edition, David Foster, 2023
- Hands-On Machine Learning with Scikit-Learn and PyTorch, 3<sup>rd</sup> Edition, Aurélien Géron
- Goodfellow, Ian J., Jean Pouget-Abadie, Mehdi Mirza, et al. 2014. “Generative Adversarial Networks.” arXiv:1406.2661. Preprint, arXiv, June 11. <https://doi.org/10.48550/arXiv.1406.2661>.
- Ho, Jonathan, Ajay Jain, and Pieter Abbeel. 2020. “Denoising Diffusion Probabilistic Models.” arXiv:2006.11239. Preprint, arXiv, December 16. <http://arxiv.org/abs/2006.11239>.
- Nichol, Alex, and Prafulla Dhariwal. 2021. “Improved Denoising Diffusion Probabilistic Models.” arXiv.Org, February 18. <https://arxiv.org/abs/2102.09672v1>.
- Saharia, Chitwan, William Chan, Saurabh Saxena, et al. 2022. “Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding.” arXiv:2205.11487. Preprint, arXiv, May 23. <https://doi.org/10.48550/arXiv.2205.11487>.

4

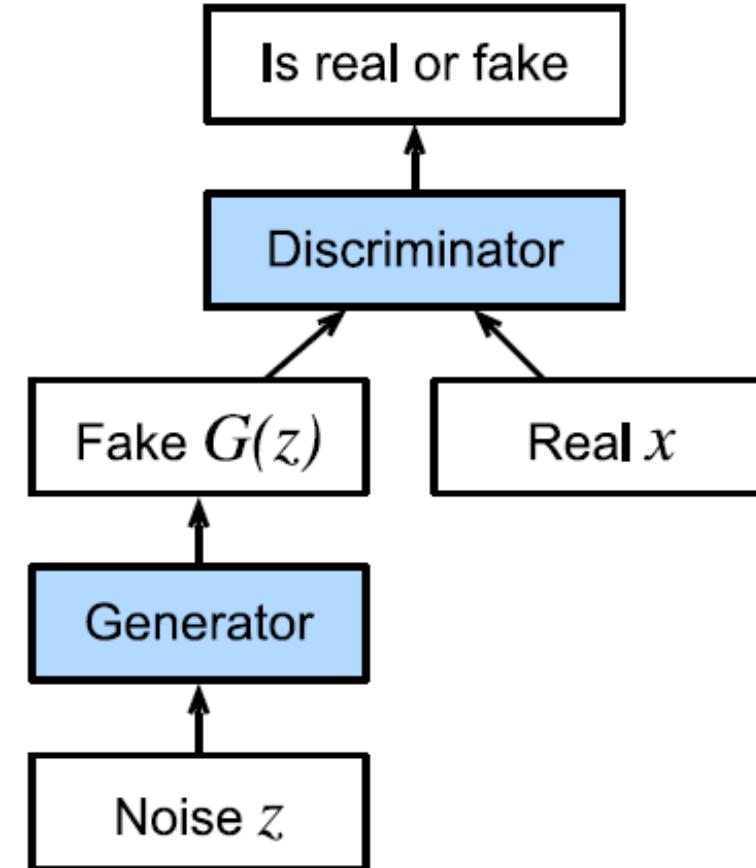
## Generative Adversarial Networks



# GAN Architecture

5

- Consists of
  - Generator network capable of generating data that looks real
  - Discriminator network capable of discriminating between fake and real data
- Iteratively train the two networks to improve them
  - Improve generator network to fool discriminator network
  - Improve discriminator network to better detect fakes



# Discriminator Loss

6

- Discriminator is typically a **binary classifier**
- Training data is labelled  $y = 1$  for true data, and  $y = 0$  for fake
- We train the discriminator  $D(x)$  to minimize **cross entropy loss**:

$$\min_D \{-y \log D(x) - (1 - y) \log(1 - D(x))\}$$

# Generator Loss

7

- Objective is to train the generator  $G(z)$  to fool the discriminator

$$D(G(z)) \approx 1$$

- Where  $z \in \mathcal{R}^d$  is some random vector sampled from  $\mathcal{N}(0,1)$
- Train  $G(z)$  to maximize cross entropy loss when  $y=0$

$$\begin{aligned} & \max_G \left\{ -(1-y) \log(1 - D(G(z))) \right\} \\ & = \max_G \left\{ -\log(1 - D(G(z))) \right\} \end{aligned}$$

- Ignore first term, because if  $y = 1$  the image is not generated

# Two Player Minimax Game

8

- Player D tries to minimize, and player G tries to maximize the comprehensive objective function

$$\min_D \max_G \left\{ -E_{x \sim Data} \log D(x) - E_{z \sim Noise} \log(1 - D(G(z))) \right\}$$

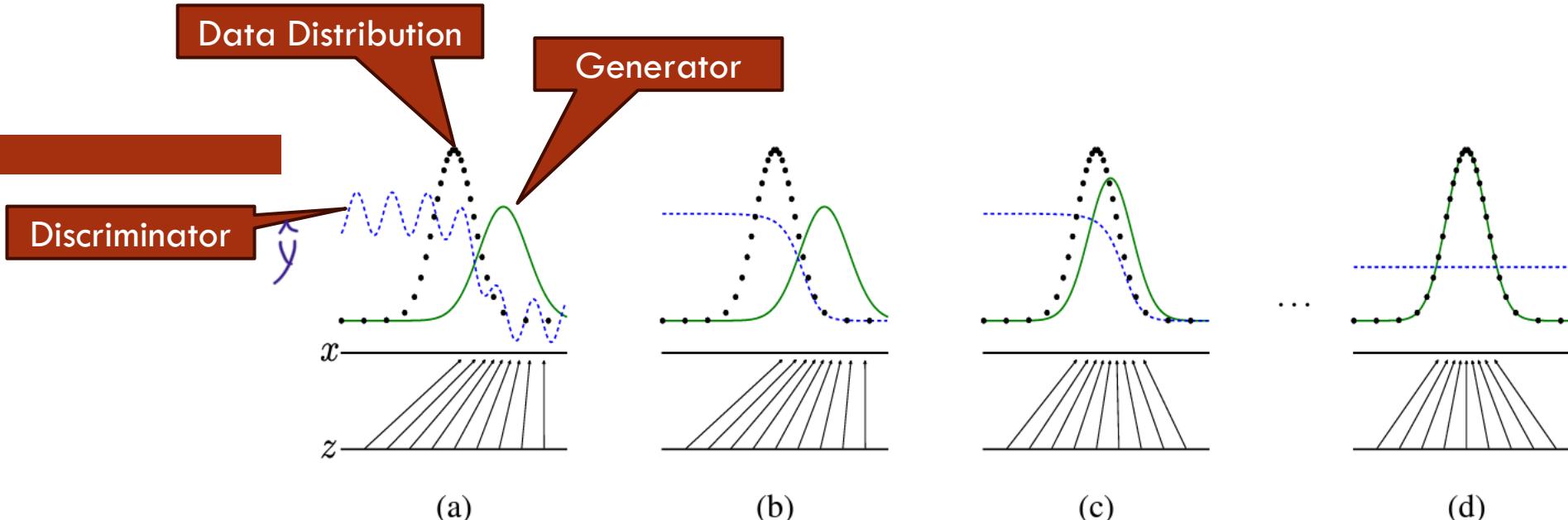


Figure 1: Generative adversarial nets are trained by simultaneously updating the discriminative distribution ( $D$ , blue, dashed line) so that it discriminates between samples from the data generating distribution (black, dotted line)  $p_{\mathbf{x}}$  from those of the generative distribution  $p_g$  (G) (green, solid line). The lower horizontal line is the domain from which  $\mathbf{z}$  is sampled, in this case uniformly. The horizontal line above is part of the domain of  $\mathbf{x}$ . The upward arrows show how the mapping  $\mathbf{x} = G(\mathbf{z})$  imposes the non-uniform distribution  $p_g$  on transformed samples.  $G$  contracts in regions of high density and expands in regions of low density of  $p_g$ . (a) Consider an adversarial pair near convergence:  $p_g$  is similar to  $p_{\text{data}}$  and  $D$  is a partially accurate classifier. (b) In the inner loop of the algorithm  $D$  is trained to discriminate samples from data, converging to  $D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$ . (c) After an update to  $G$ , gradient of  $D$  has guided  $G(\mathbf{z})$  to flow to regions that are more likely to be classified as data. (d) After several steps of training, if  $G$  and  $D$  have enough capacity, they will reach a point at which both cannot improve because  $p_g = p_{\text{data}}$ . The discriminator is unable to differentiate between the two distributions, i.e.  $D(\mathbf{x}) = \frac{1}{2}$ .

Goodfellow et al 2014

# Modified Generator Object

10

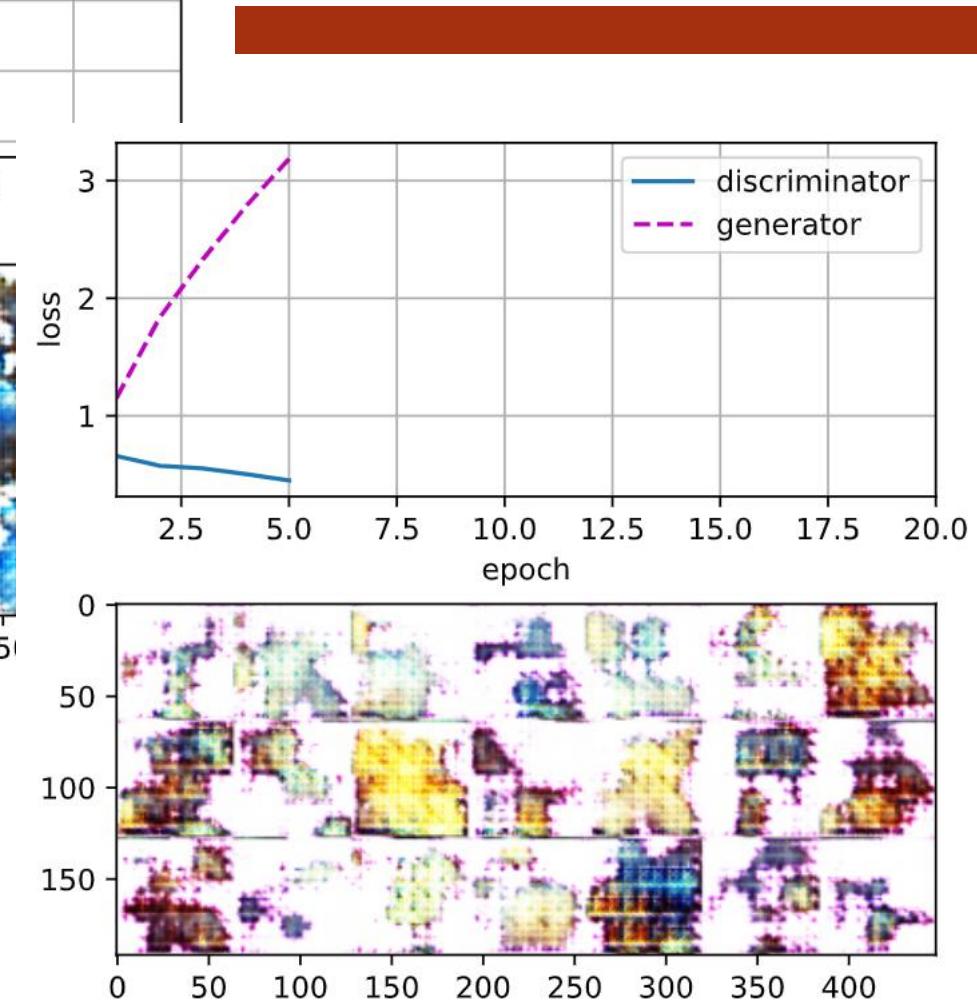
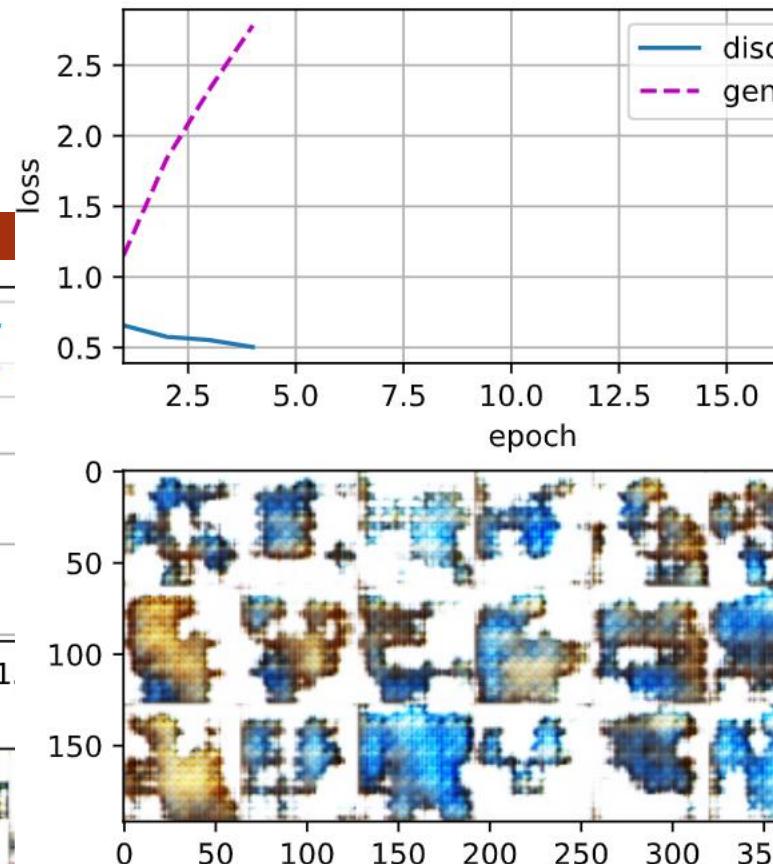
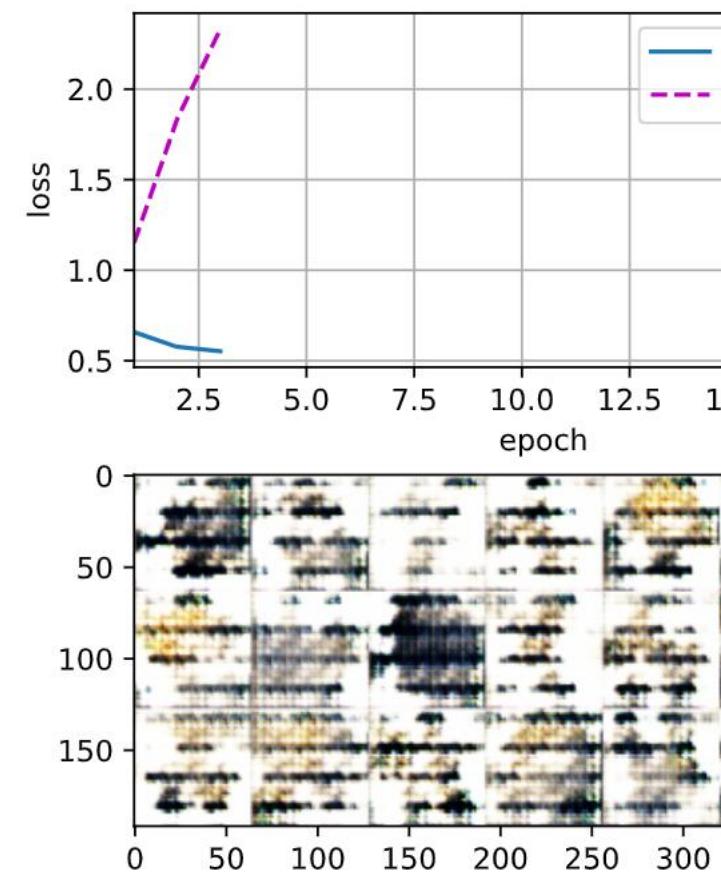
- With the current setup, the same function is being minimized by discriminator and maximized generator
- If the discriminator always successively reject generated samples with high confidence, then gradient for the generator vanishes
- Instead of  $\max_G \left\{ -\log \left( 1 - D(G(z)) \right) \right\}$   
Goodfellow suggest  $\min_G \left\{ -\log \left( D(G(z)) \right) \right\}$

# Notebook

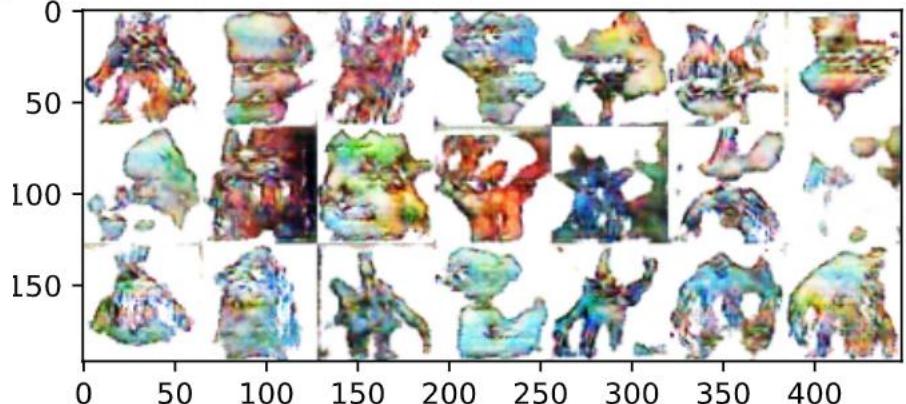
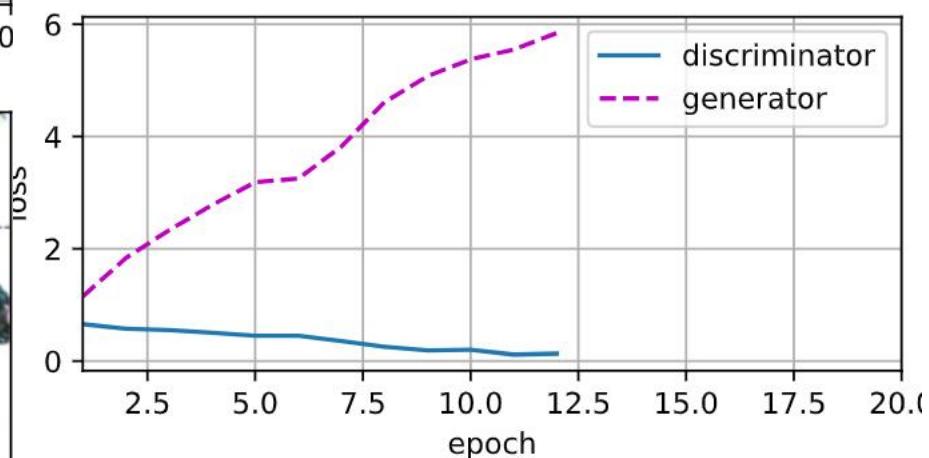
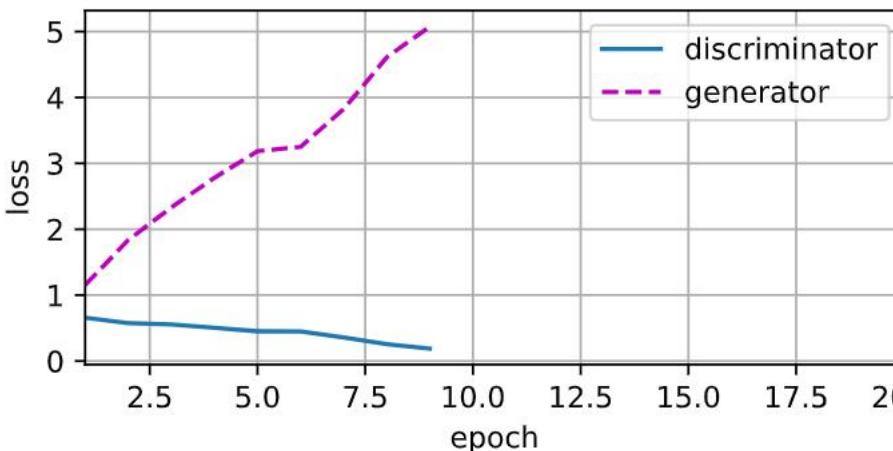
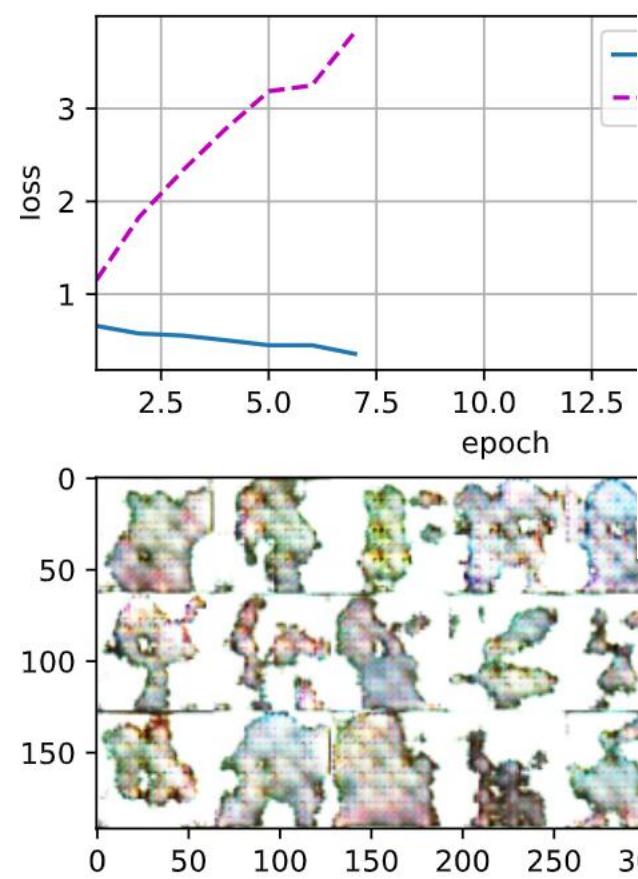
11

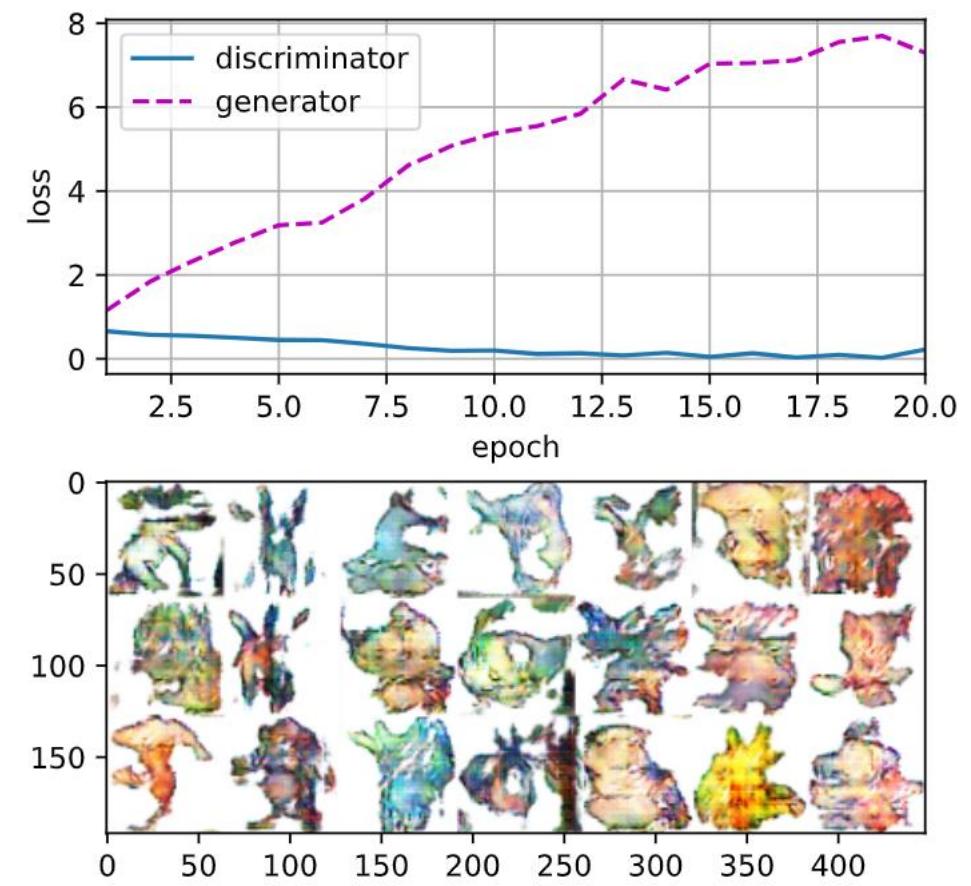
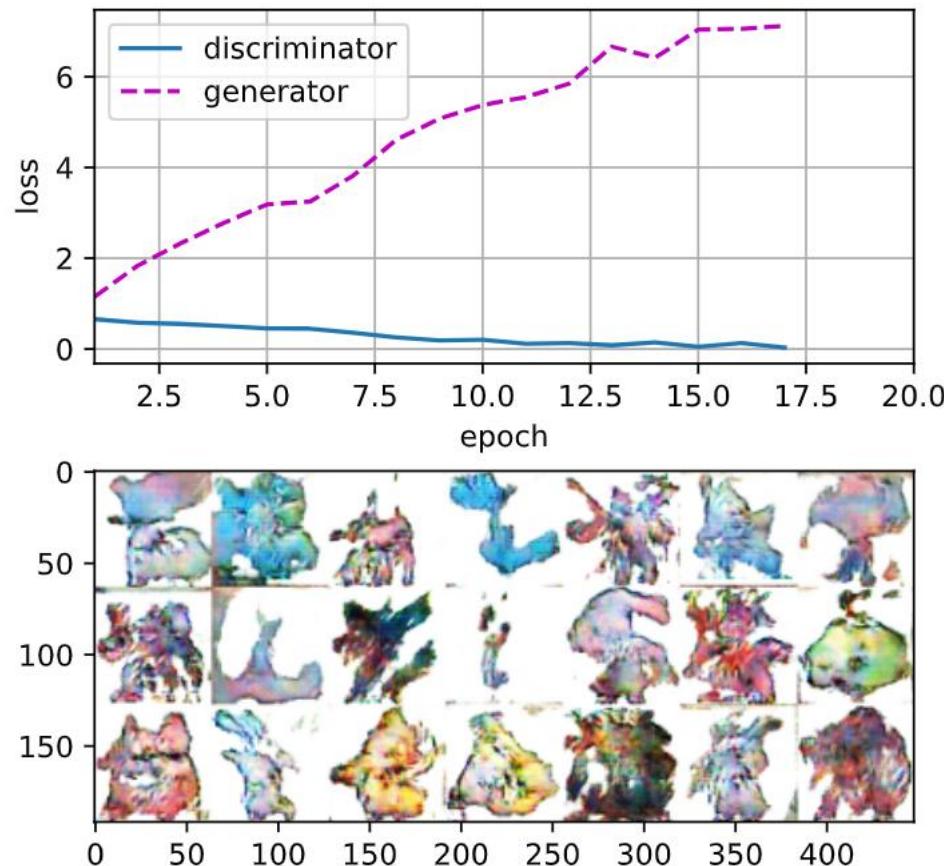
- See
  - chapter\_generative-adversarial-networks/gan.ipynb
- See Deep Convolutional GAN
  - chapter\_generative-adversarial-networks/dcgan.ipynb

12



13

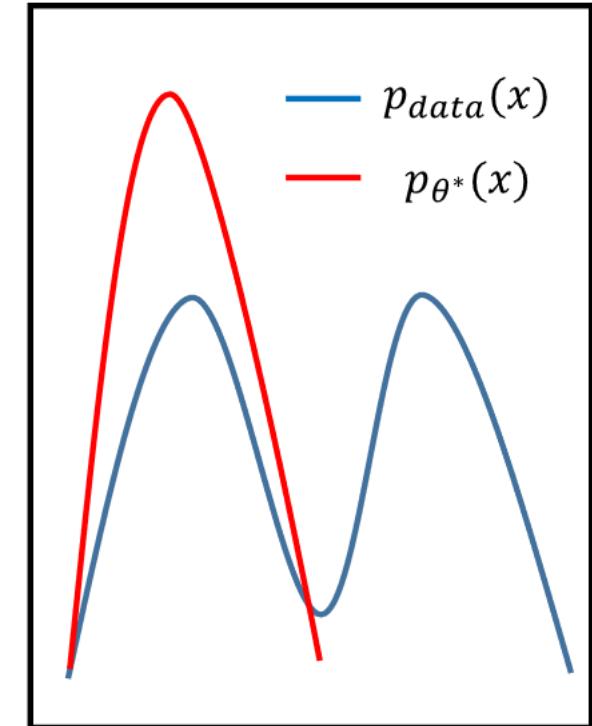




# Training Difficulties

15

- Mode collapse
  - ▣ Generator only covers one mode of the data, say cat images
- Vanishing gradients
  - ▣ Near perfect discriminator causes gradient to vanish at every data point
  - ▣ Generator cannot learn
- Convergence
  - ▣ Ideal is both networks stabilize, and generator converges to data probability distribution
  - ▣ In practice, network qualities oscillate



## Diffusion Generative Models



# Approaches Covered

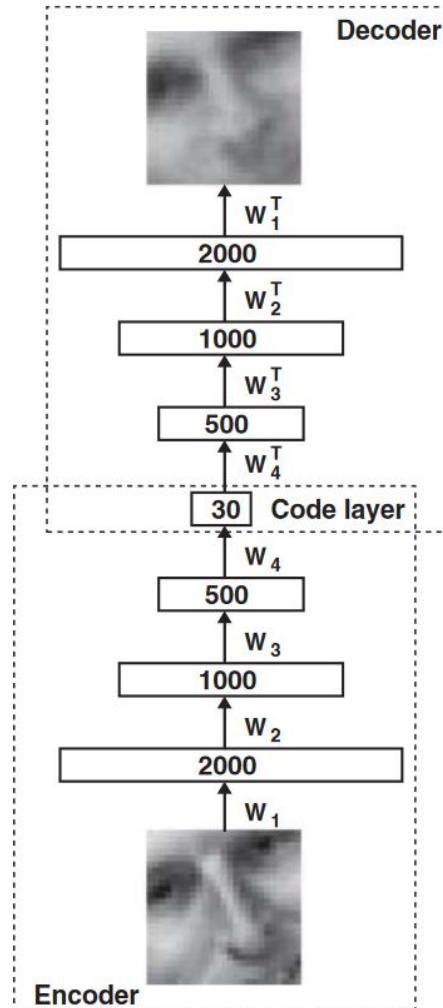
17

- Autoencoder
- Variational autoencoder
- Denoising Diffusion
- Text-to-Image Diffusion

# Autoencoder

18

- Autoencoder neural network learns to copy the input images to output images
- An autoencoder consists of
  - Encoder network that compress high-dimensional image into a lower-dimensional embedding space
  - Decoder network that decompress a given embedding vector into a high-dimensional image
- A bottleneck is introduced to force the network to learn latent representations
- Reconstruction loss penalize difference between input and output using either cross-entropy or RMSE



# Convolutional Autoencoder

19

## Encoder

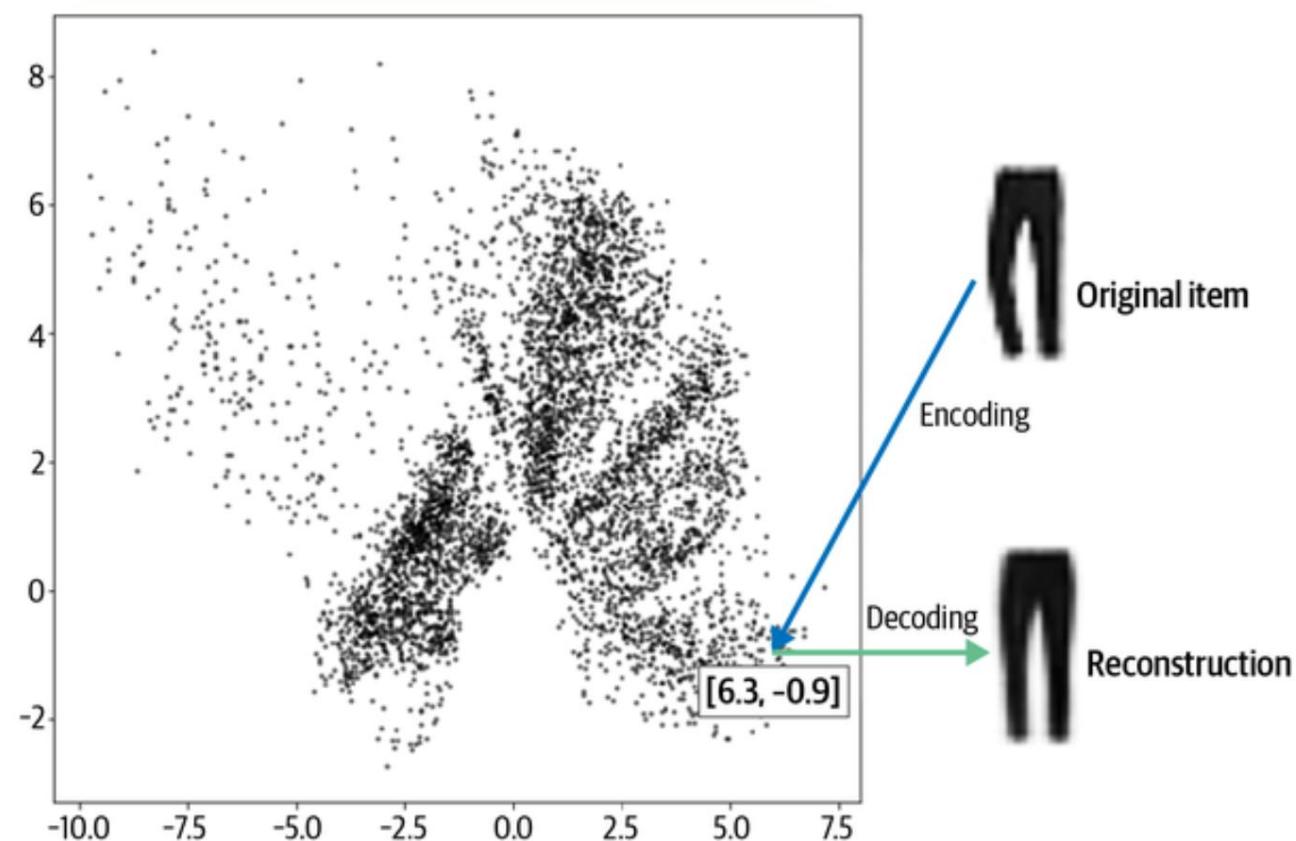
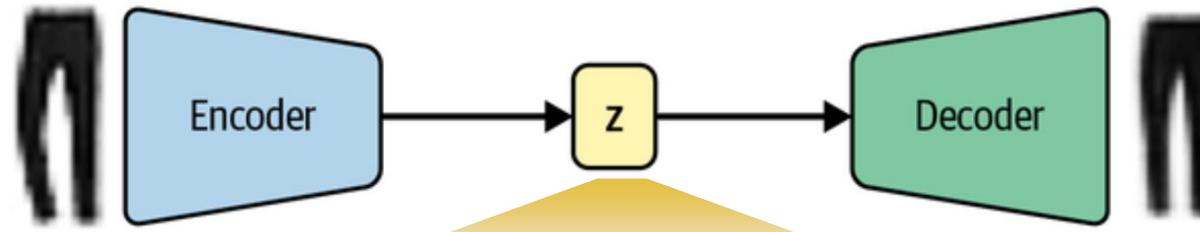
Layer (type)	Output shape	Param #
InputLayer	(None, 32, 32, 1)	0
Conv2D	(None, 16, 16, 32)	320
Conv2D	(None, 8, 8, 64)	18,496
Conv2D	(None, 4, 4, 128)	73,856
Flatten	(None, 2048)	0
Dense	(None, 2)	4,098

## Decoder

Layer (type)	Output shape	Param #
InputLayer	(None, 2)	0
Dense	(None, 2048)	6,144
Reshape	(None, 4, 4, 128)	0
Conv2DTranspose	(None, 8, 8, 128)	147,584
Conv2DTranspose	(None, 16, 16, 64)	73,792
Conv2DTranspose	(None, 32, 32, 32)	18,464
Conv2D	(None, 32, 32, 1)	289

# Latent Representation

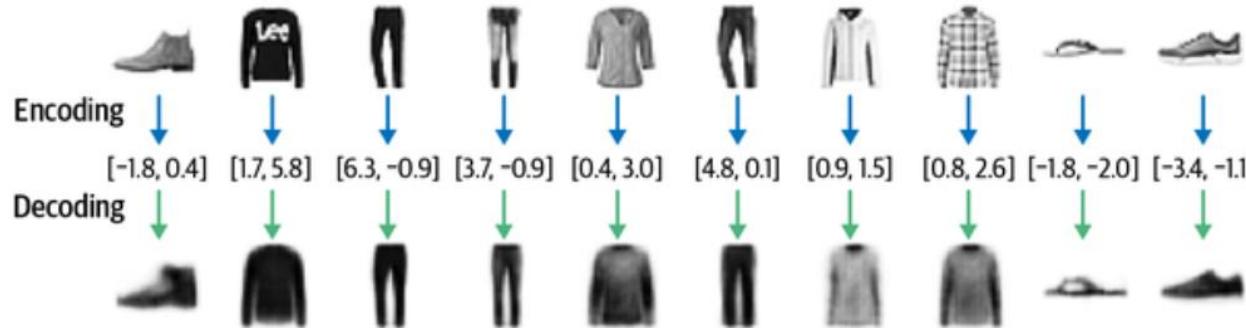
20



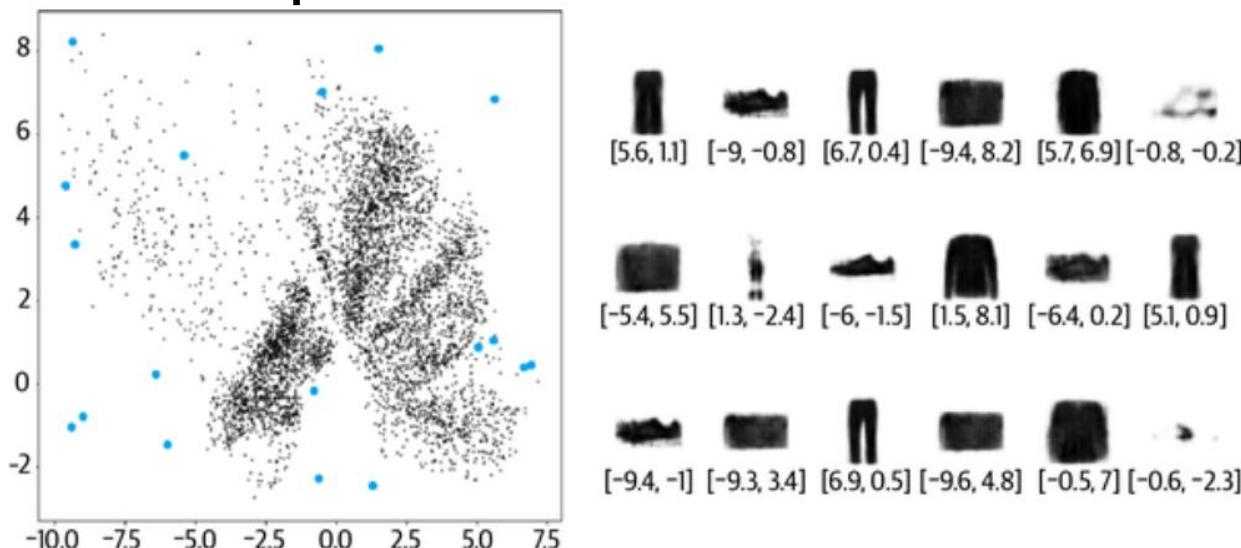
# Decoding Embedding Vectors to Generate Images

21

## □ Decoding from actual encoded vectors



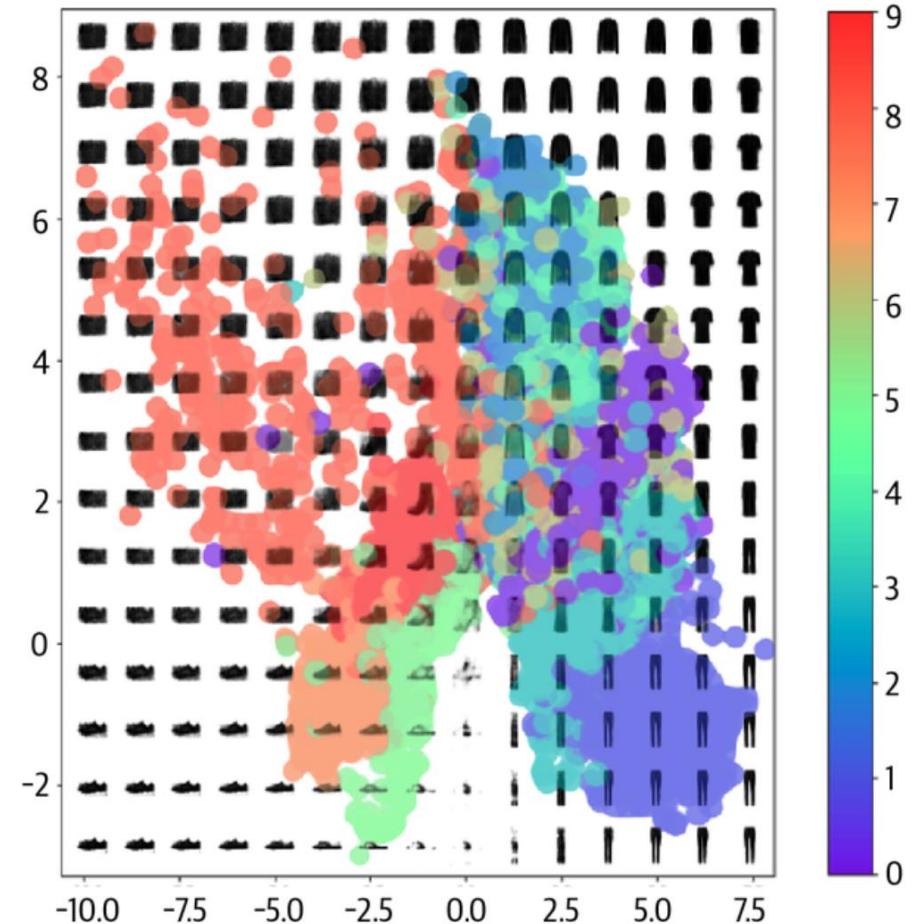
## □ Decoding from sampled vectors



# Sampling Difficulties

22

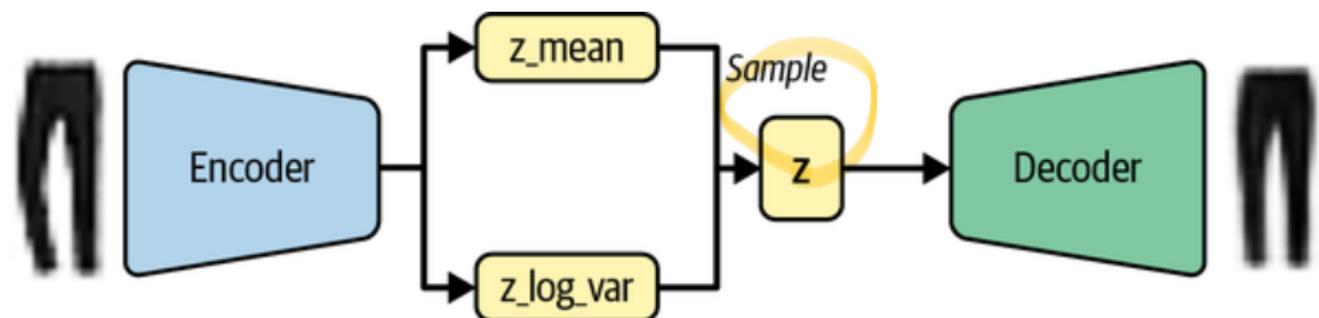
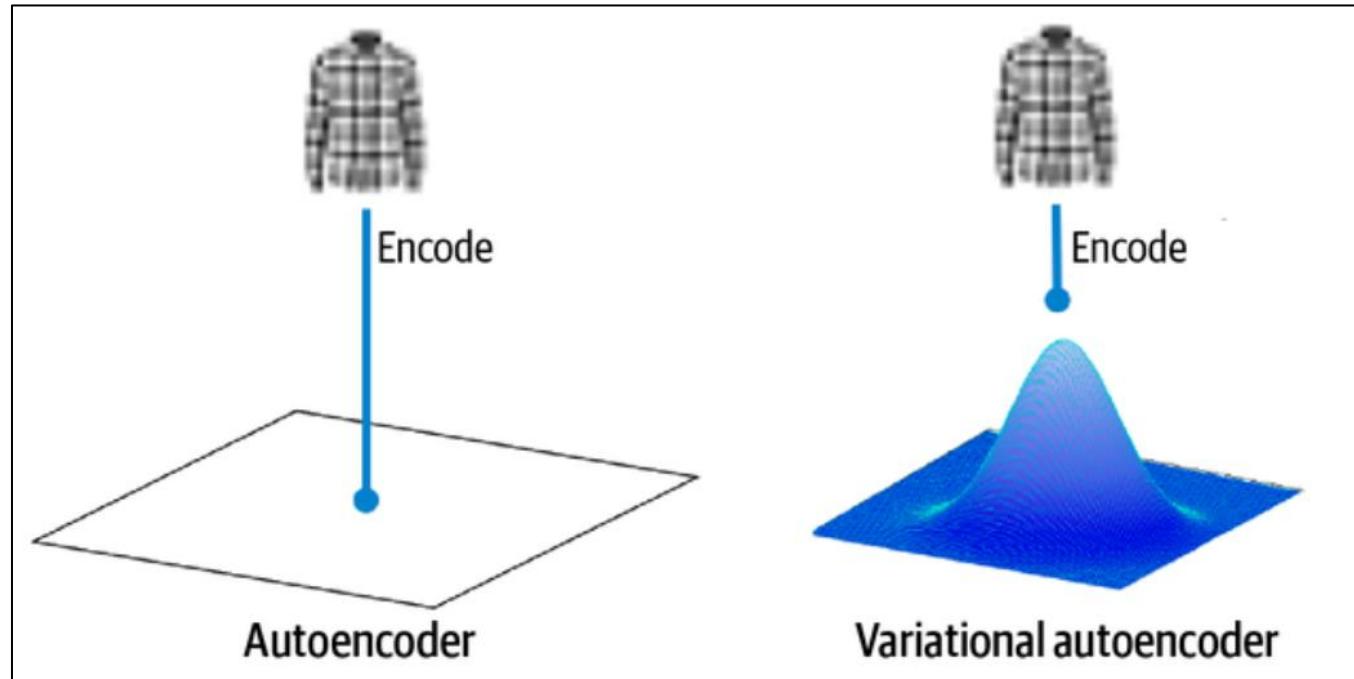
- Some clothing items are overrepresented in latent space area, e.g., bags (id=8)
- Large latent space regions with no training data, i.e., white areas
- Value range across latent space dimensions vary
- Sampling difficulties increases with higher dimensional latent space
- Want simpler probability distributions for latent space



# Variational Autoencoder

23

- Probabilistic version of autoencoders
- Encoder learns to represent each clothing type as a probabilistic region in latent space
- Give an image the encoder
  - Predicts distribution mean and (log of) variance
  - Returns a sampled point from the corresponding normal distribution
- Decoder learns to generate clothing type associate with the region

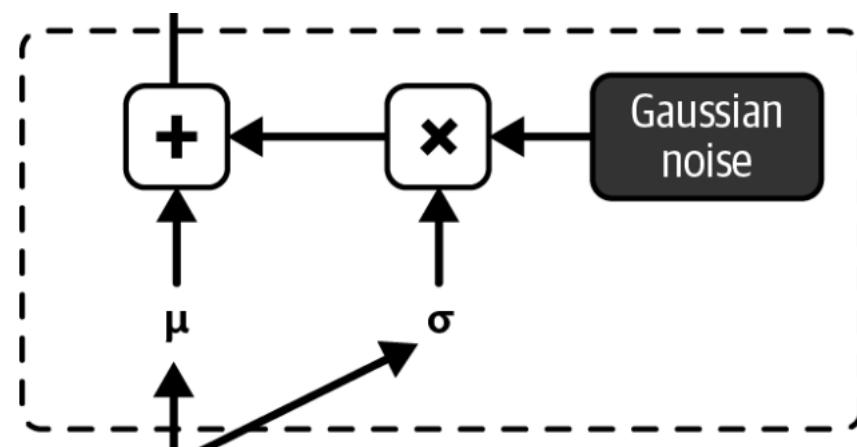


# Reparameterization Trick

24

- Even though the encoder predicts regions  $(\mu, \sigma)$ , it returns a sampled point in the region
- From a specific point, decoder generates a specific instance of clothing
  - Decoder should not generate an average image of clothing
- However, sampling from distribution is **not differentiable**
- **Reparameterization Trick**
  - Sample Gaussian noise  $\epsilon \sim \mathcal{N}(0,1)$
  - Compute embedding point:  $\mu + \sigma \times \epsilon$
  - $\epsilon$  is constant with respect to  
automatic differentiation

$\downarrow$   
*differentiable*



# Loss Function

25

- In addition to reconstruction loss, add another term to force latent distributions to look like unit Gaussians
  - Unit Gaussians are easy to sample
- Add Kullback-Leibler (KL) divergence term that measures the difference between two probability distributions

$$D_{KL}[\mathcal{N}(\mu, \sigma) \parallel \mathcal{N}(0,1)] = -\frac{1}{2}(1 + \log(\sigma^2) - \mu^2 - \sigma^2)$$

# Generating Faces with CelebA Dataset

26

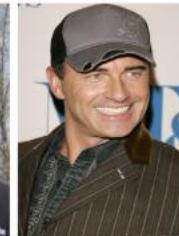
- <https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>
- Variational Autoencoder
  - Encoder with 4 CNN layers  
(Batch normalization,  
LeakyReLU)
  - 200-dimensional latent space
  - Decoder with 4 CNN  
ConvTranspose layers (Batch  
normalization, LeakyReLU)

Sample Images

Eyeglasses



Wearing Hat



Bangs



Wavy Hair



Pointy Nose



Mustache



Oval Face



Smiling



# Network Details

27

Encoder network

Decoder network

Layer (type)	Output shape	Param #	Connected to	Layer (type)	Output shape	Param #
InputLayer (input)	(None, 32, 32, 3)	0	[]	InputLayer	(None, 200)	0
Conv2D (conv2d_1)	(None, 16, 16, 128)	3,584	[input]	Dense	(None, 512)	102,912
BatchNormalization (bn_1)	(None, 16, 16, 128)	512	[conv2d_1]	BatchNormalization	(None, 512)	2,048
LeakyReLU (lr_1)	(None, 16, 16, 128)	0	[bn_1]	LeakyReLU	(None, 512)	0
Conv2D (conv2d_2)	(None, 8, 8, 128)	147,584	[lr_1]	Reshape	(None, 2, 2, 128)	0
BatchNormalization (bn_2)	(None, 8, 8, 128)	512	[conv2d_2]	Conv2DTranspose	(None, 4, 4, 128)	147,584
LeakyReLU (lr_2)	(None, 8, 8, 128)	0	[bn_2]	BatchNormalization	(None, 4, 4, 128)	512
Conv2D (conv2d_3)	(None, 4, 4, 128)	147,584	[lr_2]	LeakyReLU	(None, 4, 4, 128)	0
BatchNormalization (bn_3)	(None, 4, 4, 128)	512	[conv2d_3]	Conv2DTranspose	(None, 8, 8, 128)	147,584
LeakyReLU (lr_3)	(None, 4, 4, 128)	0	[bn_3]	BatchNormalization	(None, 8, 8, 128)	512
Conv2D (conv2d_4)	(None, 2, 2, 128)	147,584	[lr_3]	LeakyReLU	(None, 8, 8, 128)	0
BatchNormalization (bn_4)	(None, 2, 2, 128)	512	[conv2d_4]	Conv2DTranspose	(None, 16, 16, 128)	147,584
LeakyReLU (lr_4)	(None, 2, 2, 128)	0	[bn_4]	BatchNormalization	(None, 16, 16, 128)	512
Flatten (flatten)	(None, 512)	0	[lr_4]	LeakyReLU	(None, 16, 16, 128)	0
Dense (z_mean)	(None, 200)	102,600	[flatten]	Conv2DTranspose	(None, 32, 32, 128)	147,584
Dense (z_log_var)	(None, 200)	102,600	[flatten]	BatchNormalization	(None, 32, 32, 128)	512
Sampling (z)	(None, 200)	0	[z_mean, z_log_var]	LeakyReLU	(None, 32, 32, 128)	0
				Conv2DTranspose	(None, 32, 32, 3)	3,459

# Reconstructions and Generations

28

Example real faces

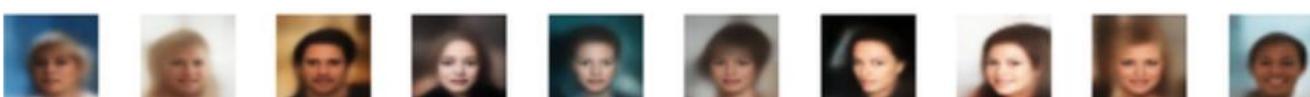
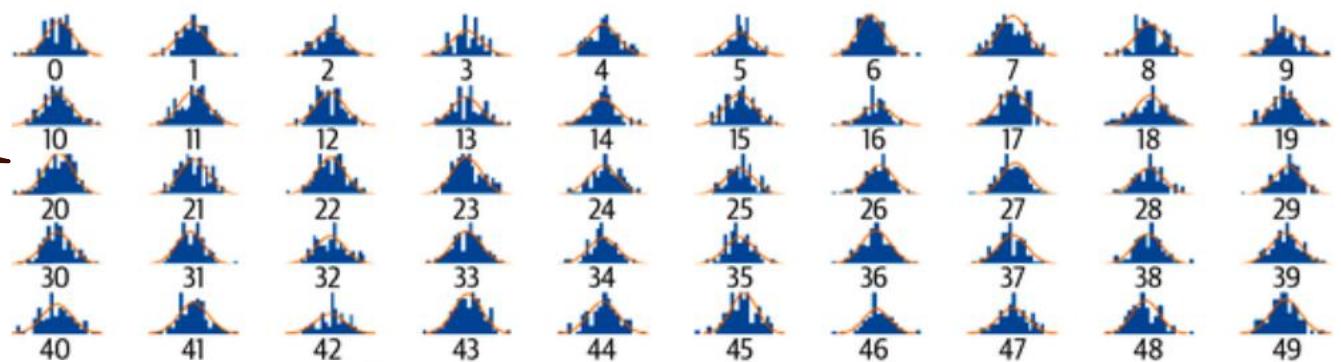


Reconstructions



Reconstructions

Histogram of latent variables  
“Gaussian Noise”



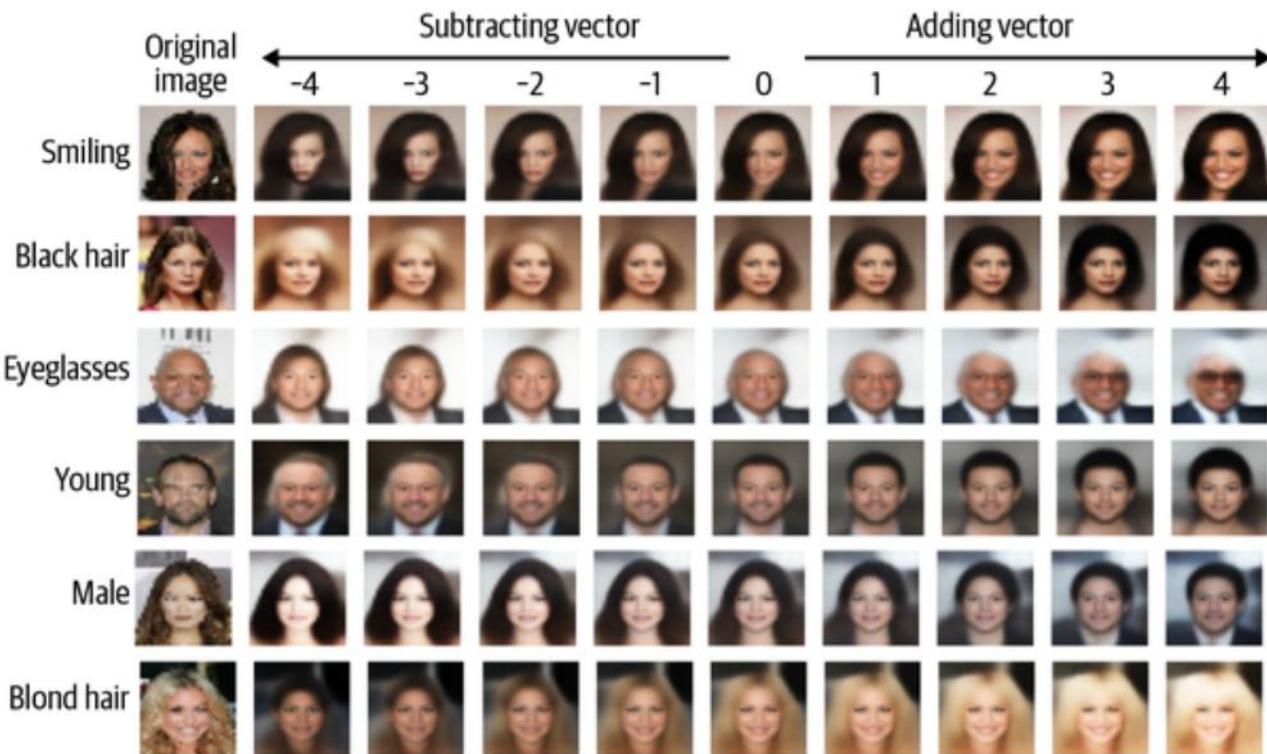
New generated Faces

# Latent Space Arithmetic

add/subtract

29

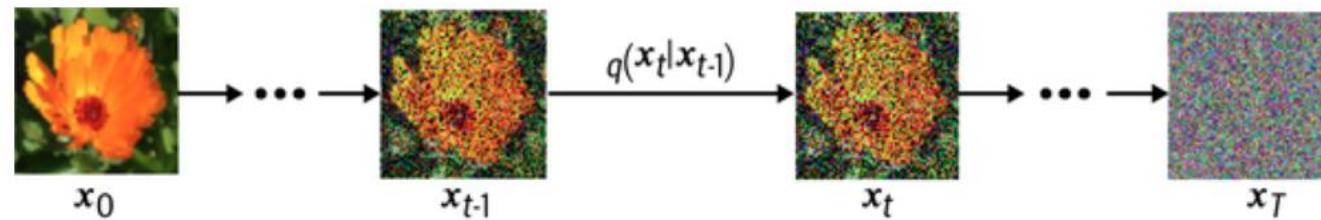
- Use linear interpolation between two faces
- Add/subtract vectors
  - Use dataset label to find average vector for attributes
  - Such as attribute vectors for smiling, black hair, eyeglasses, ...



# Denoising Diffusion Models

30

- Approach inspired by the diffusion process of **nonequilibrium thermodynamics**
- In thermodynamics diffusion is a process where particles move from higher concentration to low concentration by increase in entropy
- For images, diffusion process can be thought of as **adding noise**



- What if we can learn to **reverse this process**?

# Forward Diffusion Process

31

- Process:
  - Start with image  $x_0$  at time zero
  - Incremental add noise to create  $x_t$
  - Until at time  $T$ , image  $x_T$  is Gaussian noise (zero mean, unit variance)
- Since we want Gaussian noise in the end
  - Normalize  $x_0$  to have zero mean and unit variance
  - At each step maintain zero mean and unit variance

# Forward Diffusion Process Details

32

- To add noise and maintain zero mean unit variance, we take linear combination of previous image and Gaussian noise

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \epsilon_{t-1}$$

- Above maintains zero mean unit variance, because

$$\begin{aligned}Var(X + Y) &= Var(X) + Var(Y) \\Var(aX) &= a^2 Var(X)\end{aligned}$$

- Let the forward diffusion process be a first-order Markov process

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

# Shortcut Forward Diffusion Process

33

- Sample image  $\mathbf{x}_t$  from  $\mathbf{x}_0$  directly without computing  $\mathbf{x}_{t-1}, \mathbf{x}_{t-2}, \dots, \mathbf{x}_1$
- Possible because sum of independent Gaussians is Gaussian
- Let  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \alpha_1 \times \alpha_2 \times \dots \times \alpha_t = \prod_{i=1}^t \alpha_i$
- Shortcut equation:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

# Diffusion Schedule

34

Linear

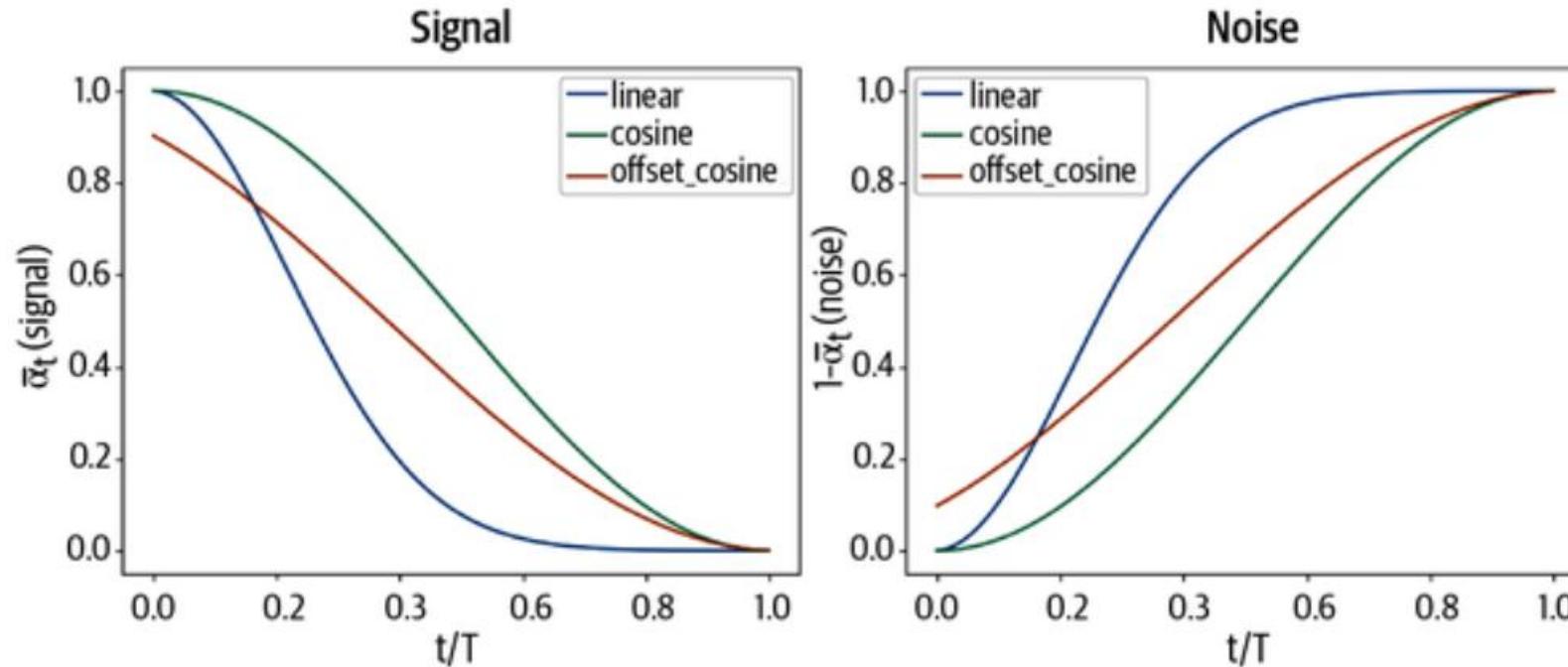


Cosine

- Original paper Ho et al, 2020 proposed linear schedule
  - Schedule is too aggressive in adding noise
  - Images from last quarter of the time are just noise
- Nichol and Dhariwal, 2021 proposed less aggressive cosine schedule

# Signal and Noise Curve Comparisons

35

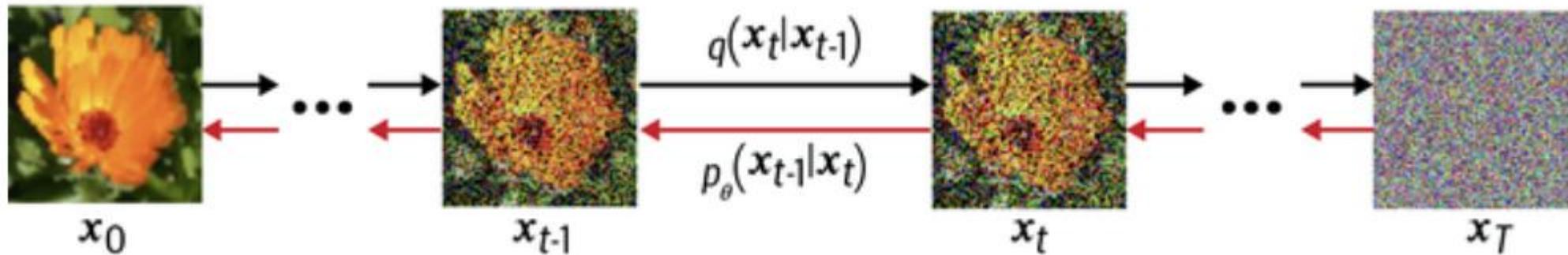


- Linear schedule adds noise the fastest
- Cosine schedule adds noise the slowest
- Foster book proposes offset cosine that is between the two schedules, but it starts with large initial noise

# Reverse Diffusion Process

36

- Build neural network  $p_\theta(x_{t-1}|x_t)$  to learn to reverse diffusion process



- Reverse diffusion process start with Gaussian noise image then generate picture image
- Like Variational Autoencoder, start with **Gaussian latent variables** then generate picture image

# Training Procedure

- Training instance creation
  - Randomly select an image  $x_0$
  - Randomly select time  $t \in [1, \dots, T]$ , say  $T = 1000$
  - Sample Gaussian noise image  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - Create noisy image at  $t$ ,  $x_t = \sqrt{\bar{\alpha}_t}x_0 + (1 - \bar{\alpha}_t)\epsilon$
  - Return  $(t, x_t, \epsilon)$
- Train network  $\epsilon_\theta$  to predict noise image with mean squared error loss
  - $\epsilon_\theta$  is neural network with parameter  $\theta$
- Why train to predict noise  $\epsilon_\theta$  instead of original image  $x_0$ ?
  - Empirical, Ho et al, 2020 tried both ways. Training to predict noise image is more stable.

---

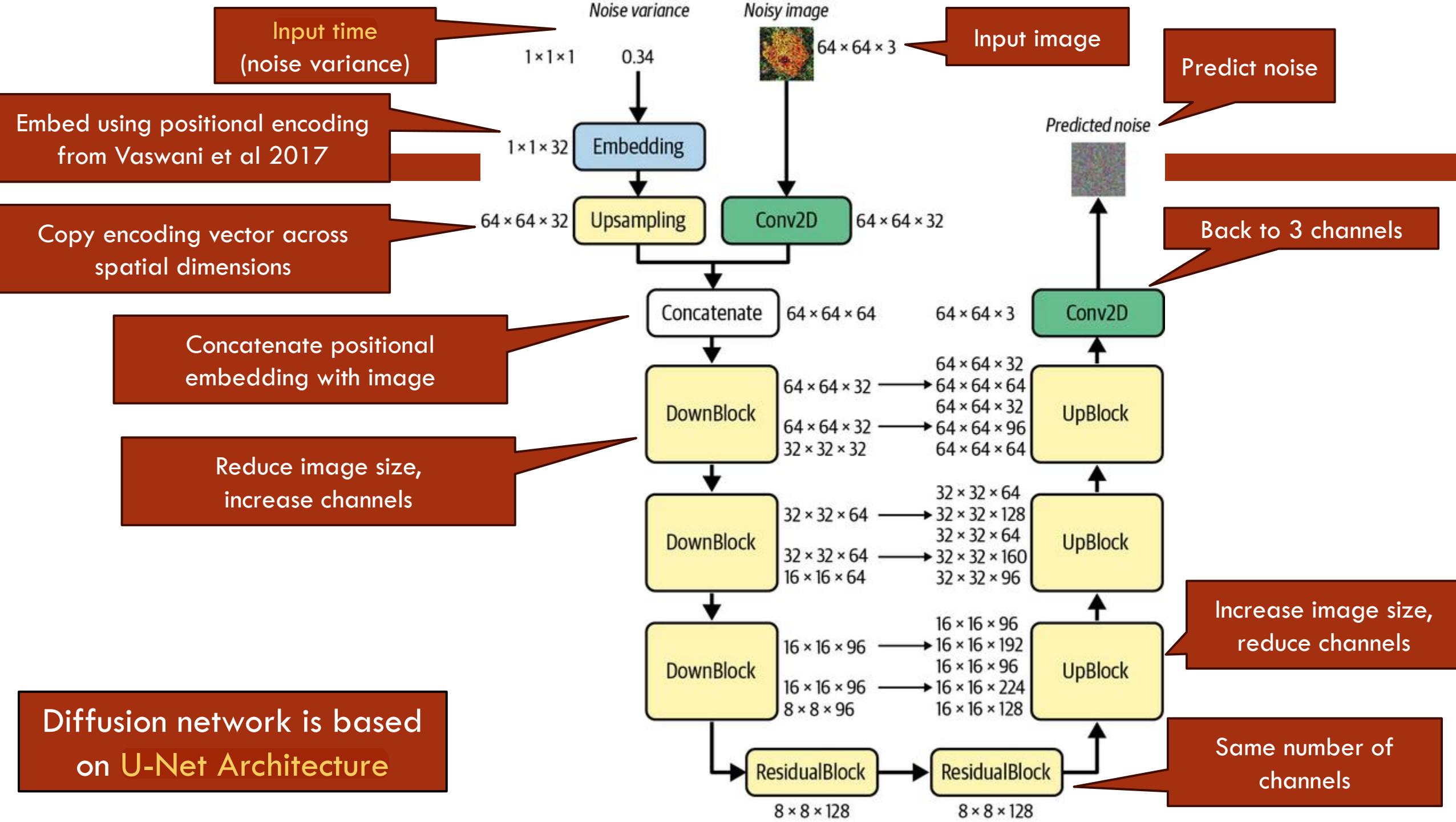
## Algorithm 1 Training

---

```

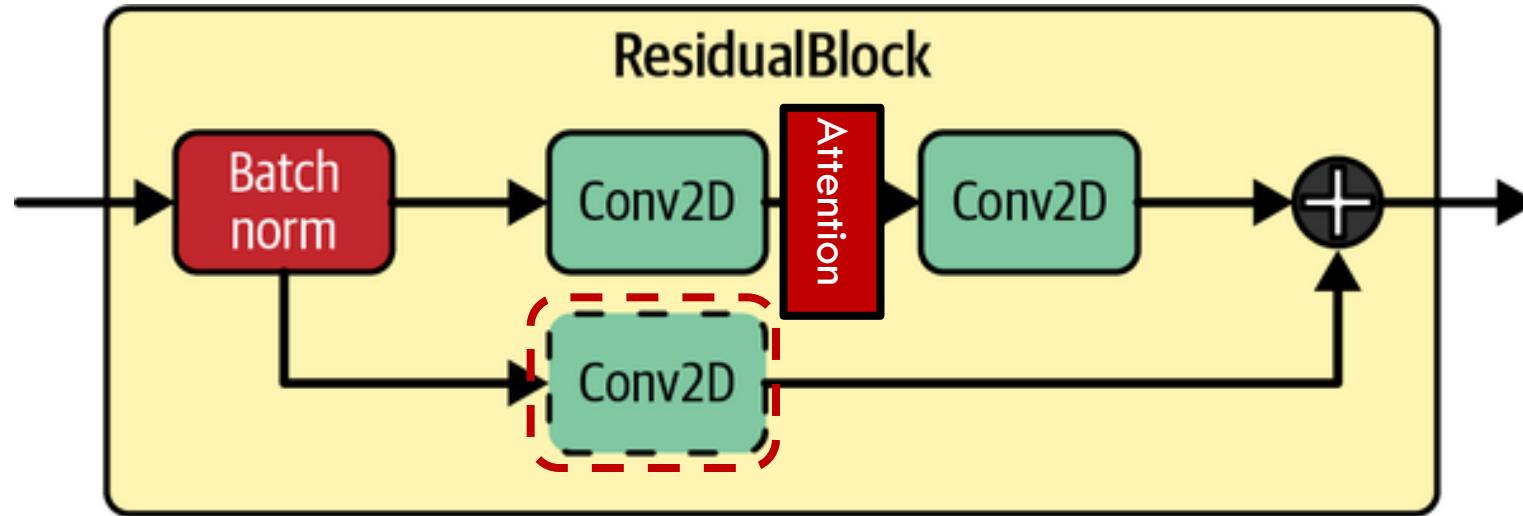
1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
      $\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$ 
6: until converged
  
```

---



# Residual Block

39

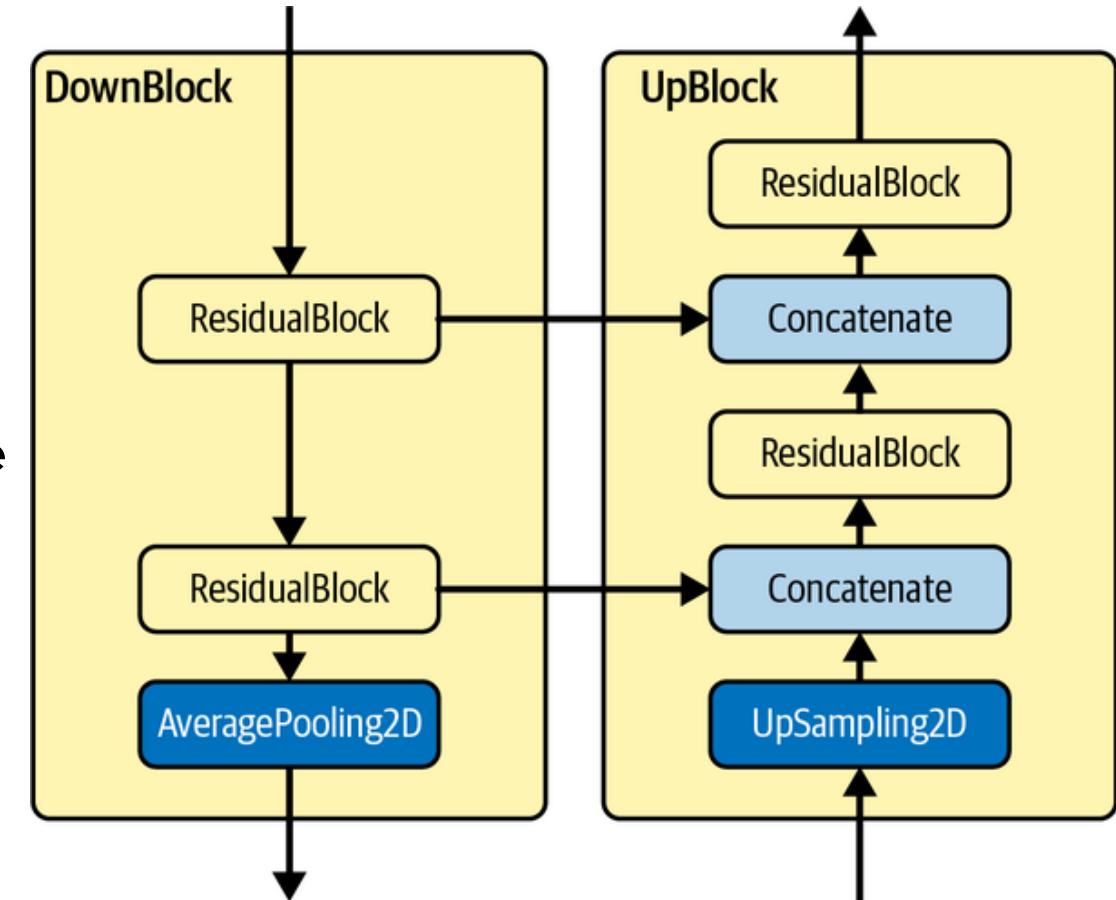


- Similar to ResNet residual blocks
- Dotted “Conv2D” is  $1 \times 1$  convolution
- Ho et al, 2020 adds **self-attention** between the two Conv2D layer for some ResidualBlocks
  - ▣ “Pixels” attend to other pixels

# DownBlocks and UpBlocks

40

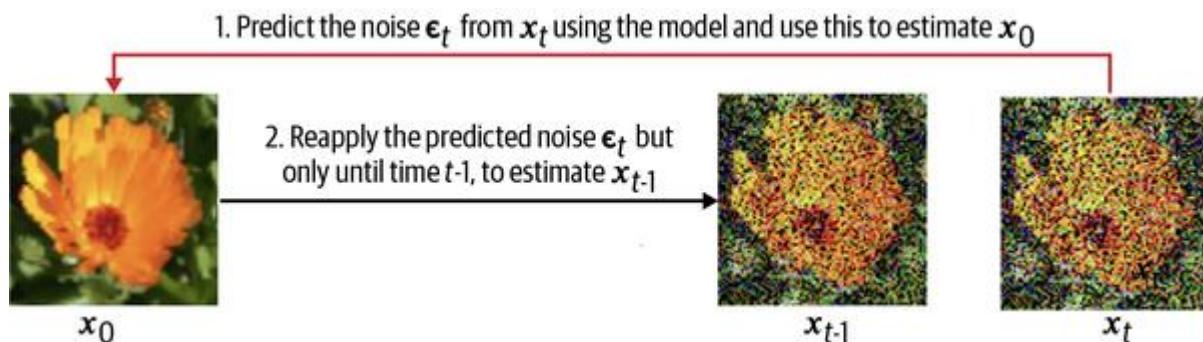
- DownBlocks and UpBlocks pairs connected by skip connections
- DownBlock applies “convolutions” and down samples
- UpBlock up samples and applies “convolutions”
- Before convolutions UpBlock concatenate corresponding tensor from DownBlock
- Rational
  - ▣ Reducing image size captures higher level semantic information
  - ▣ But spatial information is lost
  - ▣ Skip connection adds back spatial information



# Generating Images

41

- Model is trained to predict all the noise and recover the original image in one step
- However, we will take multiple steps
- Each step subtracting a small fraction of the predicted noise



---

## Algorithm 2 Sampling

---

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

---

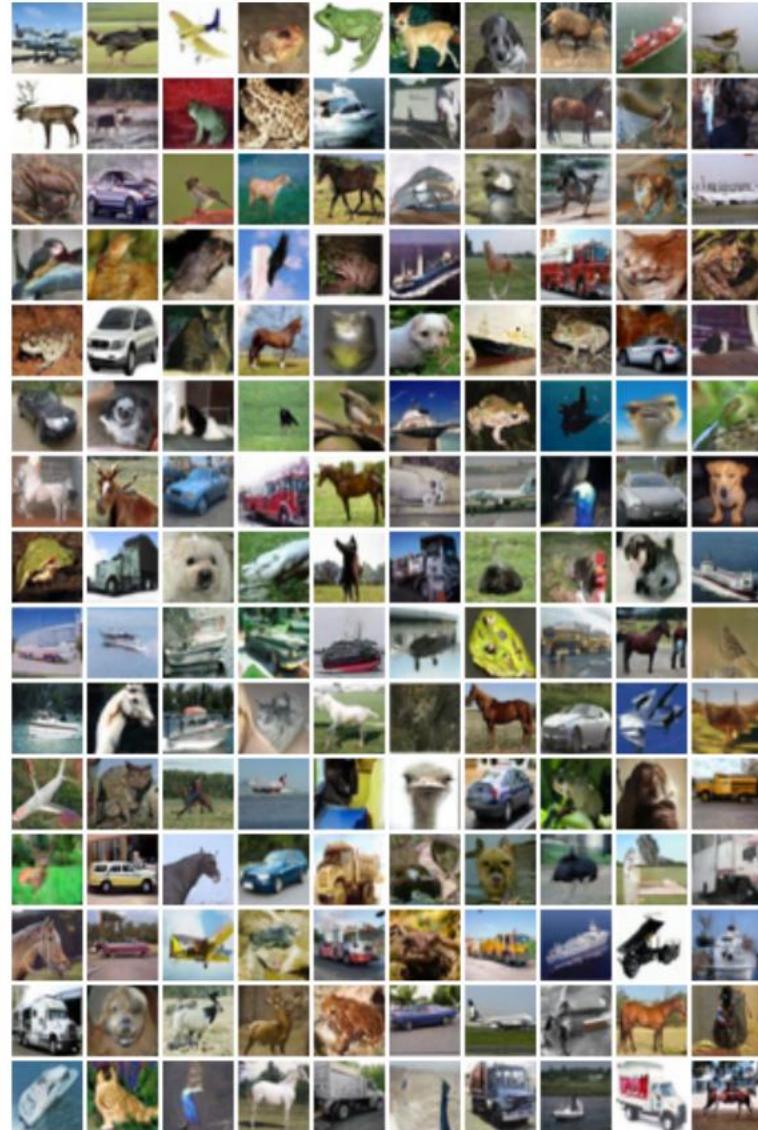


Figure 1: Generated samples on CelebA-HQ  $256 \times 256$  (left) and unconditional CIFAR10 (right)

# FID (Fréchet inception distance) Metric

43

- Distance metric used to measure the distance between the set of training images and the set of generated images
- Encode images in each set using Inception v3 net gather statistics (mean  $\mu$  and covariance  $\Sigma$ ) on the activations of the penultimate layer (2000 features)
- Use Fréchet distance to compare the difference

$$d_F((\mu_1, \Sigma_1), (\mu_2, \Sigma_2)) = \|\mu_1 - \mu_2\|_2^2 + \text{tr}(\Sigma_1 + \Sigma_2 - 2(\Sigma_1 \Sigma_2)^{\frac{1}{2}})$$

- Proposed by Heusel et al, 2018

# FID between CelebA Dataset and Its Transformations

44

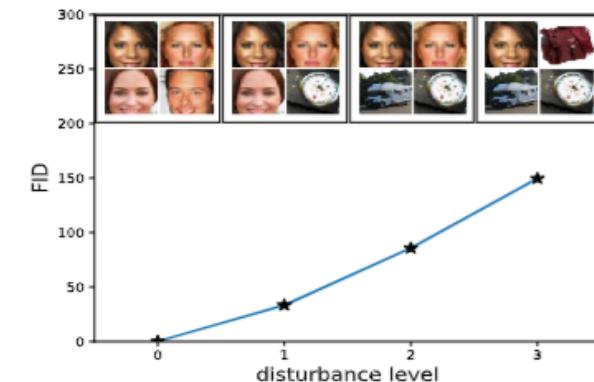
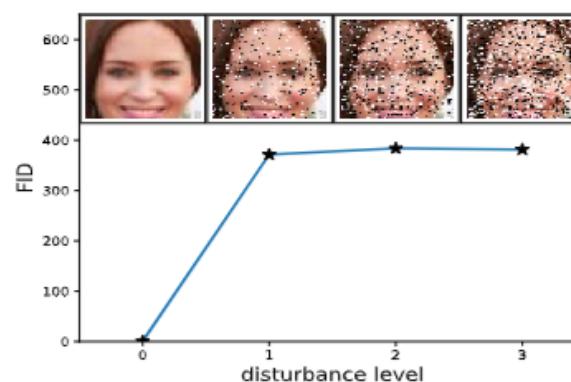
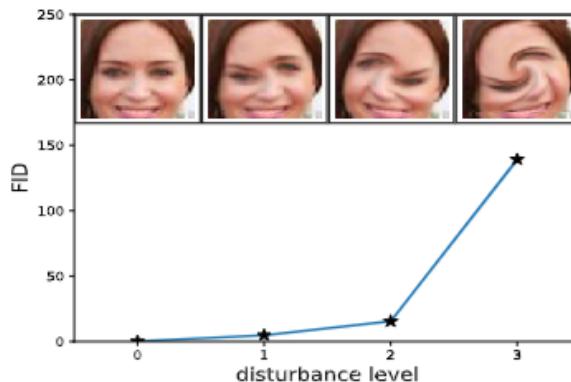
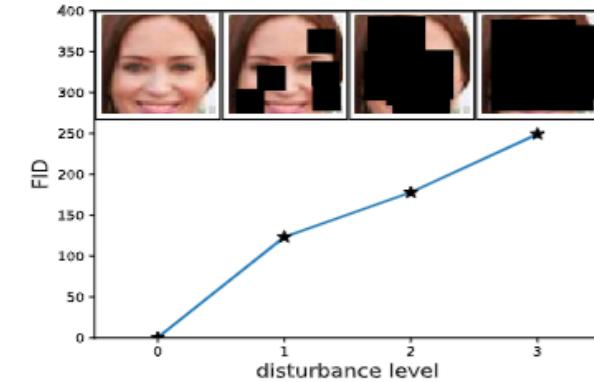
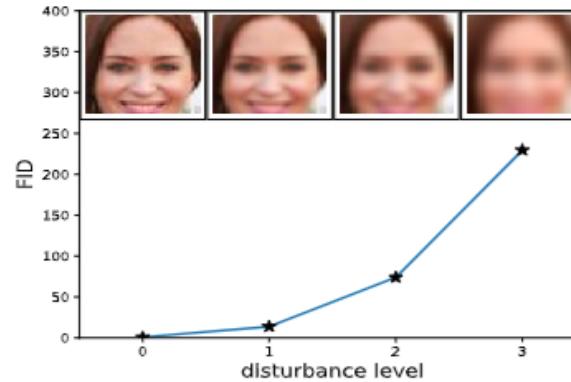
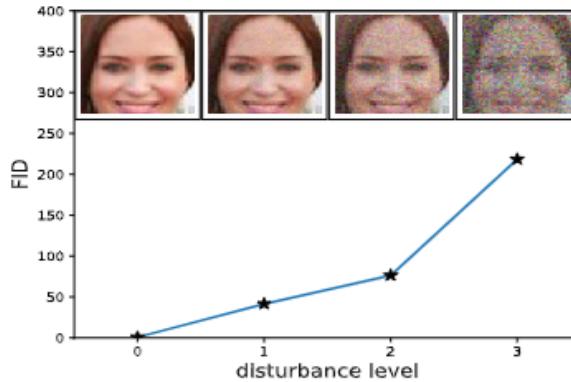
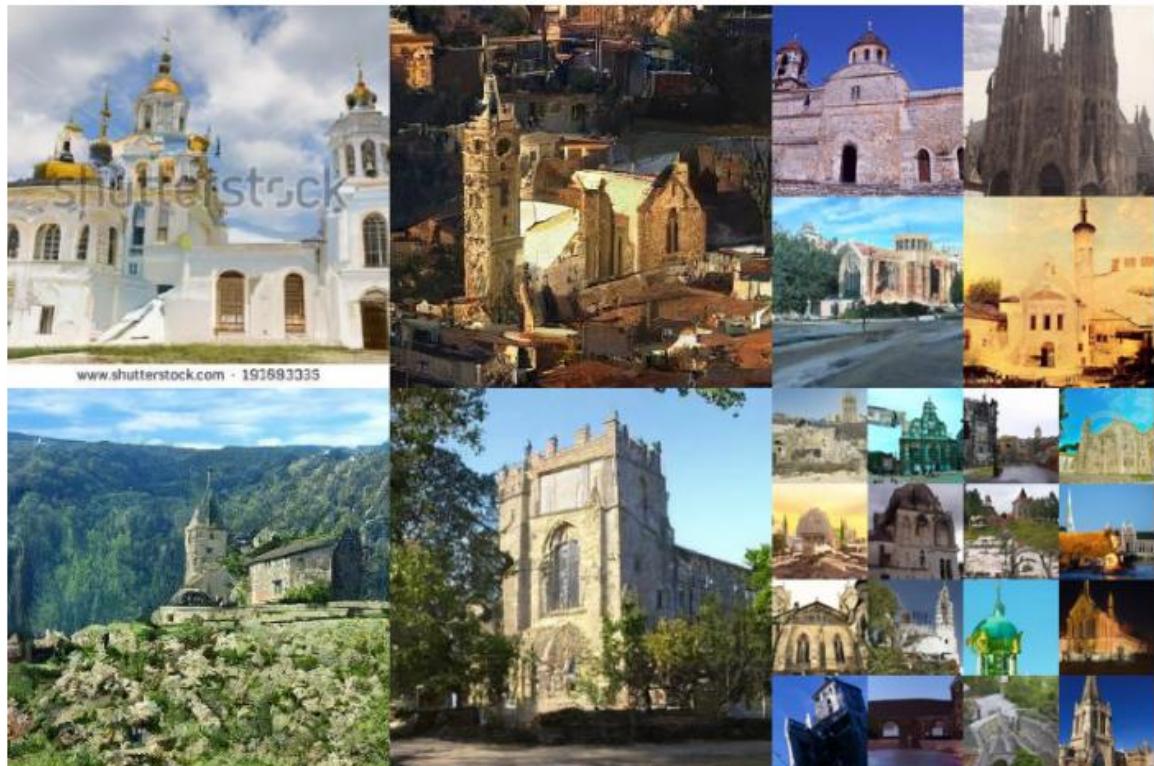


Figure 3: FID is evaluated for **upper left**: Gaussian noise, **upper middle**: Gaussian blur, **upper right**: implanted black rectangles, **lower left**: swirled images, **lower middle**: salt and pepper noise, and **lower right**: CelebA dataset contaminated by ImageNet images. The disturbance level rises from zero and increases to the highest level. The FID captures the disturbance level very well by monotonically increasing.

# FID Distance on Generate Churches and Bedrooms

45



www.shutterstock.com · 191693335

www.shutterstock.com · 191693335

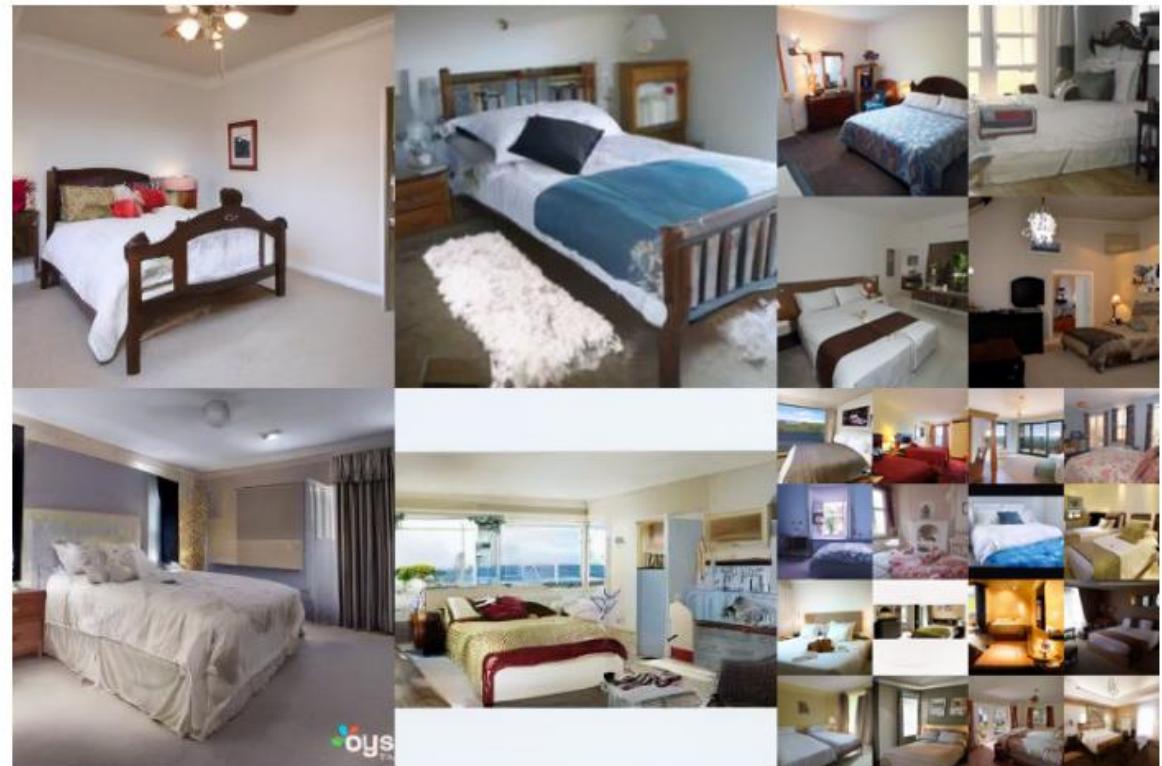


Figure 3: LSUN Church samples. FID=7.89

Figure 4: LSUN Bedroom samples. FID=4.90

# Text-to-Image Diffusion Models

46

- Sampling from a diffusion model generates a **random image** from the distribution from which the model is trained on
- In text-to-image generation we wish to guide the generation process using **text prompts**
- These models are called **conditional diffusion models**, since generation probability distribution is conditioned on the text prompts
- Example models include Stable Diffusion from Stability AI, DALL-E from OpenAI, Imagen from Google

# Imagen: Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding

47

- Saharia et al, 2022
- Approach
  - Use pretrained transformer-based LLM (T5) to encode text prompt
  - Use conditional diffusion model to generate small images
  - Use cascaded diffusion models to generate larger images
- Can generate photorealistic images as well as images various styles, such as oil painting, illustration, surreal, ...

# Sample Imagen Images

48



Sprouts in the shape of text 'Imagen' coming out of a fairytale book.



A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.



A high contrast portrait of a very happy fuzzy panda dressed as a chef in a high end kitchen making dough. There is a painting of flowers on the wall behind him.

# Samples From the Same Prompt

49



A brown bird and a blue bear.



One cat and two dogs sitting on the grass.



A sign that says 'NeurIPS'.



A small blue book sitting on a large red book.



A blue coloured pizza.



A wine glass on top of a dog.



A pear cut into seven pieces arranged in a ring.



A photo of a confused grizzly bear in calculus class.

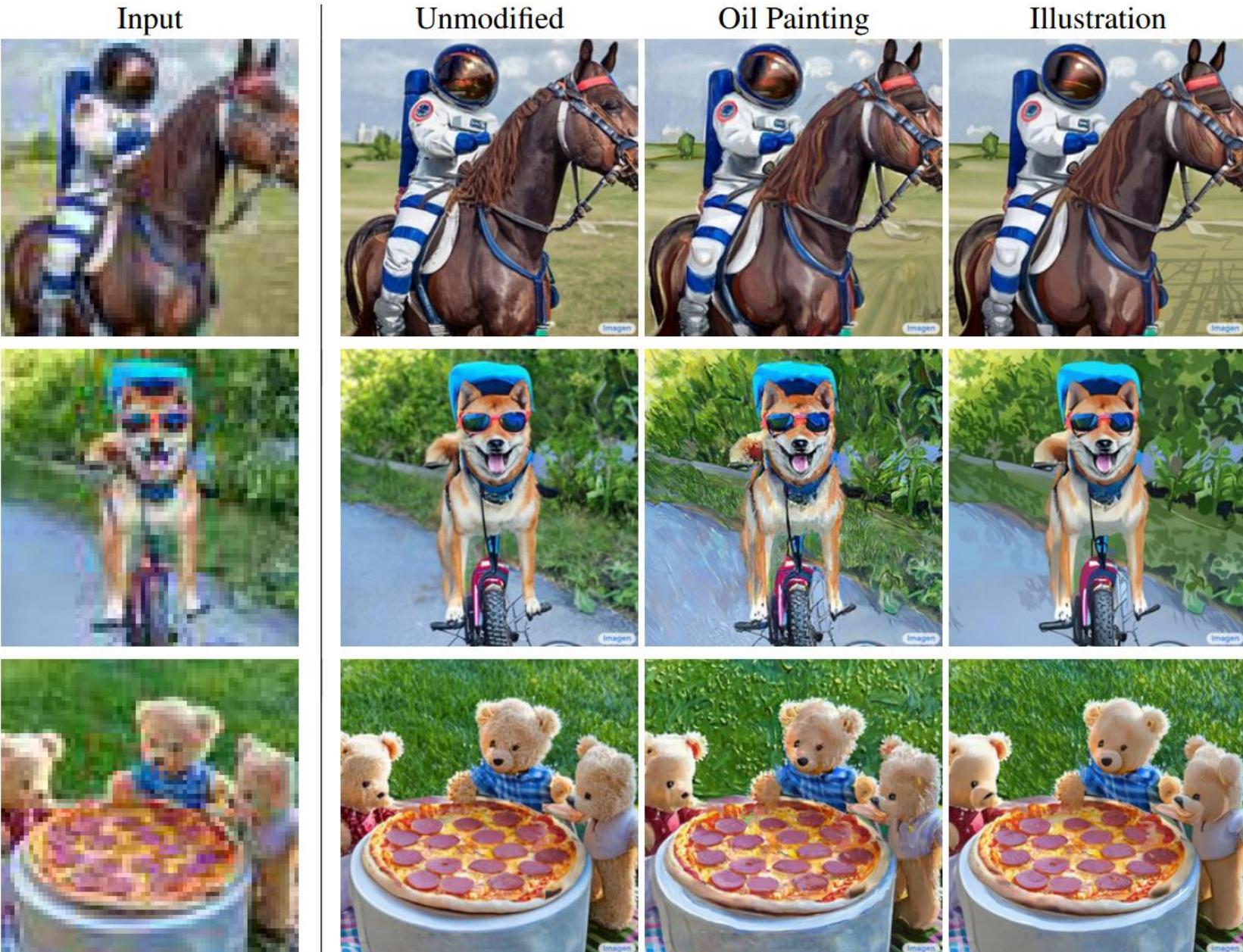


A small vessel propelled on water by oars, sails, or an engine.

# Imagen Styles

50

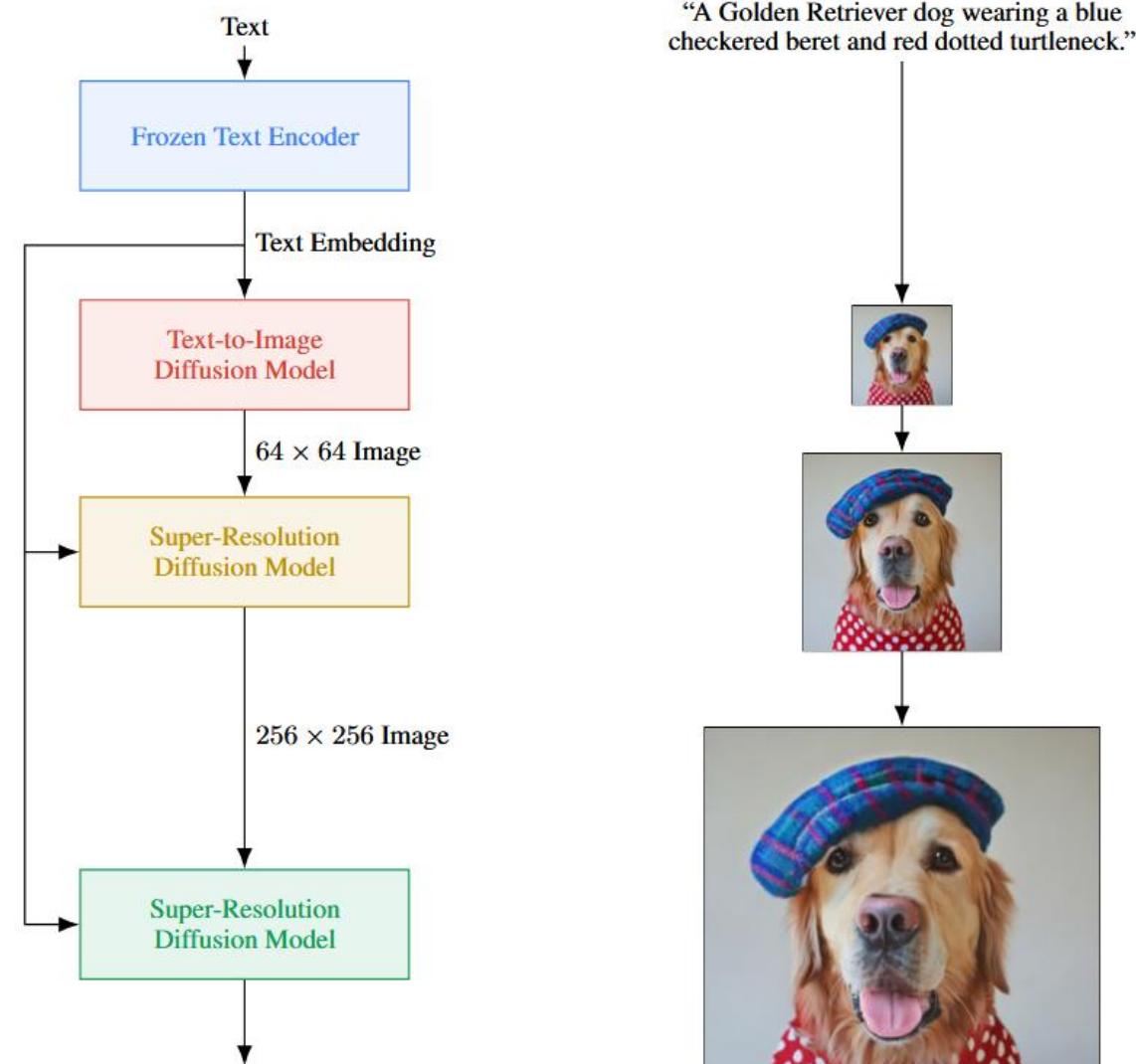
Add style to  
pixelated images



# Imagen Generation Process

51

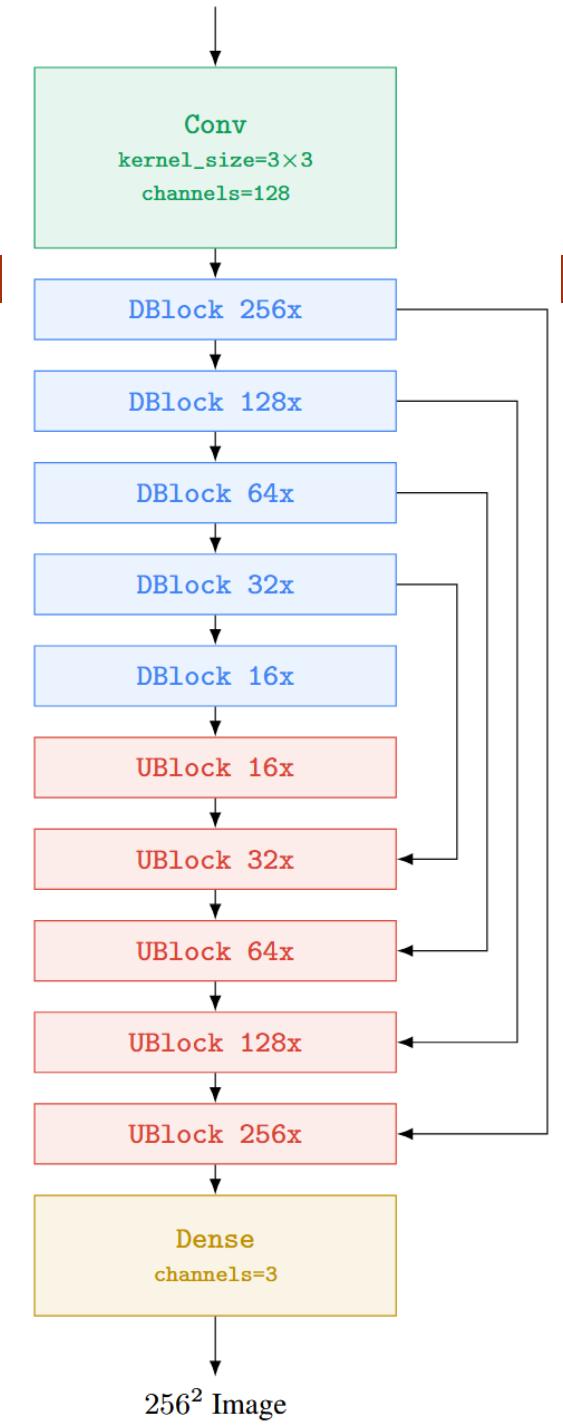
- Encode text prompt as an embedding vector with **frozen text encoder**
- Generate small **64x64 image** using **text-to-image diffusion model**
- Generate larger images using **super-resolution diffusion model**



# Imagen Architecture

52

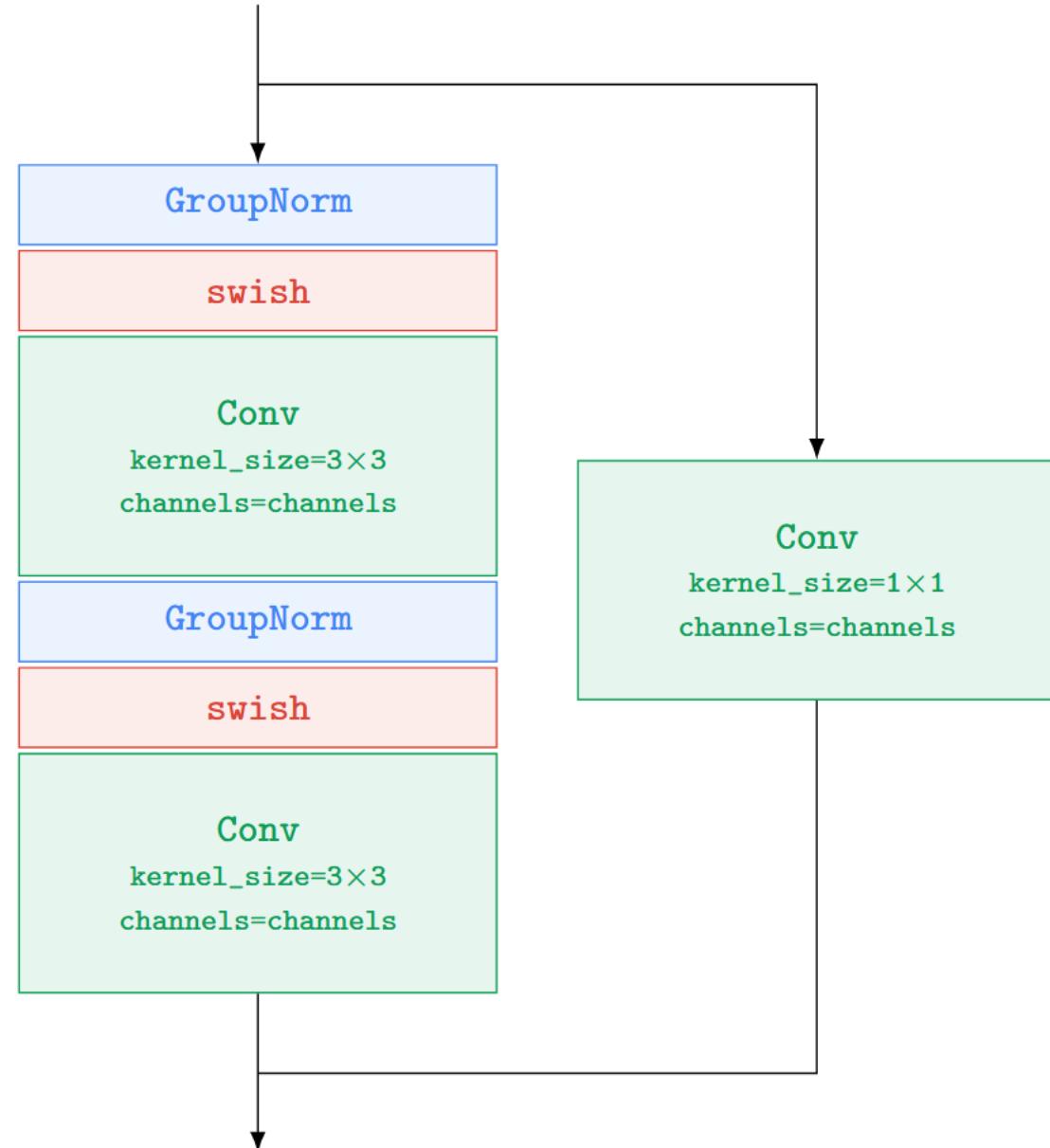
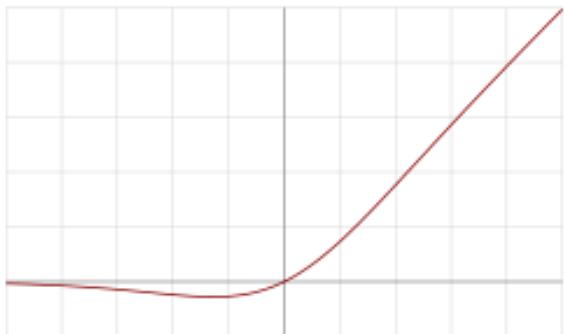
- Text-to-Image is based on U-Net
- Four pairs of down/up blocks
- Blocks have lower resolution, which speed up computation



# ResNet Blocks

53

- Down/up blocks are based on ResNet blocks
- Use group normalize, which is designed for small batches
- Use swish activation function



# Down Block

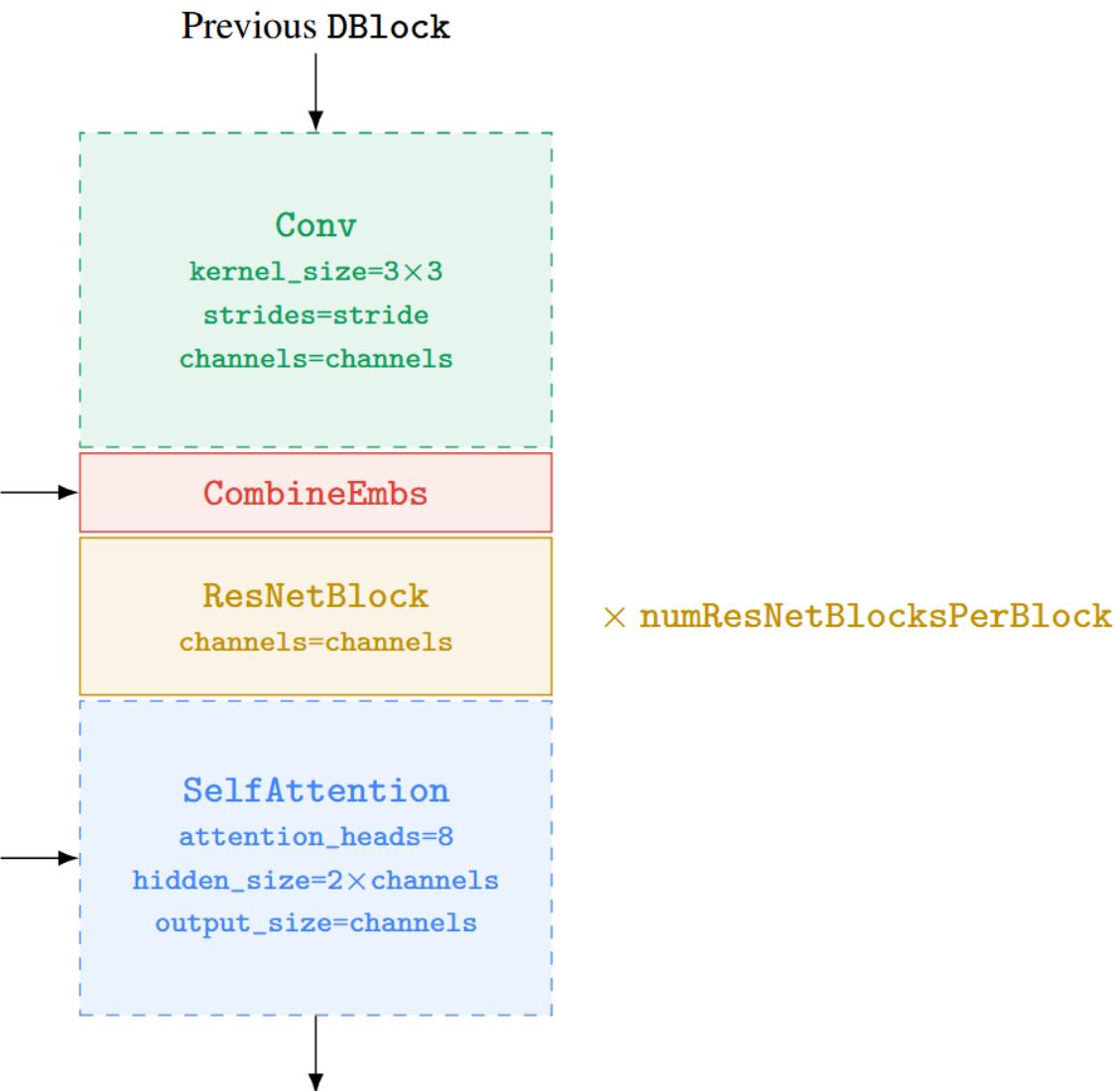
54

- **CombineEmbs** concatenate **time positional encoding** and **pooled text embedding**

- **Self attention on the “pixels”**

Conditional Embeddings  
(e.g., Time, Pooled Text Embeddings)

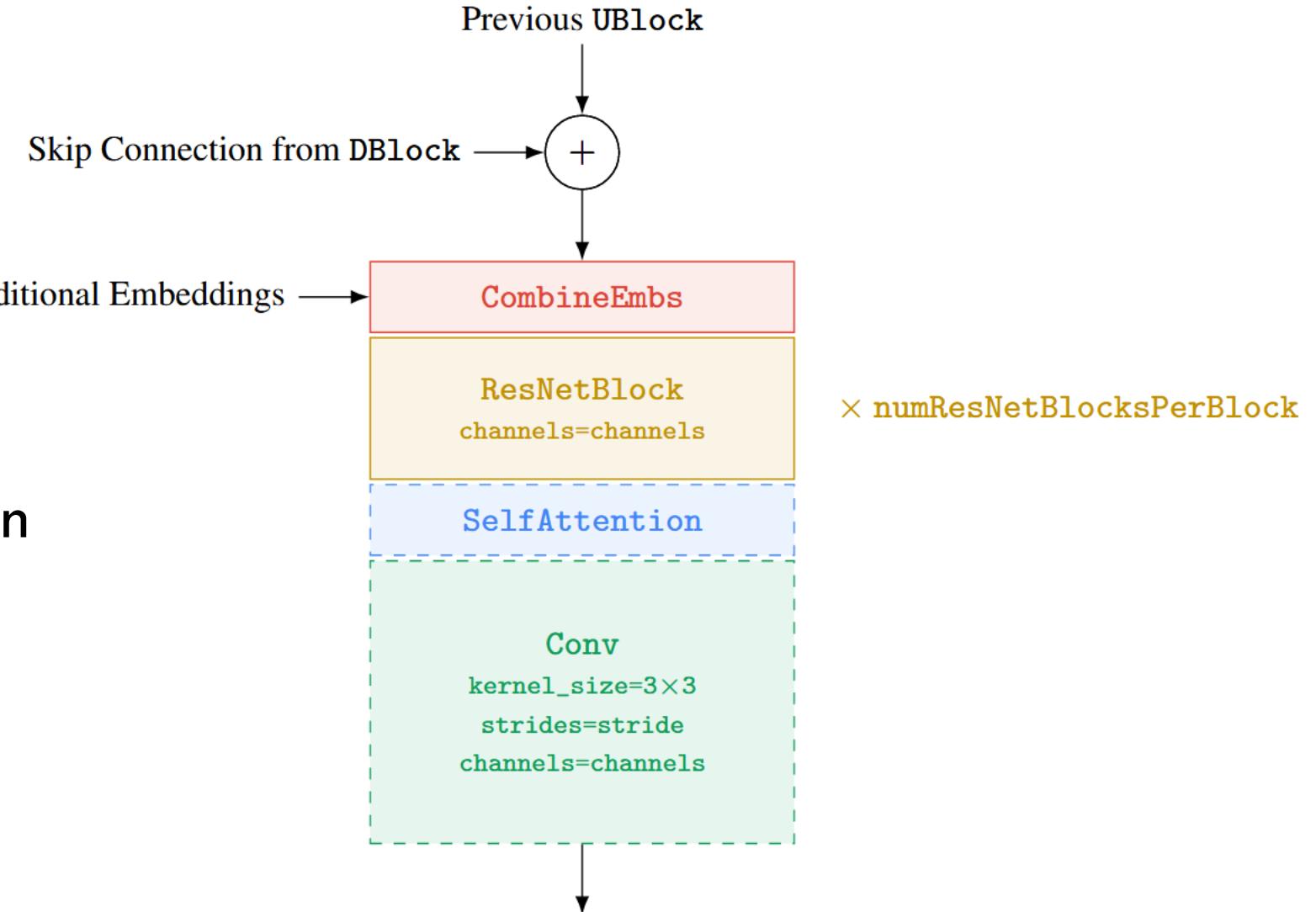
Full Contextual Text Embeddings →



# Up Block

55

- Similar to Down Block
- Includes skip connection



# Normalization Methods

56

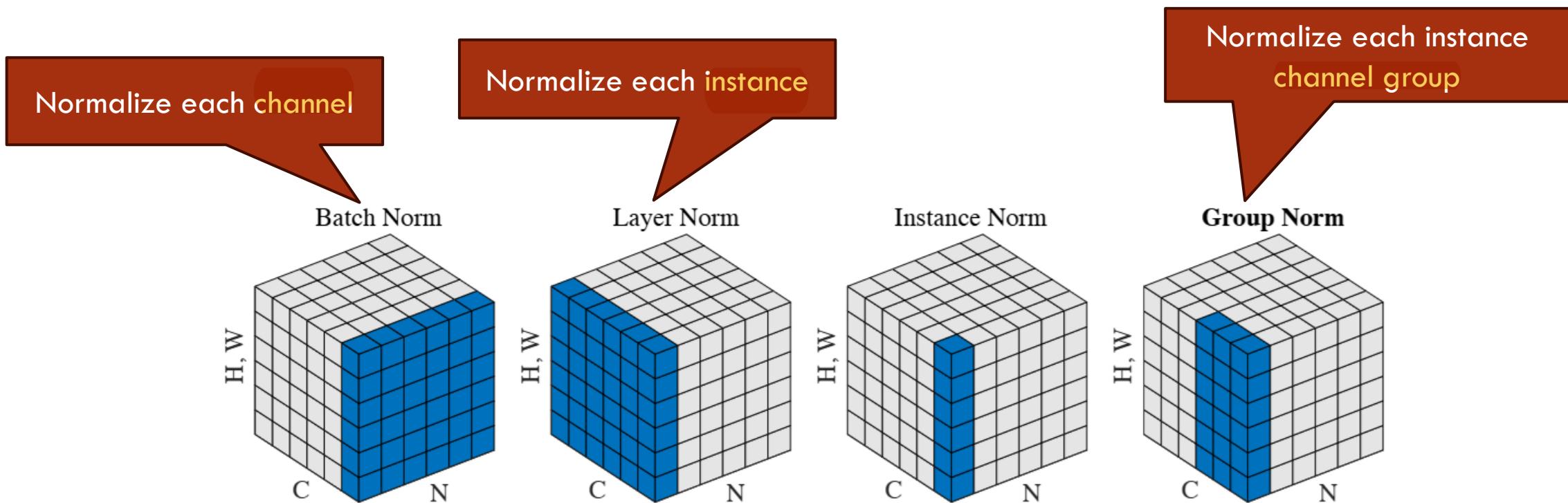


Figure 2. **Normalization methods.** Each subplot shows a feature map tensor, with  $N$  as the batch axis,  $C$  as the channel axis, and  $(H, W)$  as the spatial axes. The pixels in blue are normalized by the same mean and variance, computed by aggregating the values of these pixels.

# Normalization Method Comparison

57

Batch size: 32 Images

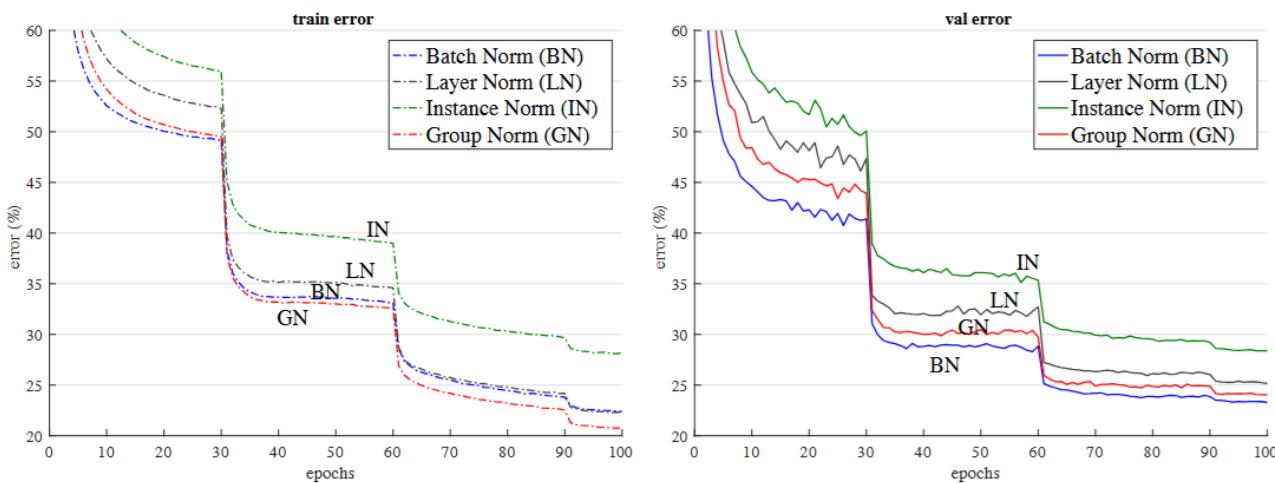


Figure 4. Comparison of error curves with a batch size of 32 images/GPU. We show the ImageNet training error (left) and validation error (right) vs. numbers of training epochs. The model is ResNet-50.

Decreasing batch size

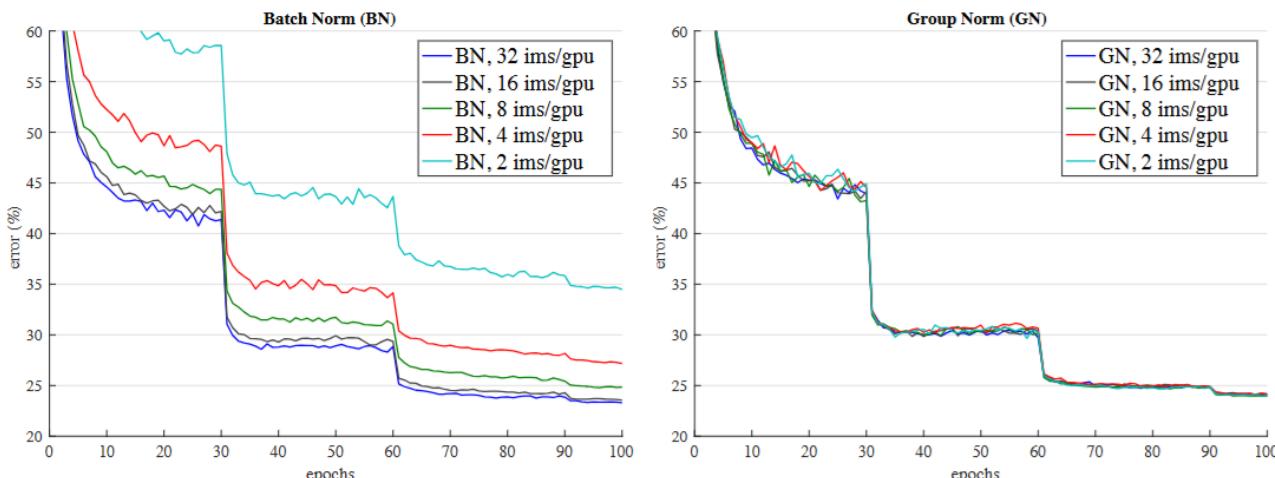


Figure 5. Sensitivity to batch sizes: ResNet-50's validation error of BN (left) and GN (right) trained with 32, 16, 8, 4, and 2 images/GPU.