

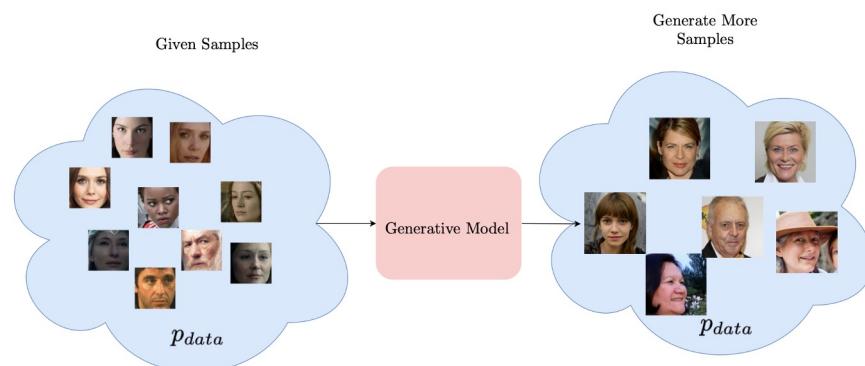
Generative Models

1

1

Generative Modelling

How does one sample from a distribution given only its samples?



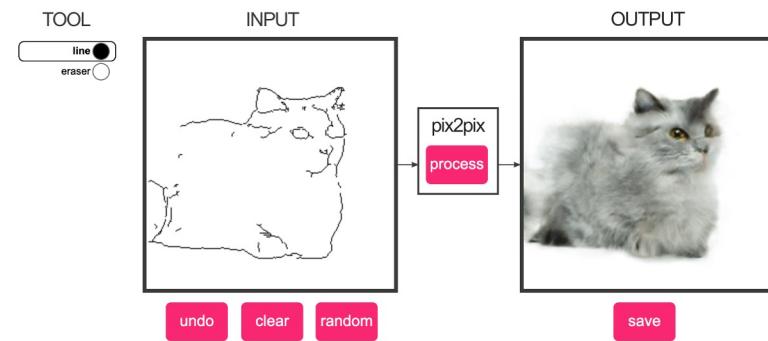
2

2

Generative Modelling: Why?

We can do some cool stuff! Like generate images from line drawings

edges2cats



Demo: <https://affinelayer.com/pixsrv/>

3

3

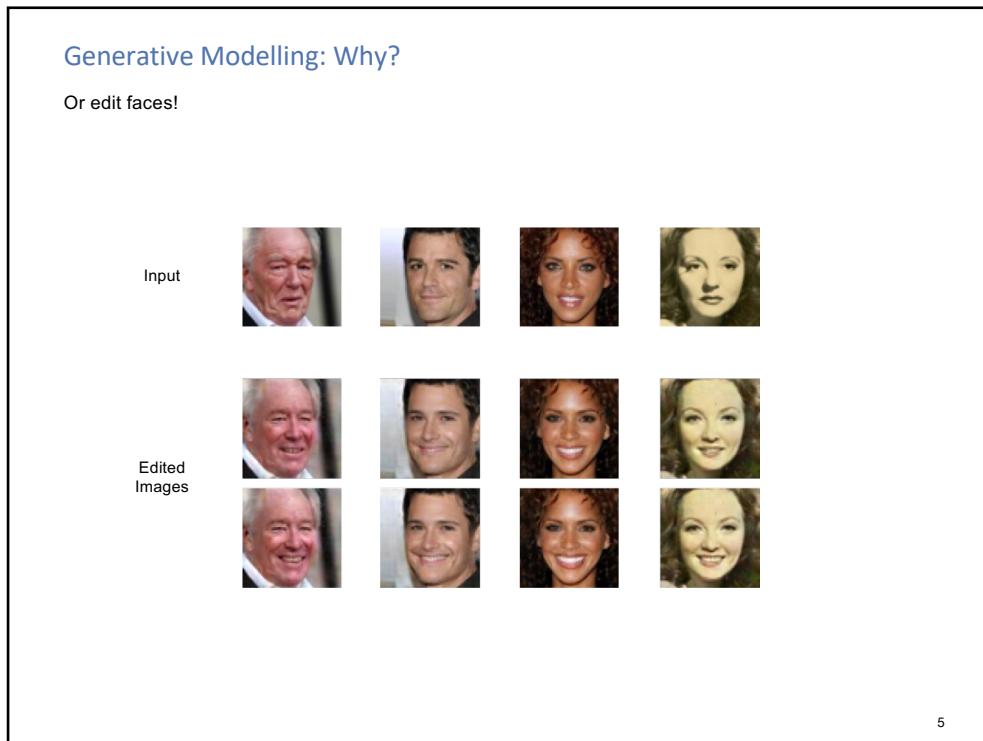
Generative Modelling: Why?

Change faces!

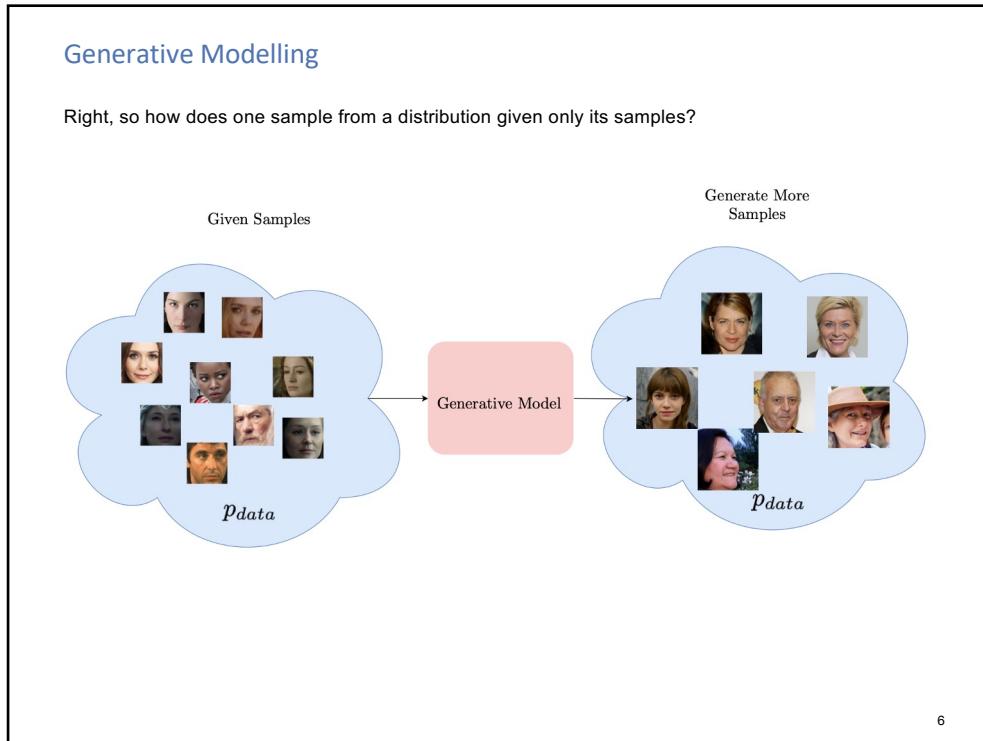


4

4



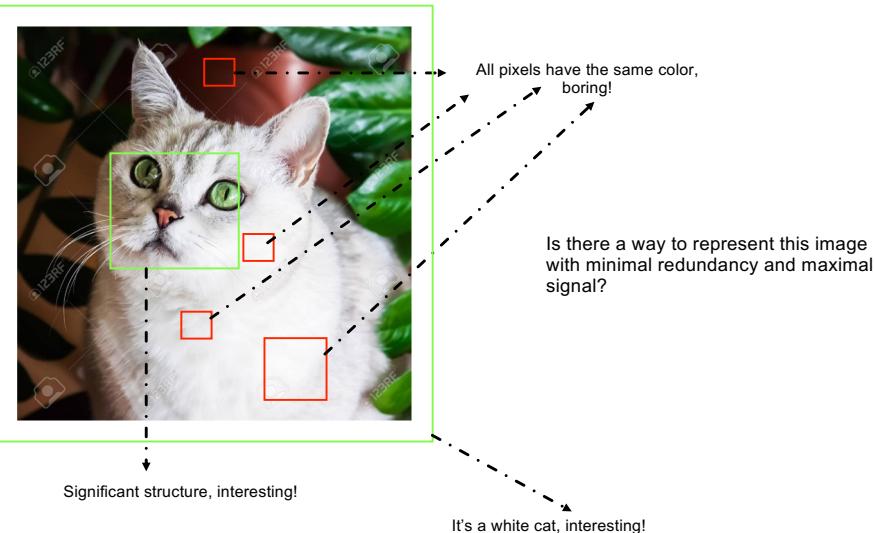
5



6

Generative Modeling: A Closer look at images

They're redundant, like very redundant



7

7

Generative Modeling: A Closer look at images

More realistically, is there a way to represent this image with minimal redundancy and a “good amount” of signal?

Yes! we do it all the time! Language!

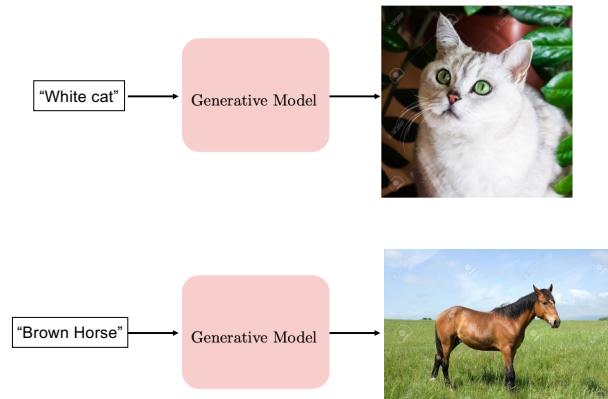


8

8

Generative Modeling: An ideal generative model

Say what you want the model to generate, and it generates it!



We can generate whatever we want!

9

9

Generative Modeling – From CSE527, 2020

"Unfortunately, life ain't that good and training such a model is hard. Especially because, text is quite a lossy representation" – Me, 2020



"White cat"

Eye Color? Head Shape? Number of whiskers?
Background?.....

10

10

Generative Modeling – From CSE527, 2020

"Unfortunately, life ain't that good and training such a model is hard. Especially because, text is quite a lossy representation" – Me, 2020



11

11

Generative Modeling – From CSE527, 2020

"Unfortunately, life ain't that good and training such a model is hard. Especially because, text is quite a lossy representation" – Me, 2020

Right.....



"A photorealistic image of a white cat on a garden ledge"

12

12

Generative Modeling

So, how do we get to DALL-E or Imagen? Lets start our journey...



"An astronaut riding a horse in a photorealistic style"



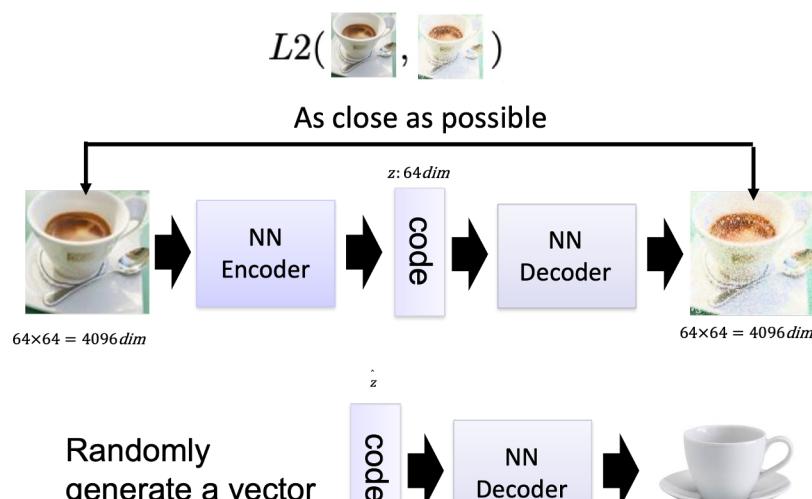
A photo of a Corgi dog riding a bike in Times Square. It is wearing sunglasses and a beach hat.

13

13

Generative Modeling: Autoencoders

Learn a representation by compressing an image as much as possible and then reconstructing it.

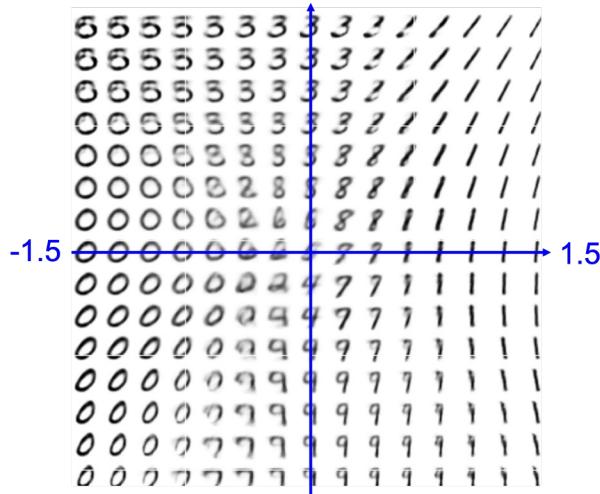


14

14

Generative Modeling: Autoencoders

But where do we sample from? A random vector may not produce a realistic image

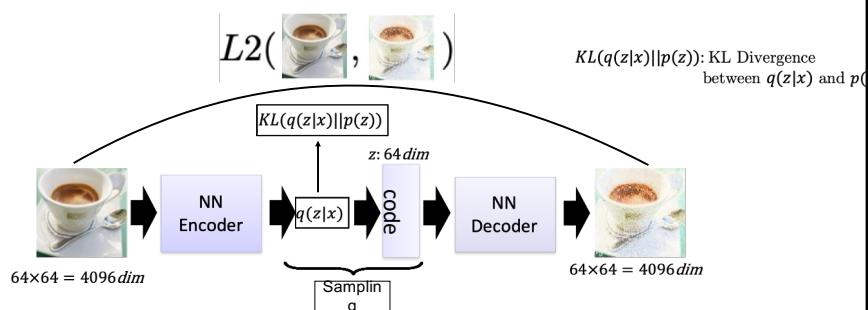


15

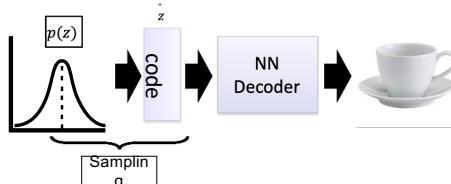
15

Generative Modeling: Variational Autoencoders

Learn a representation by compressing an image as much as possible and then reconstructing it, and push the distribution of the latent representation (i.e z) to a known prior $p(z) \sim N(0, I)$



Images can now be generated from a known prior distribution $p(z)$!



16

16

Generative Modeling: A diversion into divergences

A divergence allows one to measure the “distance” between two distributions $p(x)$ and $q(x)$. It is not a true distance measure because it may not be symmetric nor satisfy the triangle inequality. They satisfy the following properties:

- $D(p(x)||q(x)) \geq 0 \forall p(x), q(x)$
- $D(p(x)||q(x)) = 0 \text{ if and only if } p(x) = q(x)$

Therefore, even if not a true distance, a divergence can allow us to measure how similar or dissimilar two distributions $p(x)$ and $q(x)$ are!

17

17

Generative Modeling: The KL and Jenson-Shanon Divergence

A divergence we've already seen is the KL-Divergence or Kullback-Leibler divergence which is defined as follows:

$$KL(p(x)||q(x)) = \int_{-\infty}^{\infty} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx$$

Notice $KL(p(x)||q(x)) \neq KL(q(x)||p(x))$!

[This animation should give you some intuition about how the KL-Divergence works!](#)

The Jenson-Shanon Divergence (JSD) is a symmetric version of the KL-Divergence and is defined as follows

$$JSD(p(x)||q(x)) = \frac{1}{2}KL\left(p(x)||\frac{p(x) + q(x)}{2}\right) + \frac{1}{2}KL\left(q(x)||\frac{p(x) + q(x)}{2}\right)$$

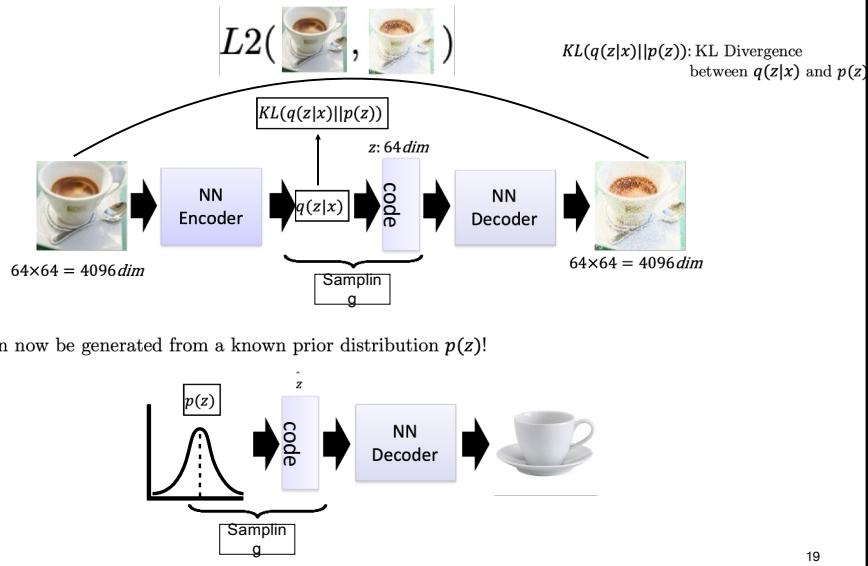
We'll be running into JSD when discussing GANs 😊

18

18

Generative Modeling: Variational Autoencoders

Learn a representation by compressing an image as much as possible and then reconstructing it, and push the distribution of the latent representation (i.e z) to a known prior $p(z) \sim N(0, I)$



Generative Modeling: Variational Autoencoders

Unfortunately, the L2 loss does not capture realism

$$L2(\boxed{7}, \boxed{7}) = L2(\boxed{7}, \boxed{7})$$



More realistic

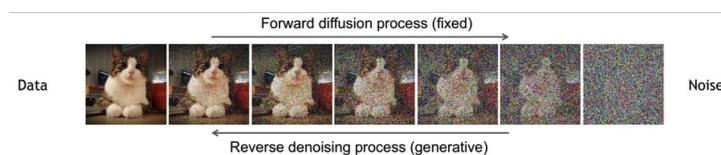


Less realistic

20

Generative Modeling: Denoising Diffusion models

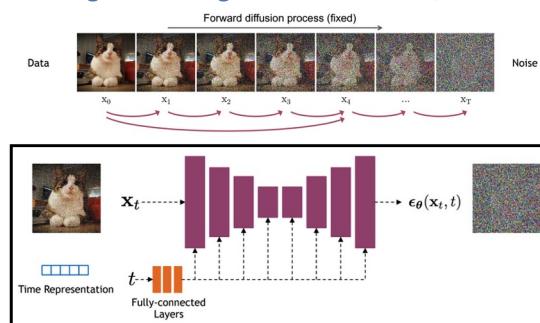
- A Denoising Diffusion model consists of two processes:
- A Forward diffusion process that gradually adds noise to the input
 - A Reverse denoising process that learns to generate data by denoising
 - The objective this minimizes is very similar to that of a VAE i.e a variational lower bound



Slides are mostly from the workshop on Denoising Diffusion models, CVPR2022: <https://drive.google.com/file/d/1DYHDh11SI9oam3O333biRYzSC0Idtmn/view>

21

Generative Modeling: Denoising Diffusion models, The forward process



Algorithm 1 Training

- Use a U-Net to predict the noise at each step
- Push the noise to be as close to gaussian as possible
- $$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$$
- $$\boldsymbol{\epsilon} \sim N(0, \mathbf{I})$$
- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on $\nabla_\theta \|\boldsymbol{\epsilon} - \epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$
 - 6: **until** converged

Slides are mostly from the workshop on Denoising Diffusion models, CVPR2022: <https://drive.google.com/file/d/1DYHDh11SI9oam3O333biRYzSC0Idtmn/view>

22

Generative Modeling: Denoising Diffusion models, The forward process

Algorithm 2 Sampling

Use the predicted noise by U-Net to invert the image

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$$

$$\epsilon \sim N(0, \mathbf{I})$$

```

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

```

Slides are mostly from the workshop on Denoising Diffusion models, CVPR2022: <https://drive.google.com/file/d/1DYHDh1ISl9oam3O333biRYzSC0ldtmn/view>

23

23

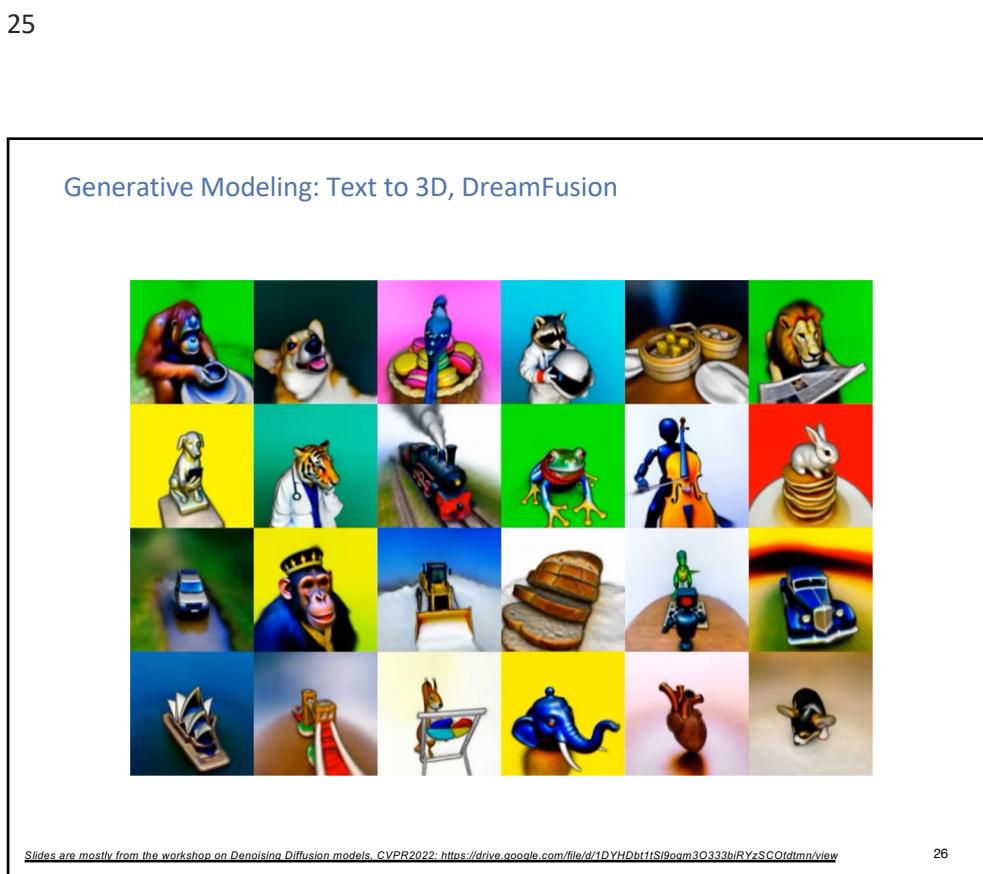
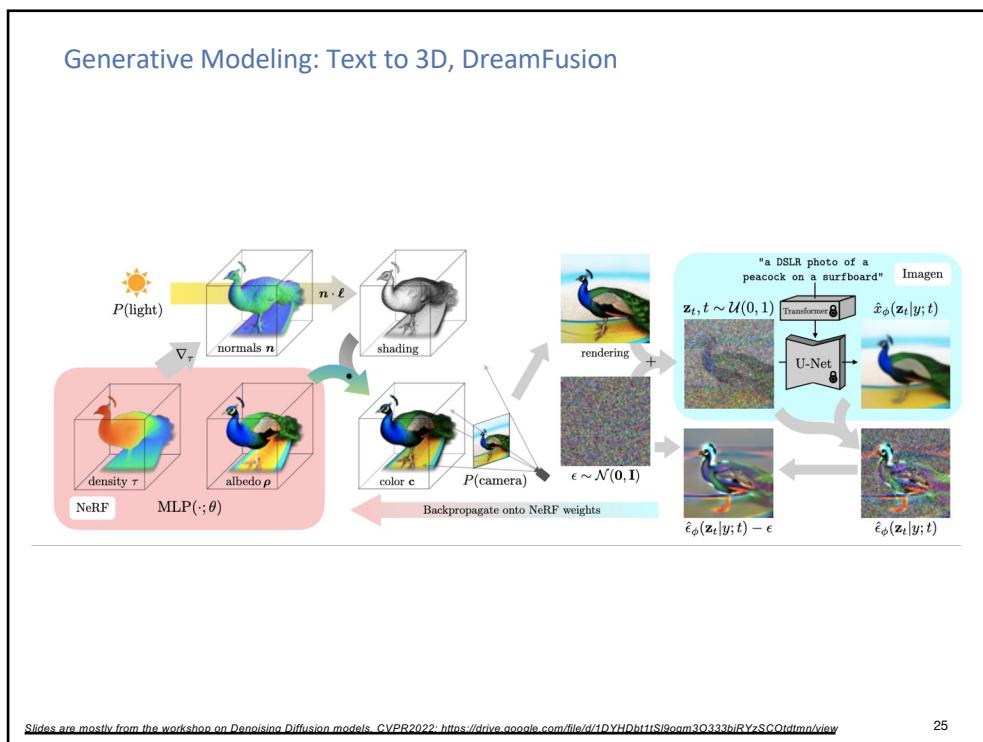
Generative Modeling: DALL-E 2

DALL-E 2

Slides are mostly from the workshop on Denoising Diffusion models, CVPR2022: <https://drive.google.com/file/d/1DYHDh1ISl9oam3O333biRYzSC0ldtmn/view>

24

12



Generative Modeling: Some Resources

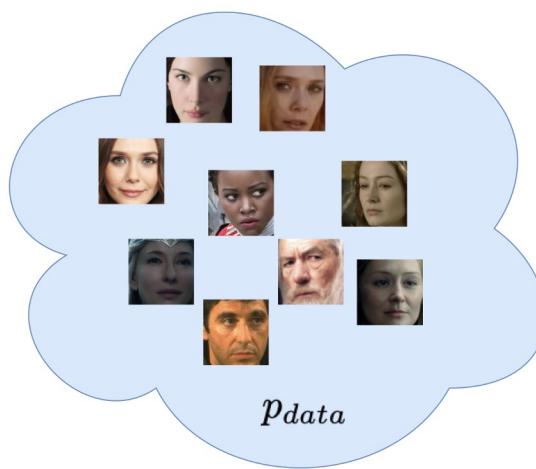
- The Deep Learning Book discusses Autoencoders along with many other generative models: <https://www.deeplearningbook.org/>
- Durk Kingma's thesis is a great read for probabilistic generative models and VAEs: <https://pure.uva.nl/ws/files/17891313/Thesis.pdf>
- Lilian Weng's blog post on the VAE is great! <https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html>
- Diffusion models: <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>
- Diffusion Models: <https://cvpr2022-tutorial-diffusion-models.github.io/>

27

27

Generative Adversarial Networks (Goodfellow et al. NIPS2014)

How does one sample from a distribution given only its samples?

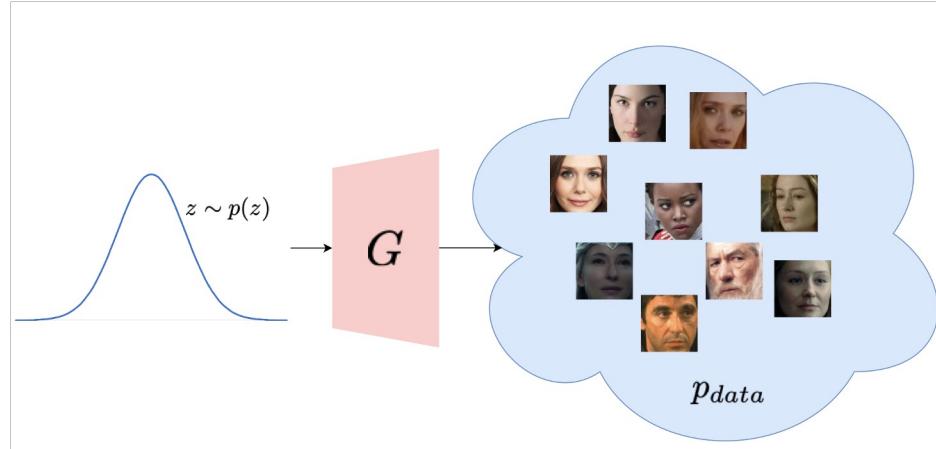


28

28

Generative Adversarial Networks

Train a “Generator” to transform samples from a gaussian to the image space!

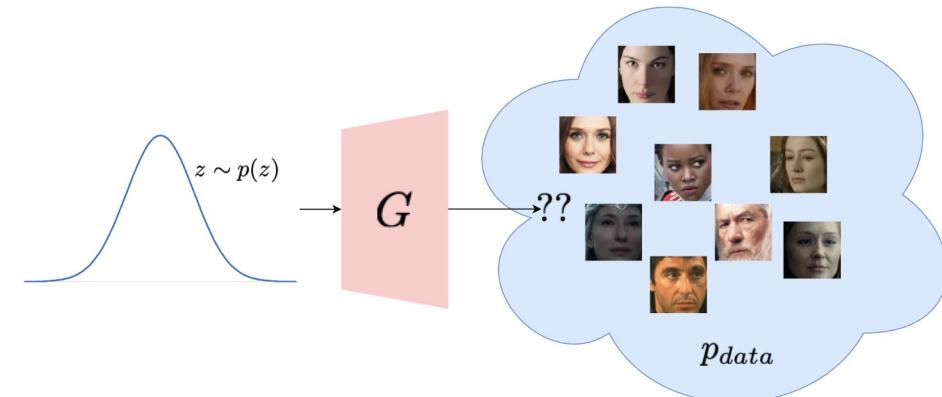


29

29

Generative Adversarial Networks

- But we don't know what a random z maps to in the Image space, in other words we do not have $(z, Image)$ pairs.
- Supervised losses (L2, L1 etc) cannot be used!

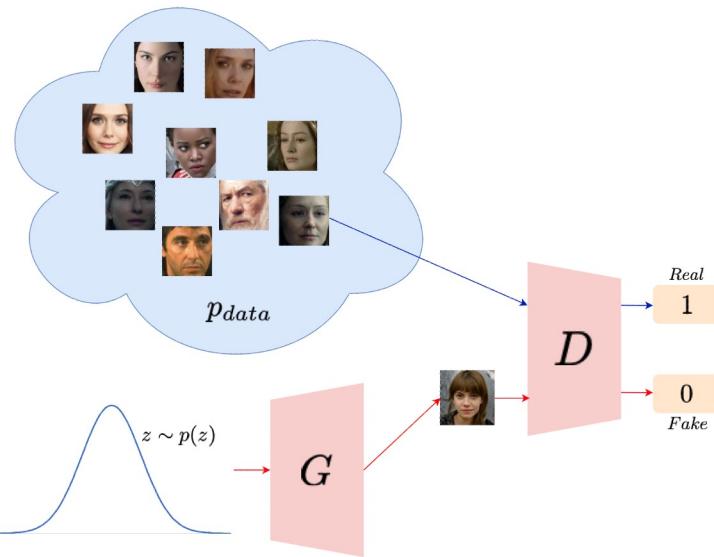


30

30

Generative Adversarial Networks

Use a “Discriminator” to distinguish between samples generated by G and the real samples.

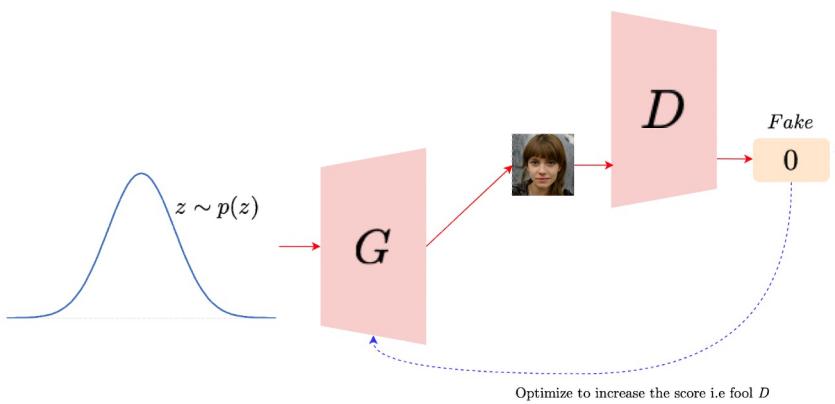


31

31

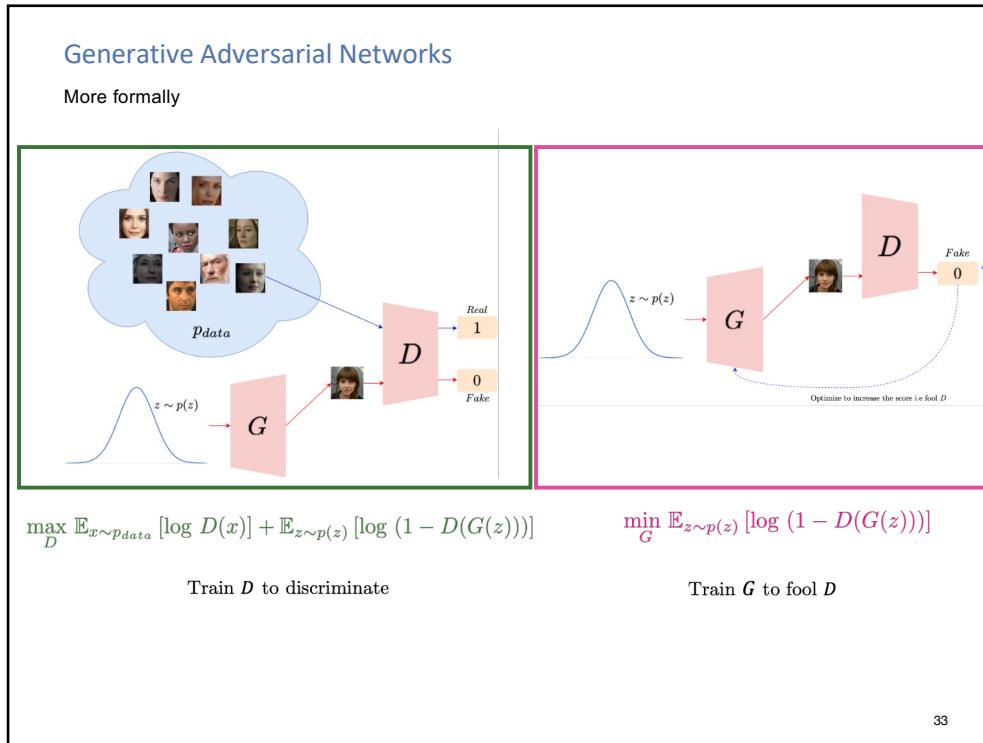
Generative Adversarial Networks

Train G to maximize D 's score. Use the gradient as a “guide” to generate realistic images.

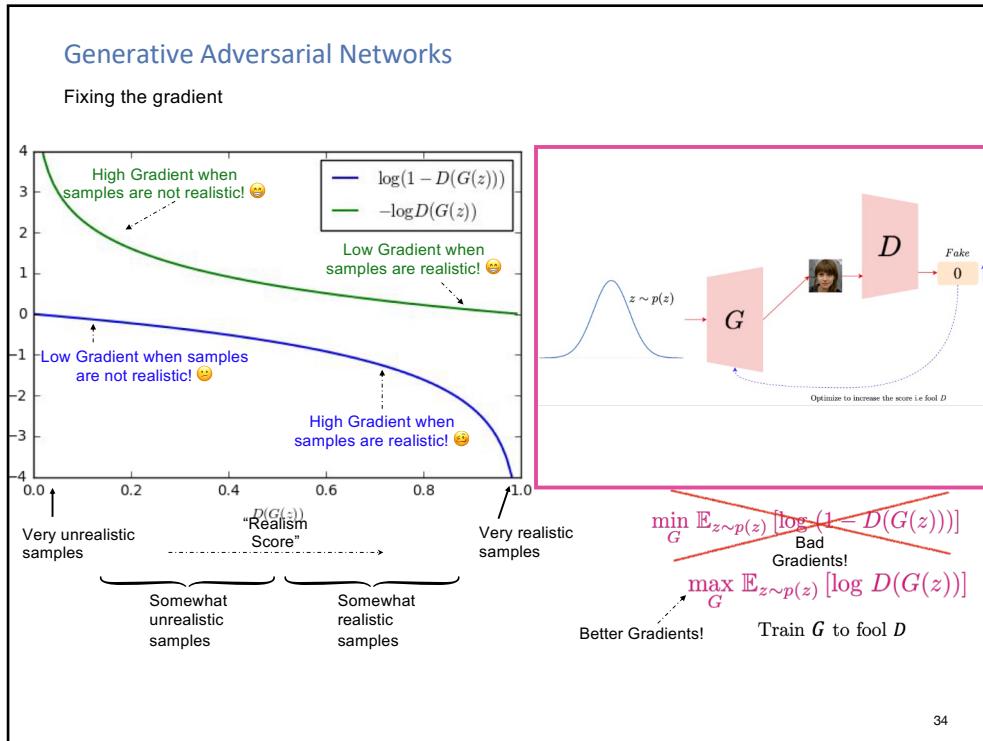


32

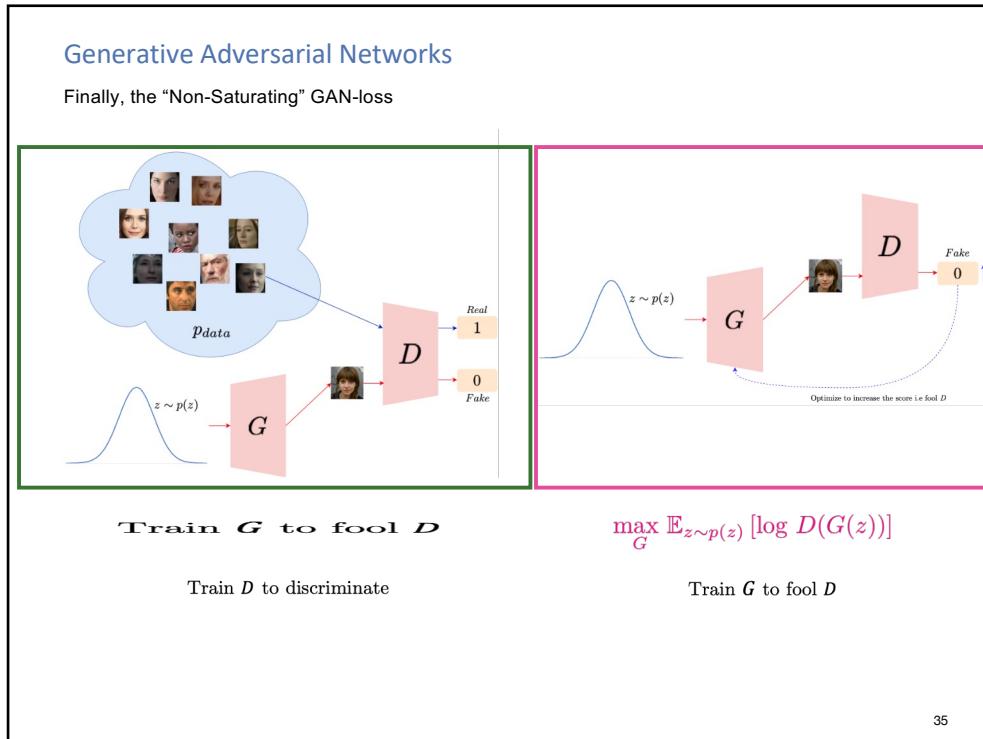
32



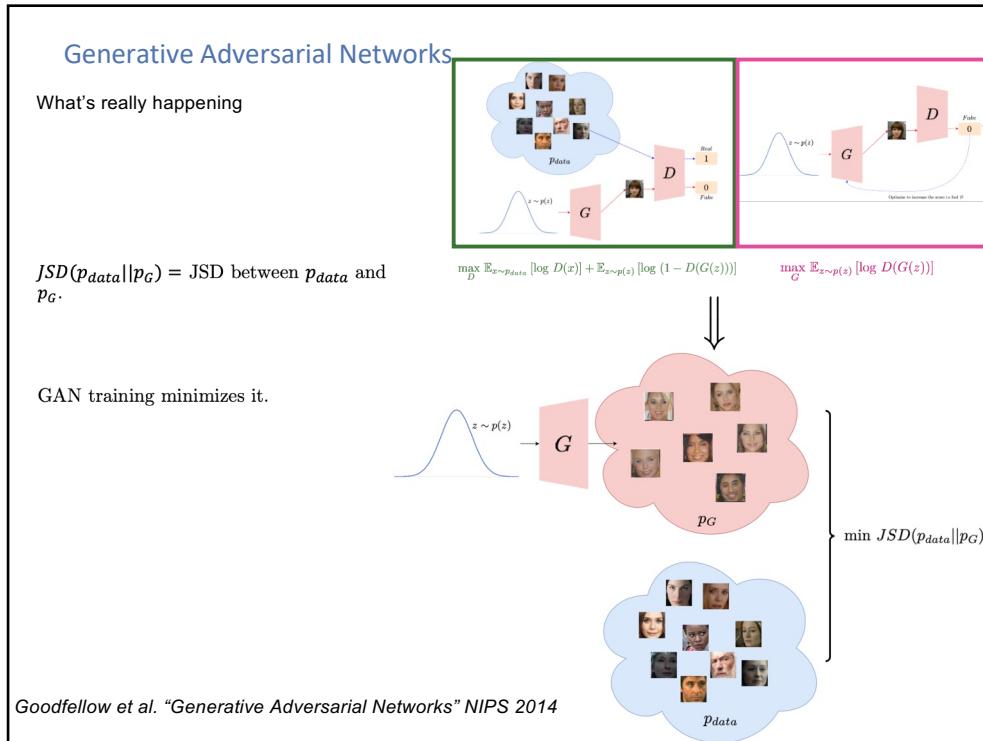
33



34



35



36

Generative Adversarial Networks

GAN samples have become really good. Samples from StyleGAN V2 (Karras et al. CVPR2020)

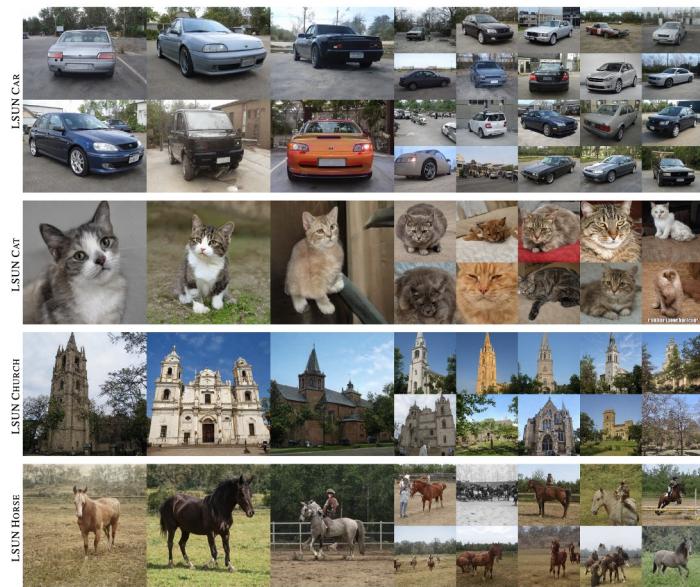


37

37

Generative Adversarial Networks

GAN samples have become really good. Samples from StyleGAN V2 (Karras et al. CVPR2020)



38

38

Generative Adversarial Networks

GAN samples have become really good. Samples from BigGAN (Brock et al. ICLR2019)



39

39

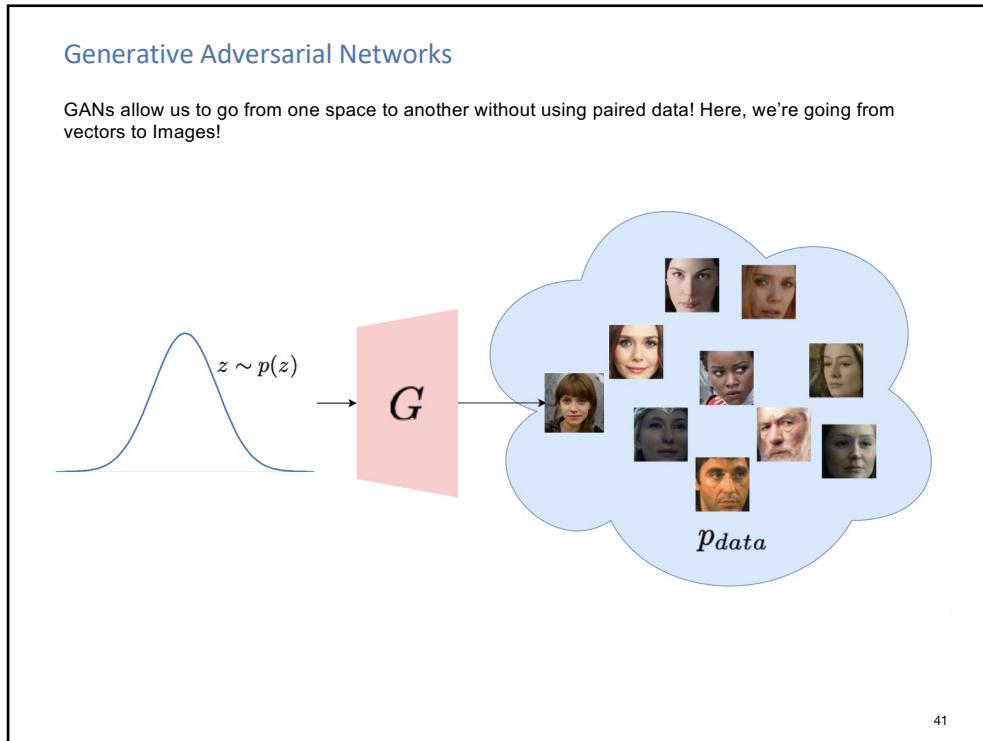
Generative Adversarial Networks: What can they be used for?

Let's see how GANs can help us for two standard CV problems: *Denoising* and *Super-Resolution*



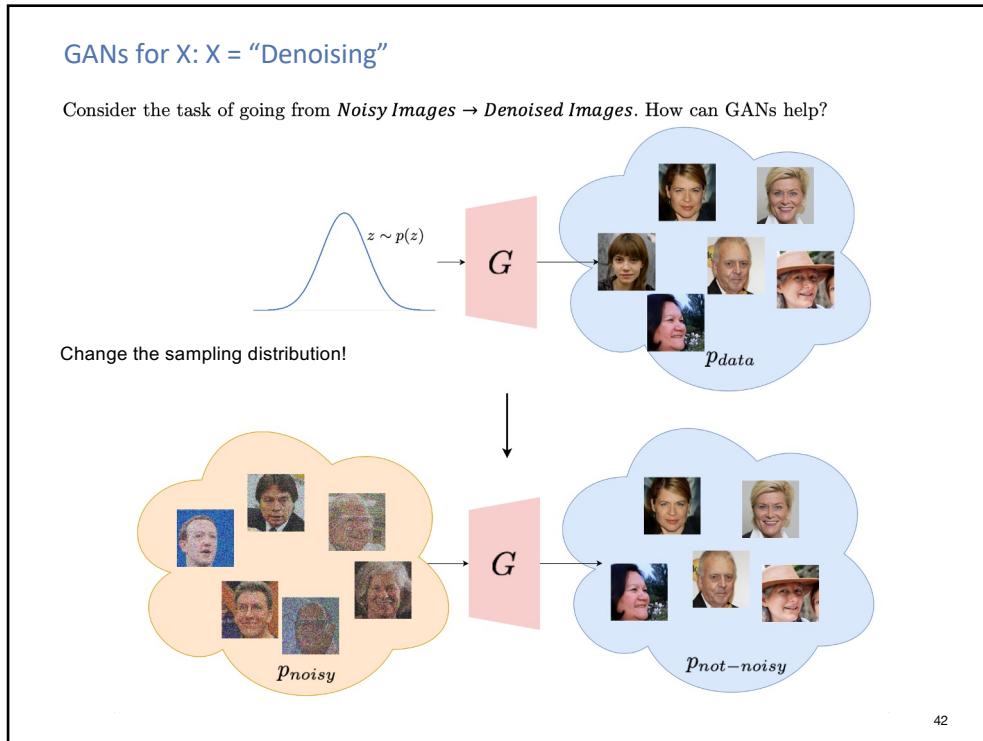
40

40



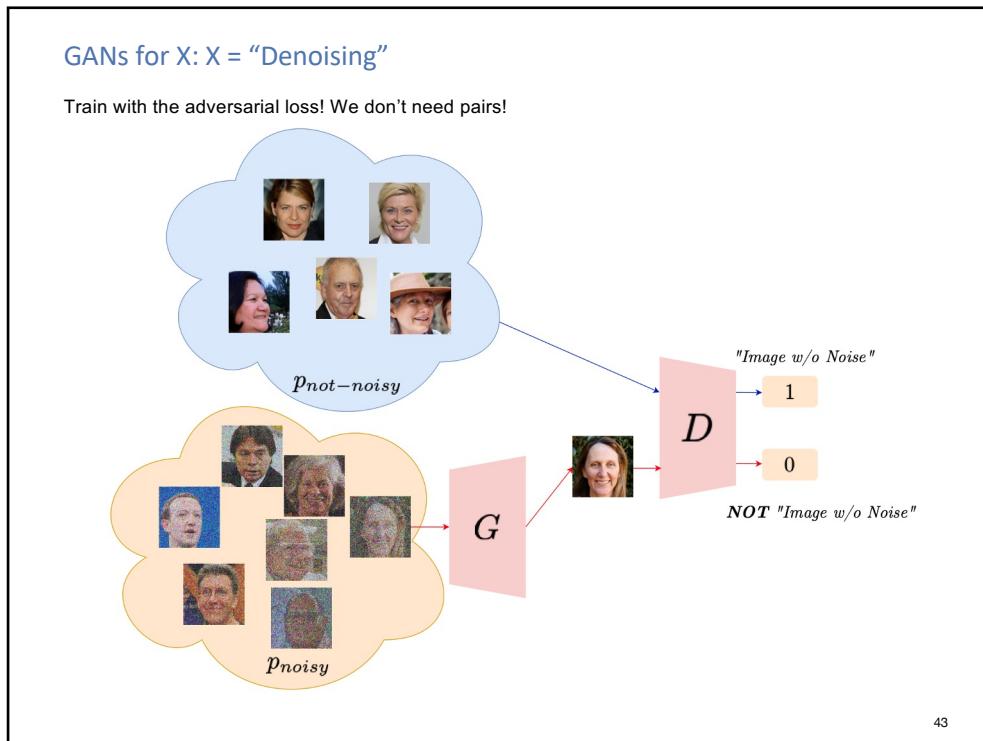
41

41

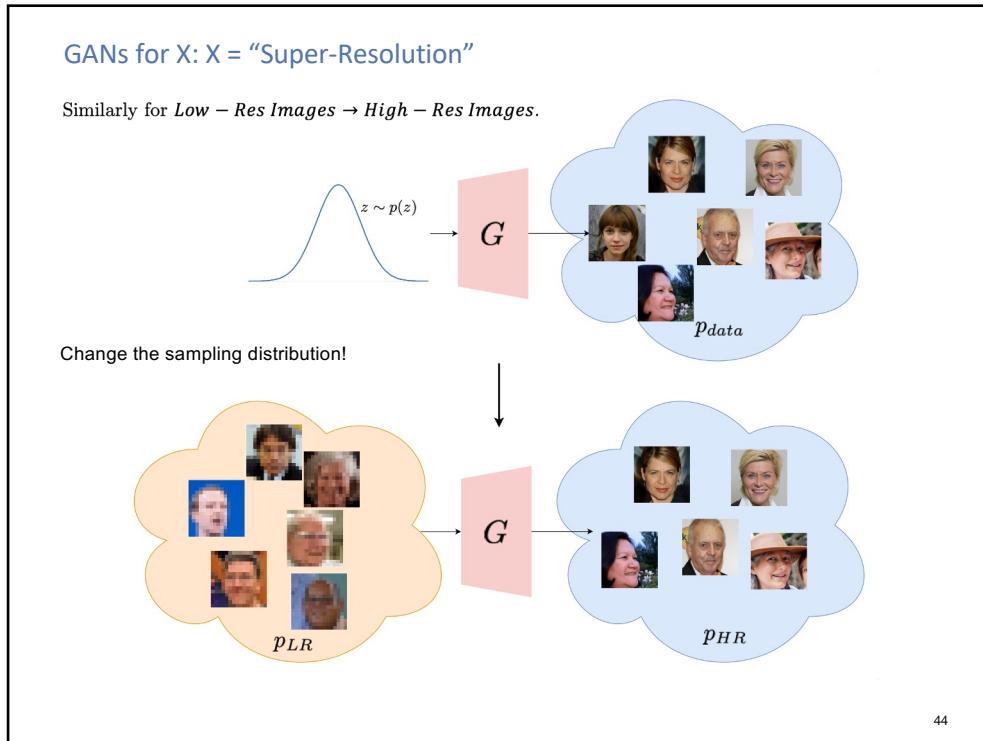


42

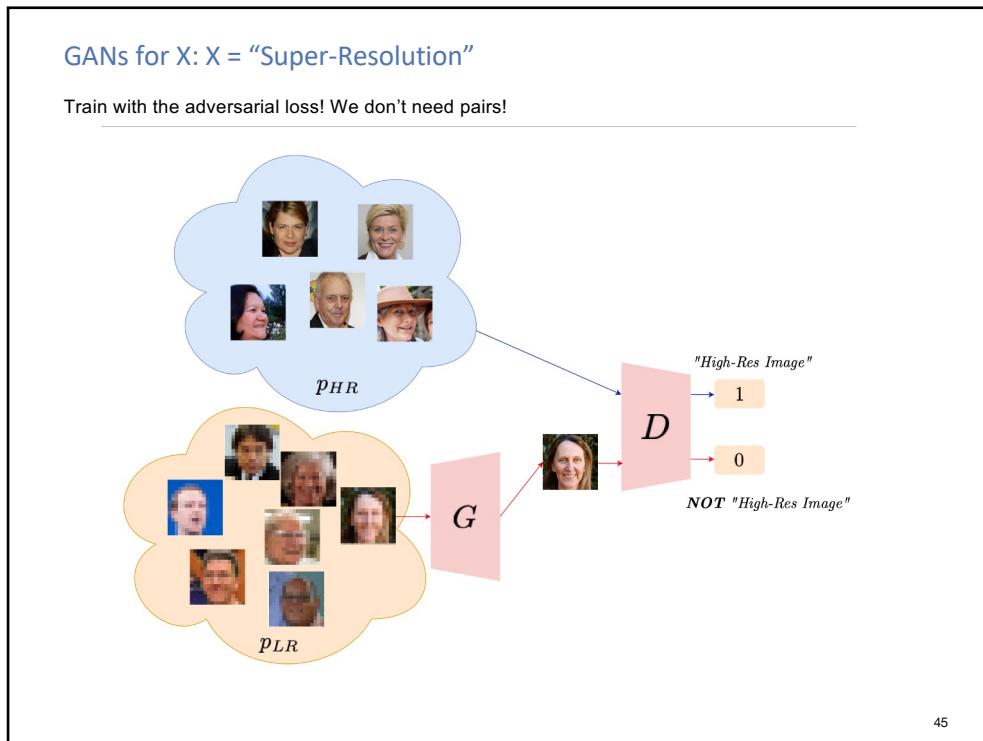
42



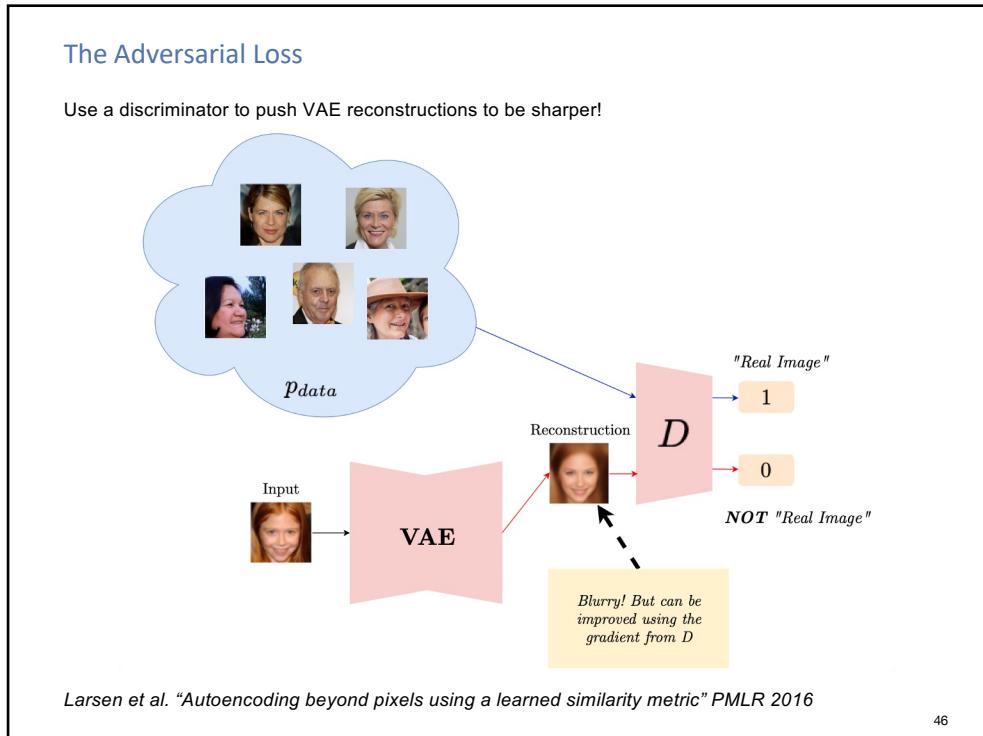
43



44



45



46

The Adversarial Loss

Use a discriminator to push VAE reconstructions to be sharper!



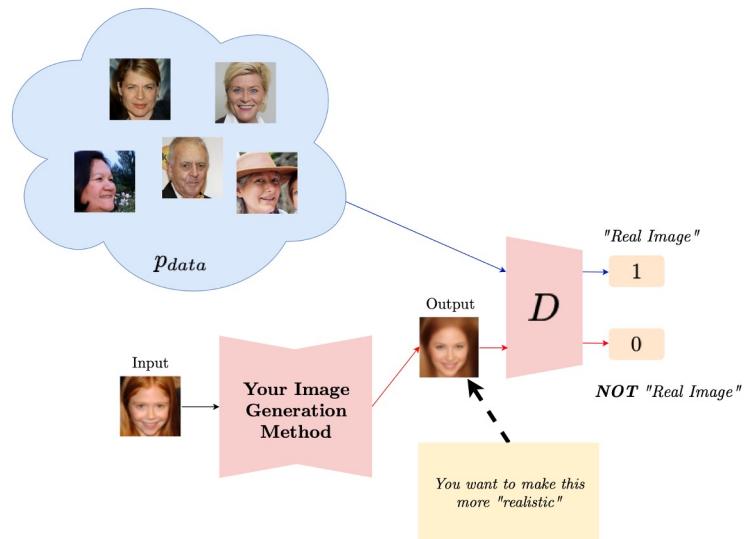
Larsen et al. "Autoencoding beyond pixels using a learned similarity metric" PMLR 2016

47

47

The Adversarial Loss

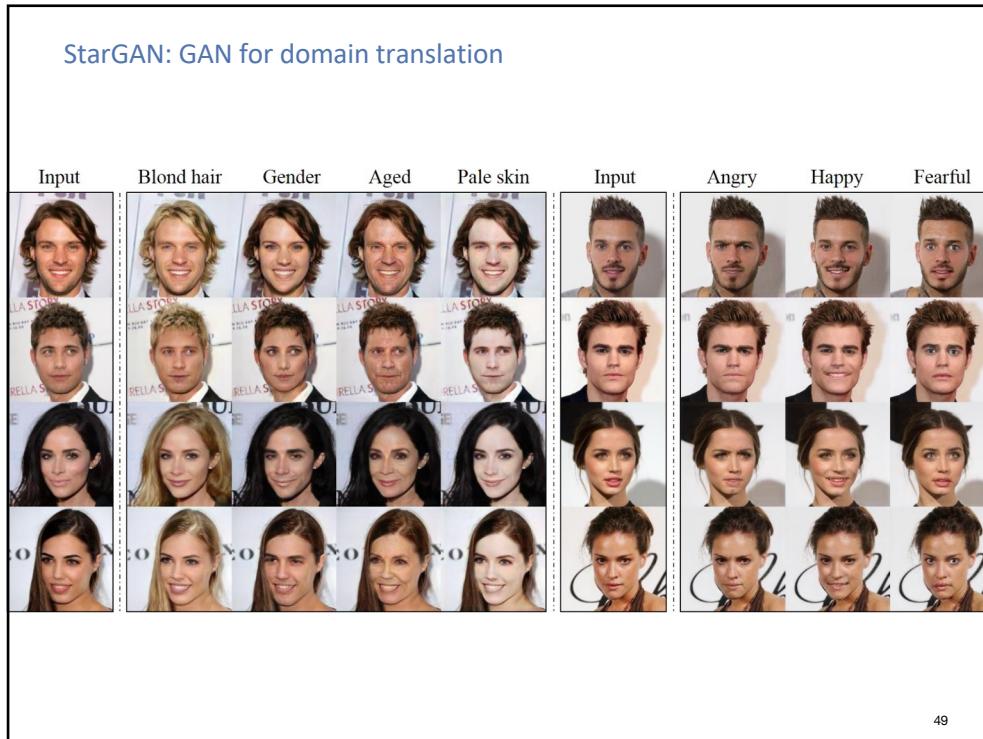
A discriminator can be used to improve to results of any image generation method!



You et al. "Your Paper Title" WhereverYouSubmit 2XXX

48

48



49

Generative Adversarial Networks: Losses Galore!

The Standard GAN Loss can be pretty unstable to train with. Consequently, multiple losses and regularizers have been proposed

- WGAN-GP (Gulrajani et al. 2017)

$$\underbrace{\mathbb{E}_{\hat{x} \sim \mathbb{P}_g} [D(\hat{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Our gradient penalty}}.$$

- Least-Squares GAN (Mao et al. 2016)

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D(x) - 1)^2] + \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)))^2]$$

$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{z \sim p_z(z)} [(D(G(z)) - 1)^2].$$

- Spectral Normalization (Miyato et al. 2018)
- R1 Regularization (Mescheder et al. 2018)

50

50

Generative Adversarial Networks: Architectures

Architectures (including activations, normalizations etc) of both the Generator and the Discriminator are also crucial to generate good samples. Some important work is given below

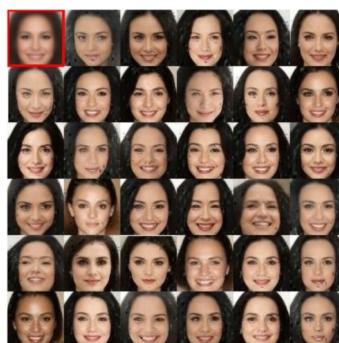
- DCGAN (Radford et al. 2016): First Convolutional Generator
- Pix2pix (Isola et al. 2017): Introduced a patch-based discriminator that discriminated on overlapping patches of the input image.
- Progressive GAN (Karras et al. 2017): First Mega-Pixel level image generation using GANs. Proposed Progressive Training
- CycleGAN (Zhu et al 2017): GANs for image-to-image translation from one domain to another
- StarGAN (Choi et al. 2017): GANs for image-to-image translation between *multiple* domains
- Spectral Normalization (Miyato et al. 2018): First good results on ImageNet classes
- StyleGAN (Karras et al. 2018): High Quality Mega-Pixel level image generation using GANs. Doesn't require progressive training. Introduced a much richer "style" latent space.
- BigGAN (Brock et al. 2019): HQ results in ImageNet Classes

51

51

Generative Adversarial Networks: The problems

- As alluded to earlier, training can be unstable
- Mode Collapse can occur (example from Richardson et al, "On GANs and GMMs", 2018). Faces are very similar!



- Given a sample x , GANs can't tell you how likely (realistic) it is i.e it cannot give you $p_{data}(x)$. The discriminator score cannot be used for this, check out *Goodfellow et al 2014* (the original GAN paper) for why.

52

52

Generative Modeling: Some Resources

- The Deep Learning Book discusses GANs along with many other generative models: [Link](#)
- [The original GAN paper is a great read!](#)
- [Scott Rome's blog digs into the details the the proofs in the GAN paper](#)
- [Lilian Weng's blog post on GANs is great too!](#)

53

53

Other Methods are catching up

VQ-VAE-2 (Razavi et al. 2019): Uses a autoregressive prior on the VAE latent space and a much deeper encoder and decoder.



1024x1024

54

54

Other Methods are catching up

GLOW (Kingma et al. 2018): Based on Normalizing Flows (Dinh et al, *NICE: Non-linear Independent Components Estimation*, 2015). Methods based on Normalizing Flows give exact likelihoods $p_{data}(x)$!



256×256

55

55