

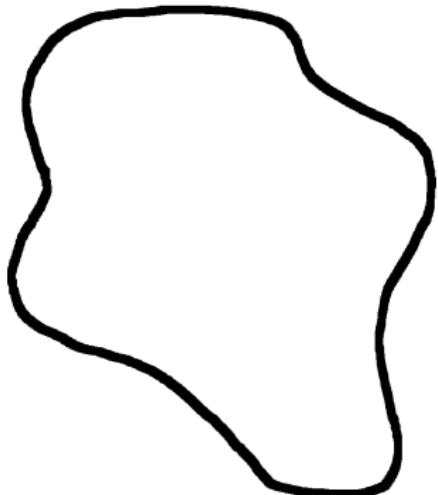
UO 607: Generative Models

Steven Walton
University of Oregon

23 October 2020

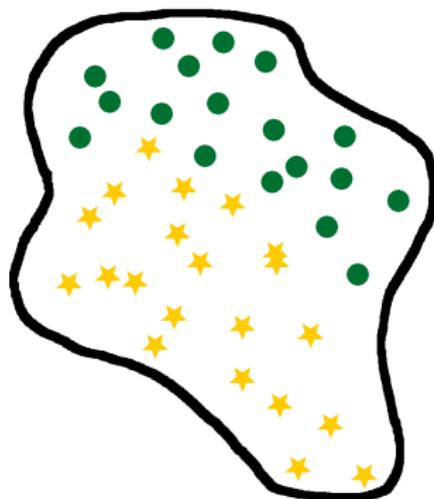
Generators

- ▶ Suppose that we have some probability space S .



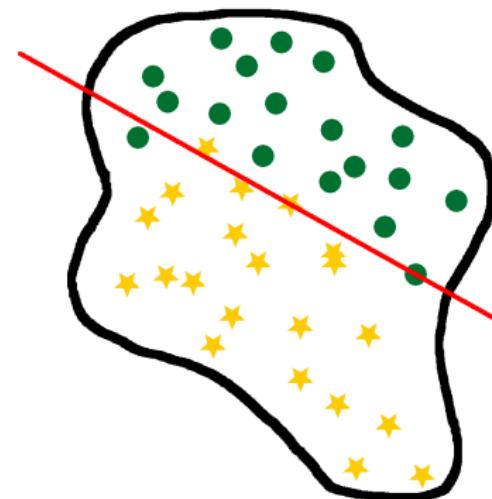
Generators

- ▶ Suppose that we have some probability space S .
- ▶ S has two variables X (instances) and Y (labels) that represent cats and dogs.



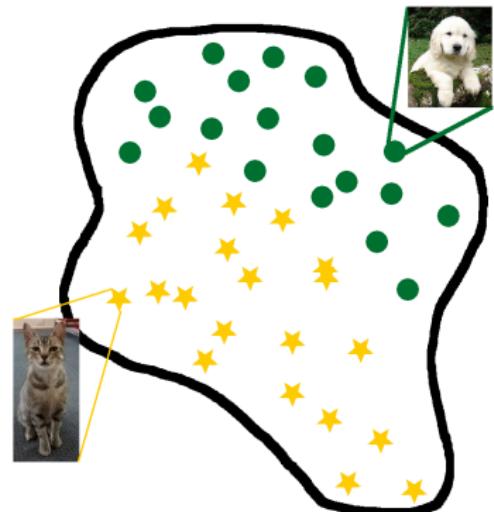
Generators

- ▶ Suppose that we have some probability space S .
- ▶ S has two variables X (instances) and Y (labels) that represent cats and dogs.
- ▶ Often we want to find a boundary between different types of data within a sample space.
(Discriminative Model)



Generators

- ▶ Suppose that we have some probability space S .
- ▶ S has two variables X (instances) and Y (labels) that represent cats and dogs.
- ▶ Often we want to find a boundary between different types of data within a sample space.
(Discriminative Model)
- ▶ We might instead want to *sample* from this space and see what images are here.

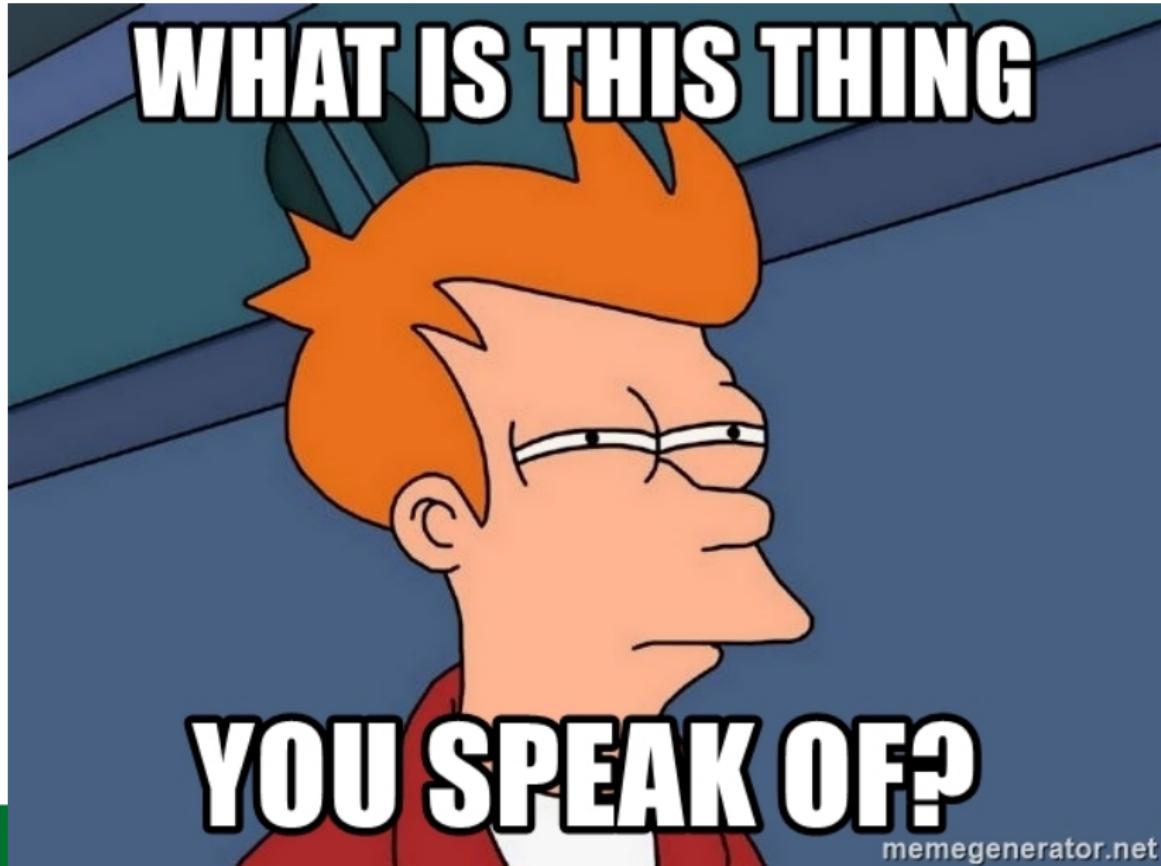


Generators

- ▶ Suppose that we have some probability space S .
- ▶ S has two variables X (instances) and Y (labels) that represent cats and dogs.
- ▶ Often we want to find a boundary between different types of data within a sample space.
(Discriminative Model)
- ▶ We might instead want to *sample* from this space and see what images are here.
- ▶ Typically we'll only work with a single classification.



What Is A Generator?



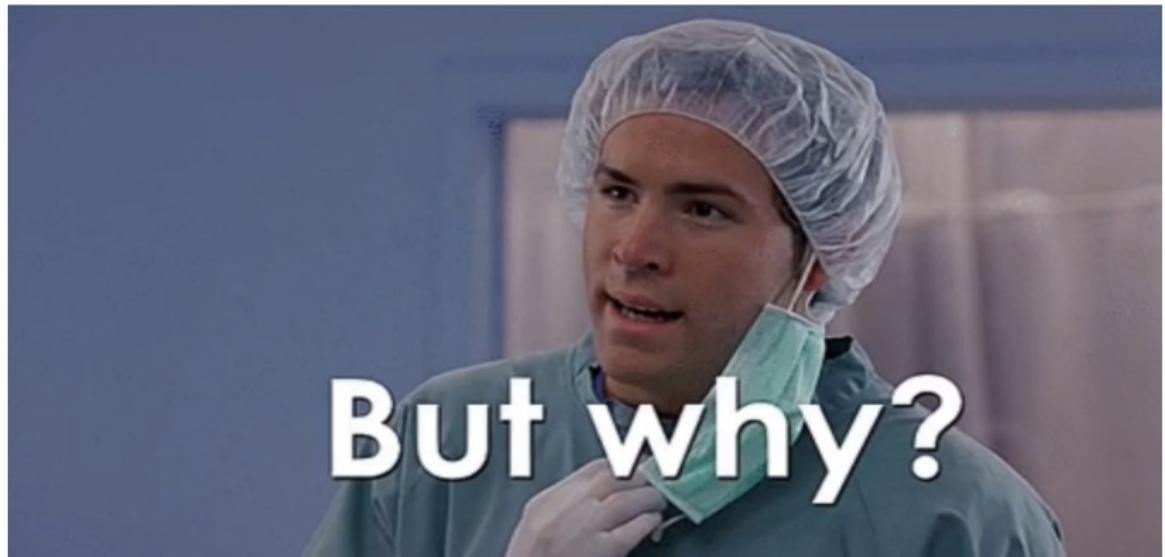
Generators

- ▶ In a Discriminative model we learn a conditional probability $P(Y = \text{label} | X = \text{image})$ (predict the label given an image).
- ▶ In a Generative model we learn to sample for $p(x)$ (produce an image).

Generators

- ▶ In a Discriminative model we learn a conditional probability $P(Y = \text{label} | X = \text{image})$ (predict the label given an image).
- ▶ In a Generative model we learn to sample for $p(x)$ (produce an image).
- ▶ In reality we learn $p_\theta(\hat{x})$ from $p_{true}(x)$

But Why?



But why?

But Why?

- ▶ We can create images that don't exist.

But Why?

- ▶ We can create images that don't exist.
- ▶ We can create new data.

But Why?

- ▶ We can create images that don't exist.
- ▶ We can create new data.
- ▶ We can upscale images/videos.

But Why?

- ▶ We can create images that don't exist.
- ▶ We can create new data.
- ▶ We can upscale images/videos.
- ▶ We can learn latent variables.

But Why?

- ▶ We can create images that don't exist.
- ▶ We can create new data.
- ▶ We can upscale images/videos.
- ▶ We can learn latent variables.
- ▶ We can uncover biases.

But Why?

- ▶ We can create images that don't exist.
- ▶ We can create new data.
- ▶ We can upscale images/videos.
- ▶ We can learn latent variables.
- ▶ We can uncover biases.
- ▶ And so much more!

But Why?: Generate Images

- ▶ We can generate cats that don't exist.
- ▶ Helps create new data within a dataset.
- ▶ Can create fake characters/animals/scenes that don't exist in real life but look realistic (art).
- ▶ Don't have to pay models or liability.



But Why?: Upscaling

- ▶ We can take lossy images and produce higher quality images.
- ▶ We can also take smaller images and make them larger.
- ▶ Helpful in compression.
- ▶ Restoration of images (also inpainting)
- ▶ Make old images/videos look new/modern.



But Why?: Learn Latent Variables

- ▶ We can learn the different latent variables in a dataset.
- ▶ For example we can learn how to make someone old!
- ▶ We could also use this as a form of compression (autoencoders)
- ▶ We can add/modify features of someone (add a mustache).



But Why?: Uncovering Bias

- ▶ Suppose we have a dataset of faces.



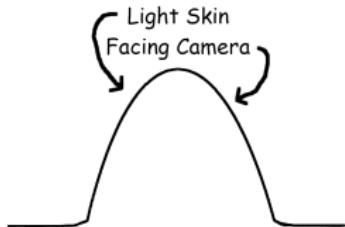
But Why?: Uncovering Bias

- ▶ Suppose we have a dataset of faces.
- ▶ Might have some distribution, let's say...



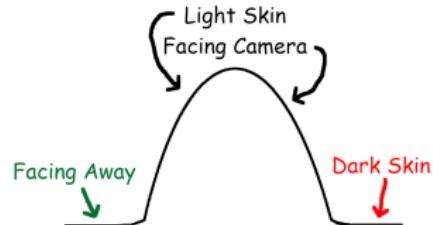
But Why?: Uncovering Bias

- ▶ Suppose we have a dataset of faces.
- ▶ Might have some distribution, let's say...
- ▶ Most of the faces fit in the main part of the distribution.



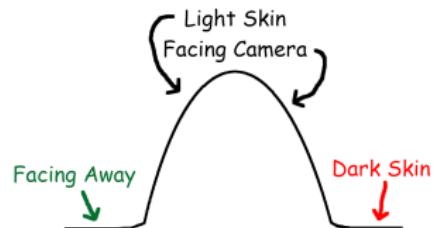
But Why?: Uncovering Bias

- ▶ Suppose we have a dataset of faces.
- ▶ Might have some distribution, let's say...
- ▶ Most of the faces fit in the main part of the distribution.
- ▶ We have outliers in the dataset.

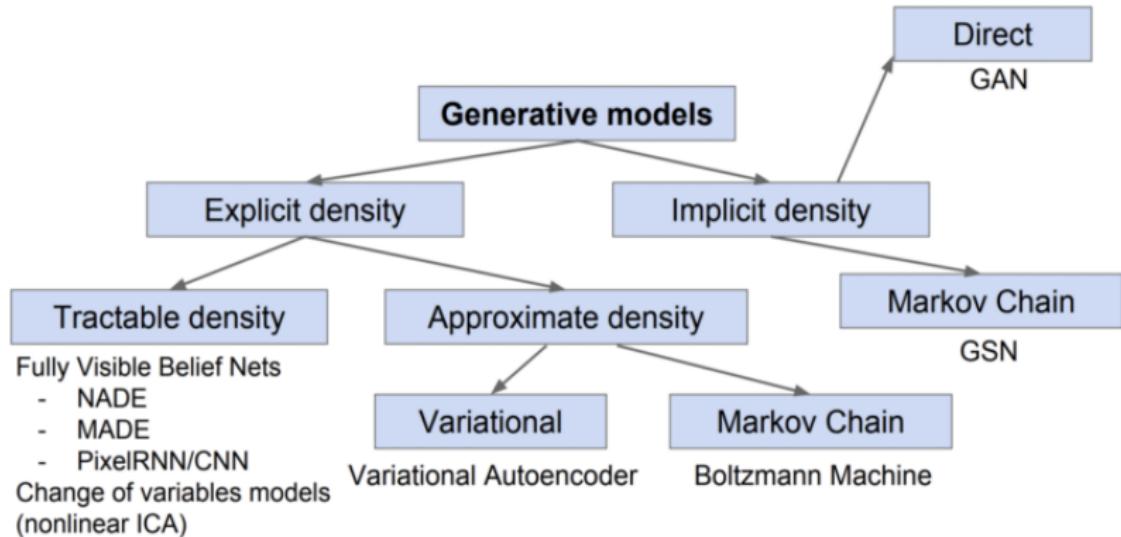


But Why?: Uncovering Bias

- ▶ Suppose we have a dataset of faces.
- ▶ Might have some distribution, let's say...
- ▶ Most of the faces fit in the main part of the distribution.
- ▶ We have outliers in the dataset.
- ▶ We can now adjust our dataset (gather new data)or potentially generate new images for a classifier if we can't fix the distribution (might not be logically possible).

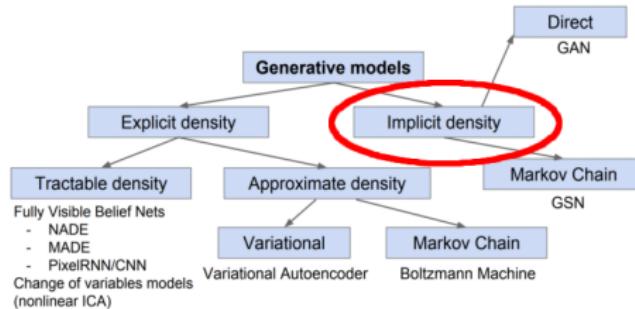


Taxonomy of Generators

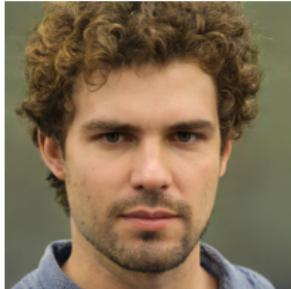


Taxonomy of Generators: Implicit Density

- ▶ We want to learn a stochastic procedure to generate data.
- ▶ We learn likelihood-free.



This X Does Not Exist



(a) This Person
Does Not Exist
(StyleGAN2)



(b) This Cat Does
Not Exist
(StyleGAN)



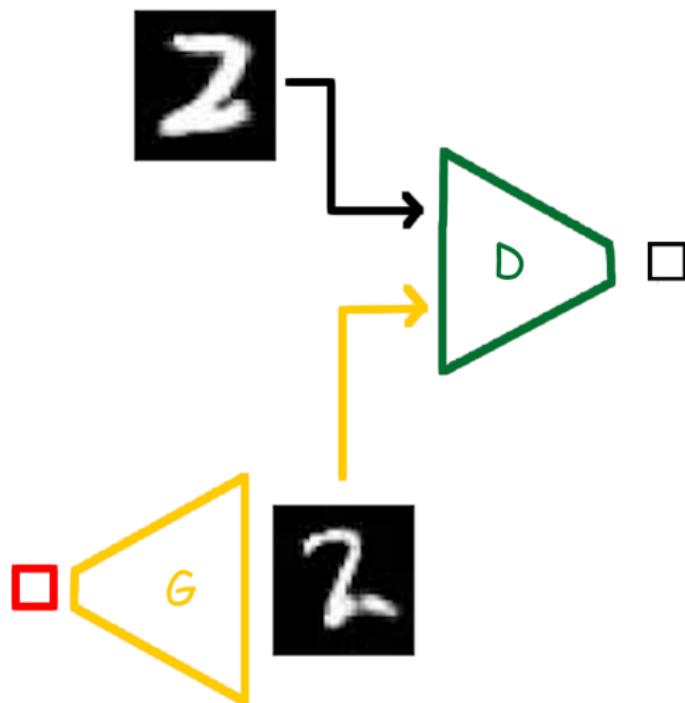
(c) This Anime Does Not Exist
(StyleGAN2 + GPT-3)

Review of the themes and plot of the first episode of the new anime "Dobon no Oji-sama," aka "The Master of Dallas." which will debut in Japan this October. Official Synopsis: Shin is a 10-year-old boy who lives in Chausugaiji, Maina — a Reservation. Aizuwakaba Takanori (Toshiyuki Morikubo) is the owner of a shop there via Crucible, a city of recycled materials and discarded break piles where children with superpowers are rubish chasers. They scavenge through the trash to find enough food to survive. A single young boy with the power of "Regeneration" is the last child living there. He is called "The King of the Fifth" and is a rock legend who is said to gather piles of discarded garbage and transform them into golden piles of discarded garbage and transform them into gold. Shin is the smallest of his peers but has incredible strength. He dreams of becoming a superhero. His best friend is Kuro, who has come to him, but has no idea how to go about doing so. One day, Kuro steals a balloon from a boy who is selling newspapers and steals the key that Shin has parked and prepared to use to open the shop. Shin is shocked and angry. While Shin is desperate to retrieve his key, he is taught by Kuro that "If you want to be a superhero, you have to give things back." After stealing a balloon from a boy with the power of destruction, he makes a decision...
The main cast includes: Toshiyuki Morikubo as Shin and Kuro; Chiaki Omigawa as Kuro; and Akinori Shibusawa as Shin's mother. With no original characters, the anime borrows its characters from abandoned island in the Pacific Ocean ruled by Children with the power of Regeneration and the drama of psychological issues of a young boy who uses his powers to survive. The animation is good, but the works depict the struggles of boy Children in totally different ways than the original anime. It is original creative. His seemingly light hearted and com-

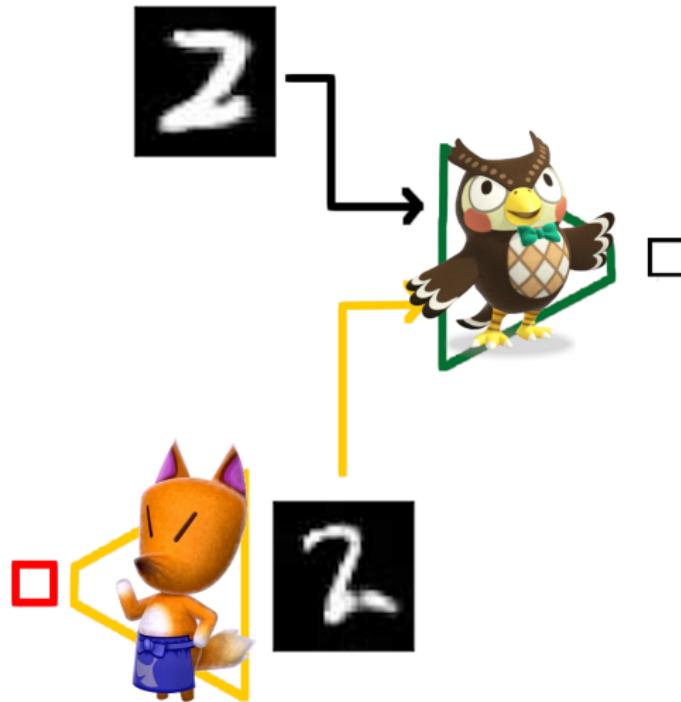
Generative Adversarial Networks



Generative Adversarial Networks



Generative Adversarial Networks



Generative Adversarial Networks

- ▶ Send noise into Generator.
- ▶ Send generated image and real images into Discriminator.
- ▶ Play mini-max game teaching the Generator to fool the Discriminator.
- ▶ When fake image can adequately fool (a well trained) discriminator we produce quality images.

Generative Adversarial Networks

- ▶ Send noise into Generator.
- ▶ Send generated image and real images into Discriminator.
- ▶ Play mini-max game teaching the Generator to fool the Discriminator.
- ▶ When fake image can adequately fool (a well trained) discriminator we produce quality images.
- ▶ In practice we train D, then G then repeat.

But Why Is This Implicit?

- ▶ We didn't learn a density function.
- ▶ We just focus on generating samples.
- ▶ Don't have an (easy) smooth transition between latent variables.
- ▶ Likelihood is not tractable.
- ▶ "We're training them to memorize, not generalize" - Ian Goodfellow

Latent Space



Latent Space



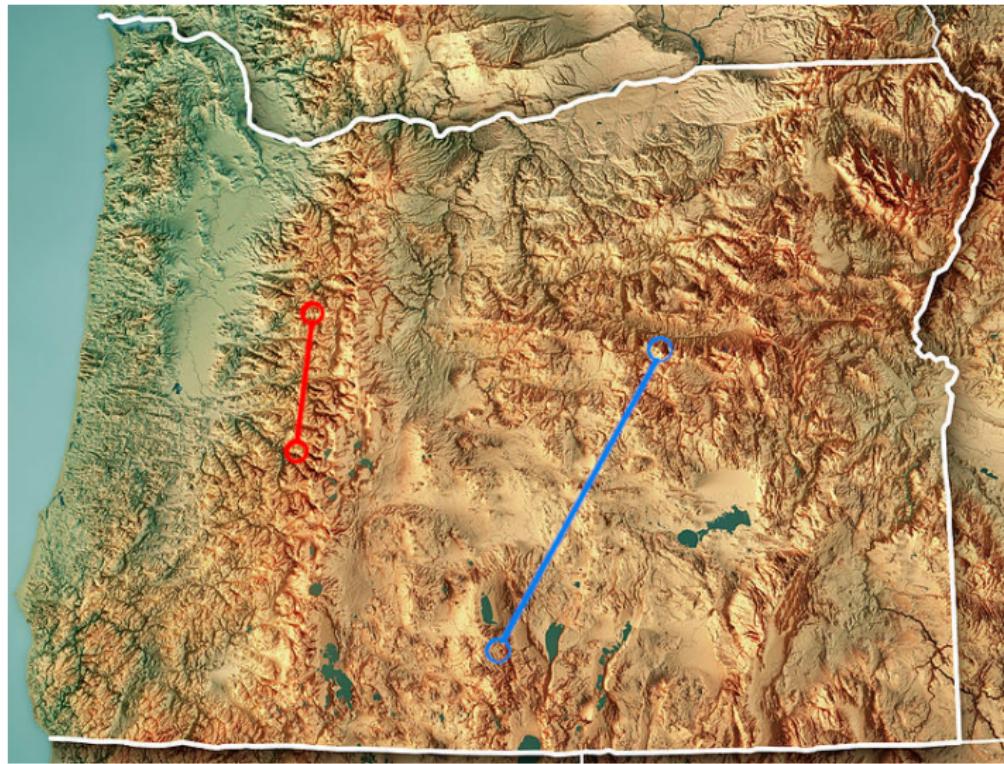
Latent Space



Latent Space



Latent Space



Latent Space

TSNE visualization of vanilla GAN latent space.

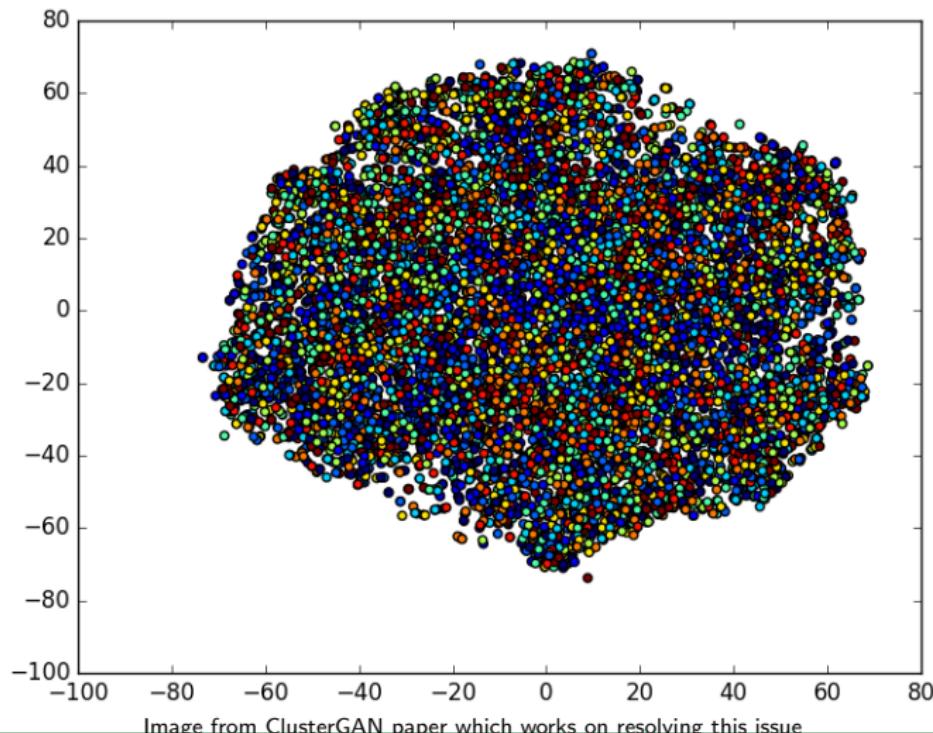


Image from ClusterGAN paper which works on resolving this issue

So Why Use GANs?

- ▶ Because they are quick.
- ▶ They produce high quality images.
- ▶ Can still do many tasks.
- ▶ They're well researched.

StyleGAN2

Let's look at a modern GAN and how it works.

Problems With StyleGAN

- ▶ GANs typically have “blobs” that are mysteriously produced.
- ▶ Often get “GAN Monsters” where parts of images get wildly distorted.
- ▶ Certain characteristics are too common. Eyes forward, smile alignment.

Blobs and More Blobs



Figure 1. Instance normalization causes water droplet -like artifacts in StyleGAN images. These are not always obvious in the generated images, but if we look at the activations inside the generator network, the problem is always there, in all feature maps starting from the 64x64 resolution. It is a systemic problem that plagues all StyleGAN images.

Monsters



(a) Low PPL scores



(b) High PPL scores

Figure 4. Connection between perceptual path length and image quality using baseline StyleGAN (config A) with LSUN CAT. (a) Random examples with low PPL ($\leq 10^{\text{th}}$ percentile). (b) Examples with high PPL ($\geq 90^{\text{th}}$ percentile). There is a clear correlation between PPL scores and semantic consistency of the images.

Phase Characteristics

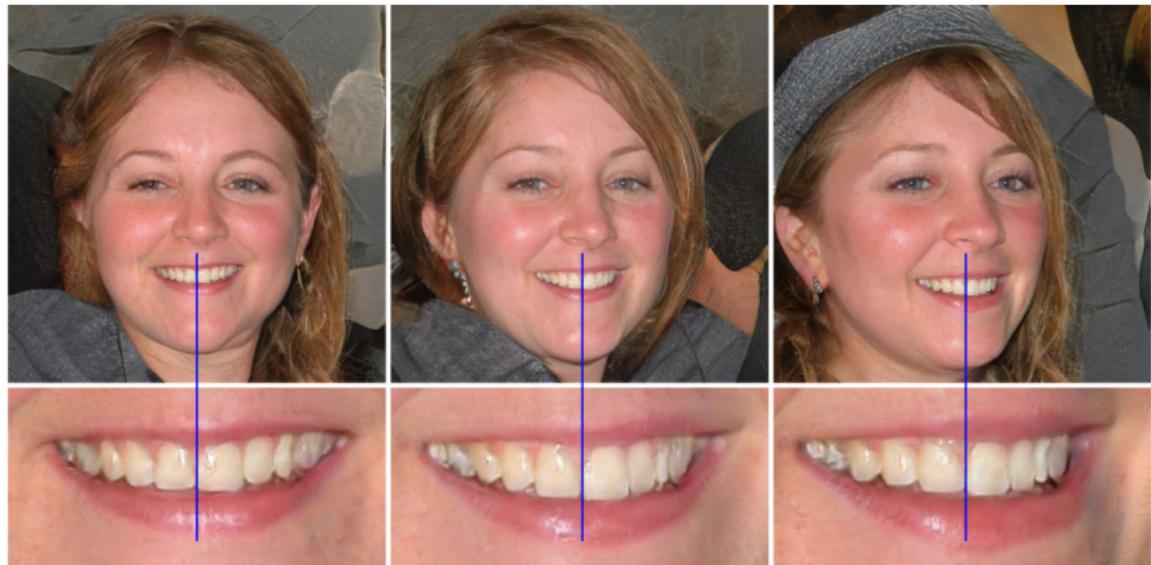
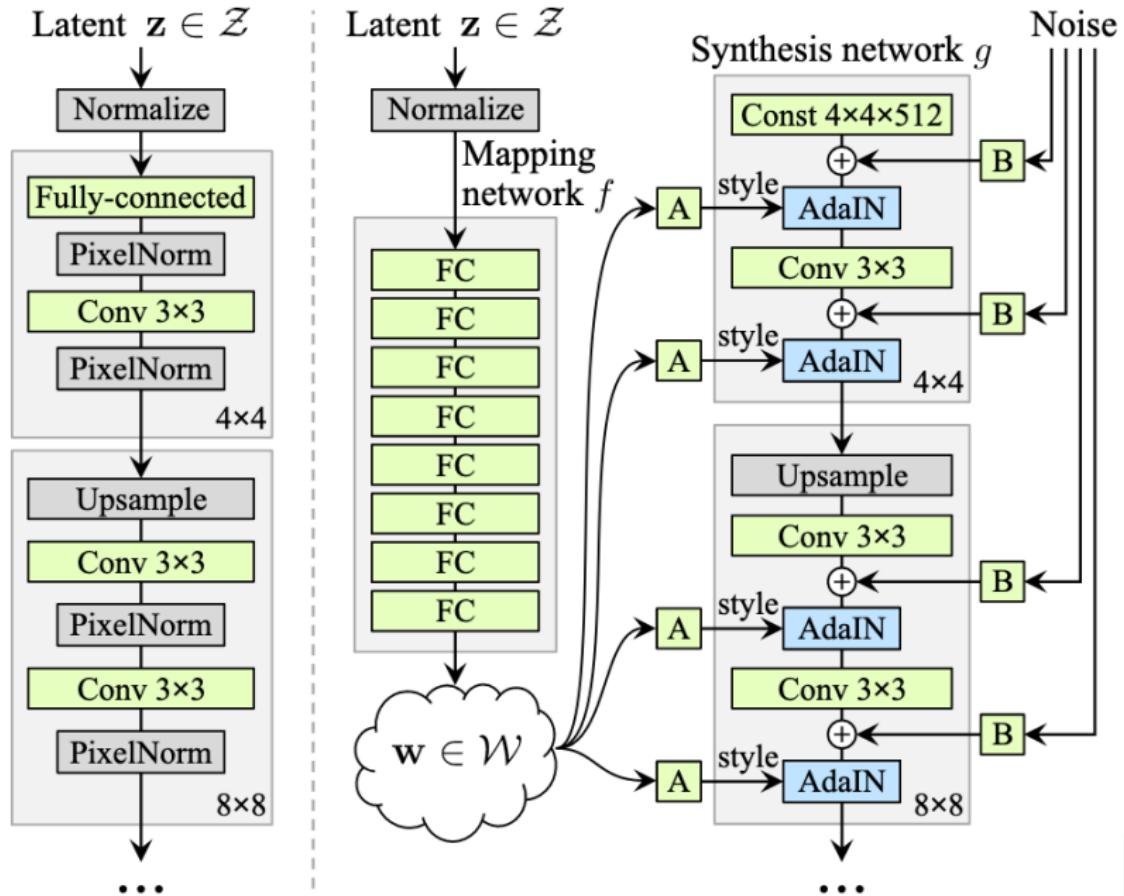


Figure 6. Progressive growing leads to “phase” artifacts. In this example the teeth do not follow the pose but stay aligned to the camera, as indicated by the blue line.

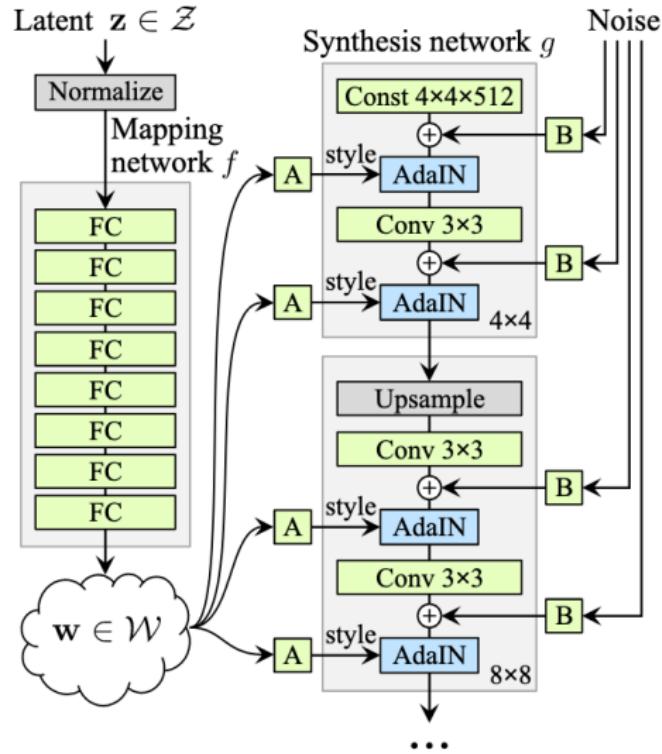
StyleGAN (OG)



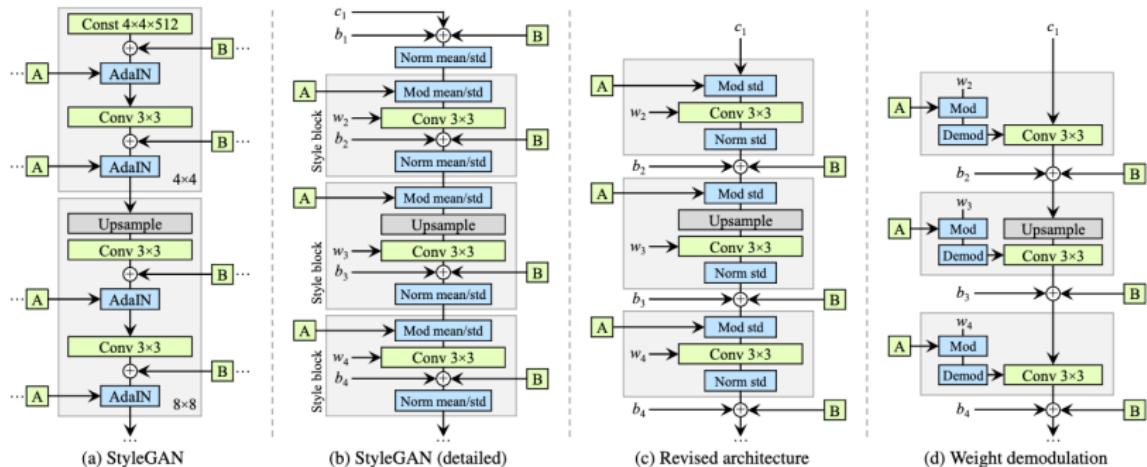
StyleGAN (OG)

- ▶ Uses upscaling convolutions, doubling.
- ▶ Learns a “style space” (w)
- ▶ Introduces Adaptive Instance Normalization (AdaIN)

$$\text{AdaIN}(x_i, y) = y_{s,i} \frac{x_i - \mu(x_i)}{\sigma(x_i)} + y_{b,i}$$



StyleGAN vs StyleGAN 2



StyleGAN2: What Changed

- ▶ Removed AdaIN. Removing this removes the droplet/blobs.
- ▶ Lazy Regularization: Split loss and regularization terms and regularize every 16 minibatches.
- ▶ Noticed low Perceptual Path Length (PPL) results had better images so introduced a new regularizer to account for this.
- ▶ Use skip generator and residual discriminator instead of progressive growing.
- ▶ Increase network size (helps get from 512x512 to 1024x1024).

StyleGAN2 Results

Dataset	Resolution	StyleGAN (A)		StyleGAN2 (F)	
		FID	PPL	FID	PPL
LSUN CAR	512×384	3.27	1485	2.32	416
LSUN CAT	256×256	8.53	924	6.93	439
LSUN CHURCH	256×256	4.21	742	3.86	342
LSUN HORSE	256×256	3.83	1405	3.43	338

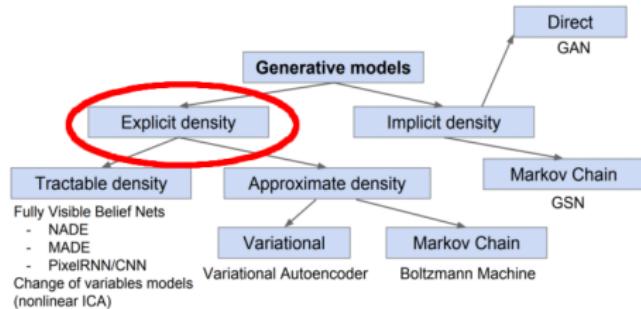
Table 3. Improvement in LSUN datasets measured using FID and PPL. We trained CAR for 57M images, CAT for 88M, CHURCH for 48M, and HORSE for 100M images.

StyleGAN2 Results



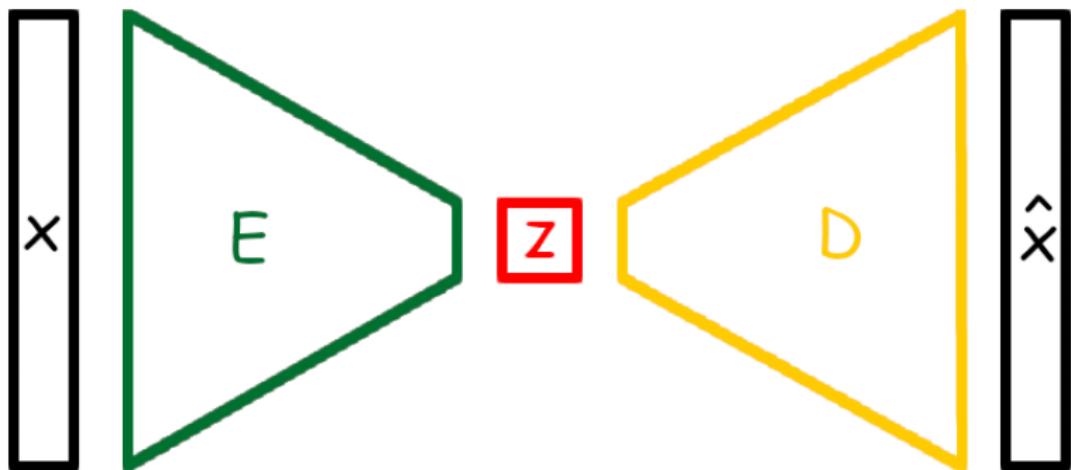
Taxonomy of Generators: Explicit Density

- ▶ We want to learn a distribution of data.
- ▶ We assume some prior about the distribution of the data.
- ▶ We learn the log-likelihood of the function $\log p_\theta(x)$



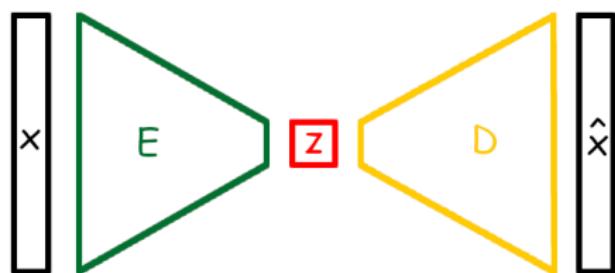
source: ian goodfellow

AutoEncoders



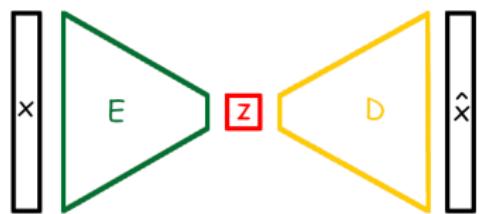
AutoEncoders

- ▶ Want to learn latent variable z from x . $\mathbb{R}^n \mapsto \mathbb{R}^m$ where $m < n$
- ▶ Compressed representation of the data.
- ▶ Since we don't have training data for z , we construct a decoder to do unsupervised learning.



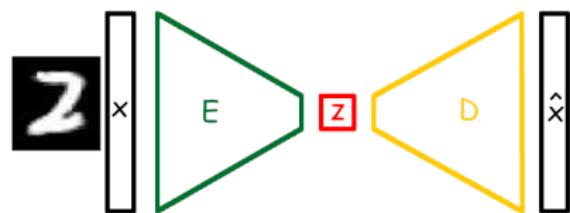
AutoEncoders

- ▶ Want to take x and produce \hat{x}



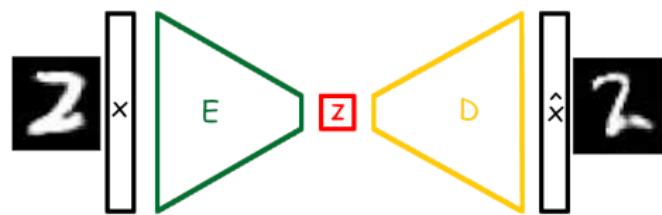
AutoEncoders

- ▶ Want to take x and produce \hat{x}
- ▶ Sample x from dataset and encode to latent space z .



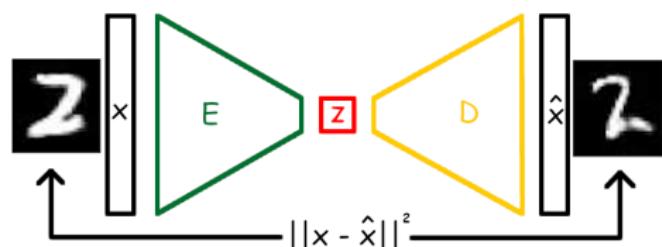
AutoEncoders

- ▶ Want to take x and produce \hat{x}
- ▶ Sample x from dataset and encode to latent space z .
- ▶ Decode to sample from z and generate \hat{x}



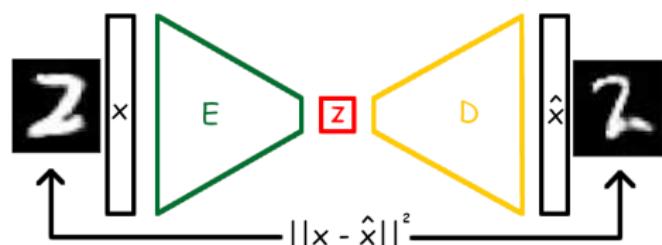
AutoEncoders

- ▶ Want to take x and produce \hat{x}
- ▶ Sample x from dataset and encode to latent space z .
- ▶ Decode to sample from z and generate \hat{x}
- ▶ Regularize $||x - \hat{x}||^2$



AutoEncoders

- ▶ Want to take x and produce \hat{x}
- ▶ Sample x from dataset and encode to latent space z .
- ▶ Decode to sample from z and generate \hat{x}
- ▶ Regularize $||x - \hat{x}||^2$
- ▶ Gives us a unsupervised training.



Pick Our Latent Space Size Carefully

Dimensionality of latent space →
reconstruction quality

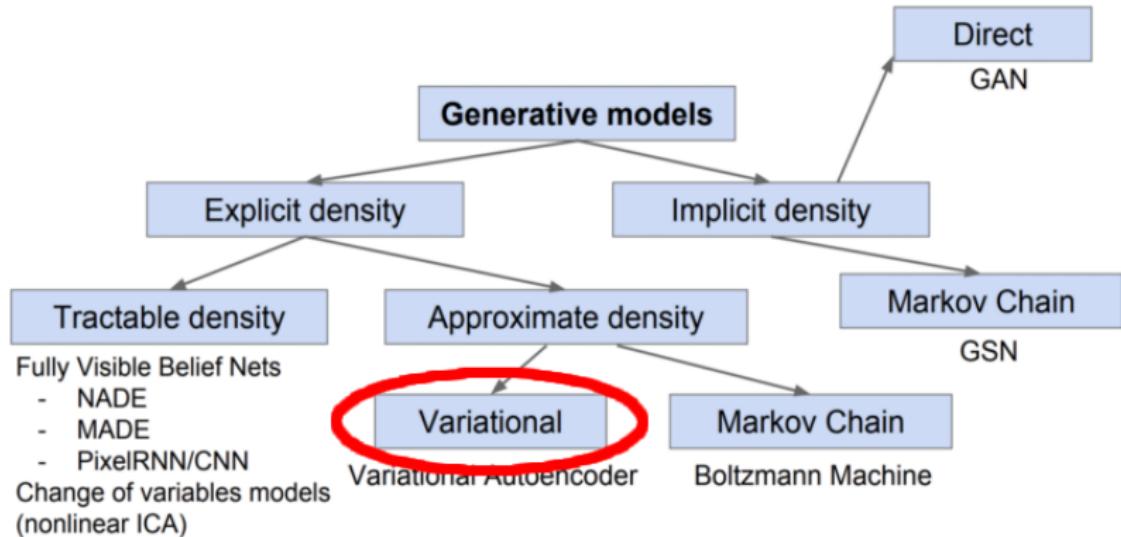
Autoencoding is a form of compression!
Smaller latent space will force a larger training bottleneck

2D latent space	5D latent space	Ground Truth
7 2 1 0 4 1 4 9 9 9 9	7 2 1 0 4 1 4 9 9 9	7 2 1 0 4 1 4 9 5 9
0 6 9 0 1 5 9 7 8 9	0 0 9 0 1 5 9 7 3 4	0 6 9 0 1 5 9 7 8 4
9 6 6 5 4 0 7 9 0 1	9 6 6 5 4 0 7 4 0 1	9 6 6 5 4 0 7 4 0 1
3 1 3 0 7 2 7 1 2 1	3 1 3 0 7 2 7 1 2 1	3 1 3 4 7 2 7 1 2 1
1 7 4 2 3 5 1 2 9 4	1 7 4 2 3 5 1 2 9 4	1 7 4 2 3 5 1 2 4 4
6 3 5 5 6 0 4 1 9 8	6 3 5 5 6 0 4 1 9 8	6 3 5 5 6 0 4 1 9 5
7 8 4 2 9 9 6 4 3 0	7 8 4 3 7 4 6 4 3 0	7 8 4 3 7 4 6 4 3 0
7 0 2 9 1 9 3 2 9 7	7 0 2 9 1 7 3 2 9 7	7 0 2 9 1 7 3 2 9 7
9 6 2 7 8 4 7 3 6 1	9 6 2 7 8 4 7 3 6 1	9 6 2 7 8 4 7 3 6 1
3 6 4 3 1 4 1 7 6 9	3 6 9 3 1 4 1 7 6 9	3 6 9 3 1 4 1 7 6 9

The BIG Problem!

We can't generate new and unique samples!

Taxonomy of Generators



source: ian goodfellow

Variational Autoencoder (VAE)

- ▶ We need a way to generate **new** samples that aren't seen in the dataset.
- ▶ Need to add some variance and noise (make non-deterministic).
- ▶ Still have to be able to do back propagation.

VAE: Breaking Down the Problem

- ▶ Suppose there is some hidden (latent) variable z that produces observation x , $p(z|x)$. **We want to learn this.**

VAE: Breaking Down the Problem

- ▶ Suppose there is some hidden (latent) variable z that produces observation x , $p(z|x)$. **We want to learn this.**
- ▶ $p(z|x) = \frac{p(x|z)p(z)}{p(x)}$

VAE: Breaking Down the Problem

- ▶ Suppose there is some hidden (latent) variable z that produces observation x , $p(z|x)$. **We want to learn this.**
- ▶ $p(z|x) = \frac{p(x|z)p(z)}{p(x)}$
- ▶ $p(x) = \int p(x|z)p(z)dz$ (intractable!)

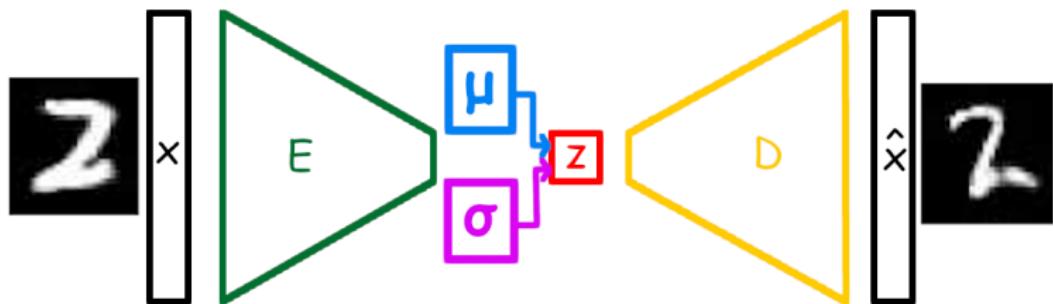
VAE: Breaking Down the Problem

- ▶ Suppose there is some hidden (latent) variable z that produces observation x , $p(z|x)$. **We want to learn this.**
- ▶ $p(z|x) = \frac{p(x|z)p(z)}{p(x)}$
- ▶ $p(x) = \int p(x|z)p(z)dz$ (intractable!)
- ▶ We can approximate $p(z|x)$ with another distribution $q(z|x)$ that is tractable.

VAE: Breaking Down the Problem

- ▶ Suppose there is some hidden (latent) variable z that produces observation x , $p(z|x)$. **We want to learn this.**
- ▶ $p(z|x) = \frac{p(x|z)p(z)}{p(x)}$
- ▶ $p(x) = \int p(x|z)p(z)dz$ (intractable!)
- ▶ We can approximate $p(z|x)$ with another distribution $q(z|x)$ that is tractable.
- ▶ We'll then need to measure the difference between p and q .

Variational Autoencoder (VAE)



Variational Autoencoder (VAE)

- ▶ Encoder computes $p_\phi(z|x)$
(reconstruction loss, like autoencoder)



Variational Autoencoder (VAE)

- ▶ Encoder computes $p_\phi(z|x)$
(reconstruction loss, like autoencoder)
- ▶ Decoder computes $q_\theta(x|z)$
(regularization term)



Variational Autoencoder (VAE)

- ▶ Encoder computes $p_\phi(z|x)$
(reconstruction loss, like autoencoder)
- ▶ Decoder computes $q_\theta(x|z)$
(regularization term)
- ▶ $L(\phi, \theta, x) = ||x - \hat{x}||^2 + KL(p_\phi(z|x)||p(x))$



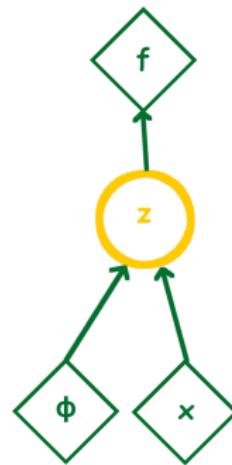
(VAE) Regularization Term

- ▶ $L(\phi, \theta, x) = \|x - \hat{x}\|^2 + KL(p_\phi(z|x)||p(z))$
- ▶ We are using the KL-Divergence between the *inferred latent distribution* ($p_\phi(z|x)$) and the *fixed prior on latent distribution* ($p(z)$).
- ▶ What prior? Why not Gaussian?
 $p(z) = \mathcal{N}(\mu = 0, \sigma^2 = 1)$
- ▶ Encourages encoding to distribute evenly around latent space and penalizes when memorizing.
- ▶ $KL(p_\phi(z|x)||p(z)) = -\frac{1}{2} \sum (\sigma_j + \mu_j^2 - 1 - \log \sigma_j)$



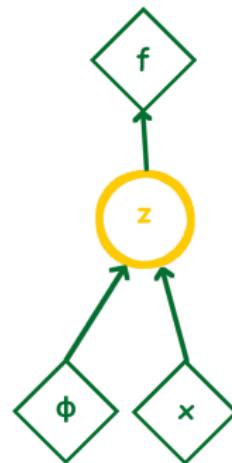
(VAE) Reparameterization

- ▶ Ops! We can't backprop through this! z is a stochastic node.



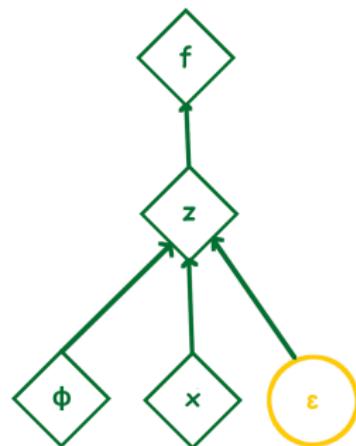
(VAE) Reparameterization

- ▶ Ops! We can't backprop through this! z is a stochastic node.
- ▶ We can reparameterize.



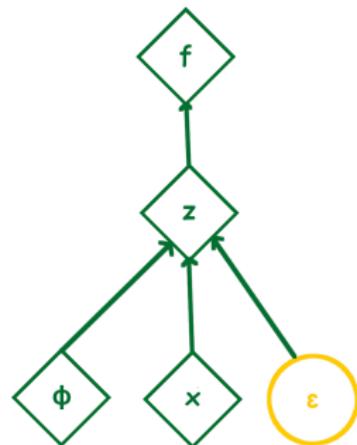
(VAE) Reparameterization

- ▶ Ops! We can't backprop through this! z is a stochastic node.
- ▶ We can reparameterize.
- ▶ $z = \mu + \sigma \odot \varepsilon$
- ▶ We let μ and σ be fixed vectors and $\varepsilon \sim \mathcal{N}(0, 1)$



(VAE) Reparameterization

- ▶ Ops! We can't backprop through this! z is a stochastic node.
- ▶ We can reparameterize.
- ▶ $z = \mu + \sigma \odot \varepsilon$
- ▶ We let μ and σ be fixed vectors and $\varepsilon \sim \mathcal{N}(0, 1)$
- ▶ Now we can backprop through ϕ and x .
- ▶ Now we can adjust ε and create new images similar to the ones we learned from.



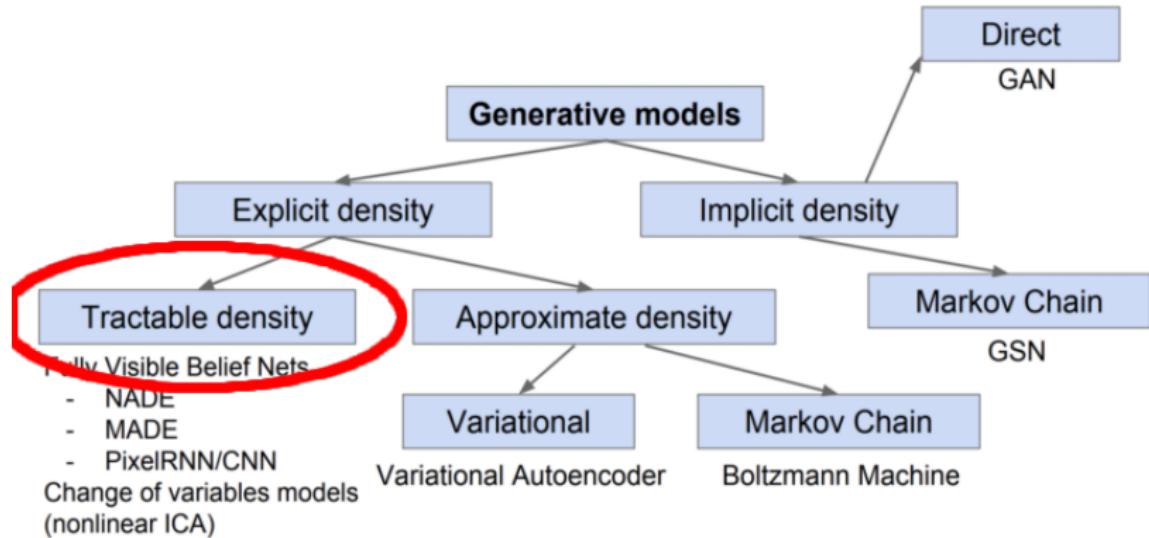
Approximate Density Estimation

- ▶ We didn't actually learn $p(x)$ but rather $q(x)$.
- ▶ We approximated the density of $p(x)$ with a tractable distribution.
- ▶ KL Divergence is always ≥ 0 .

Good and Bad of a VAE

- ▶ They are very easy to train.
- ▶ Often not high quality images.
- ▶ You learn an approximate density function (have log-likelihood)

Taxonomy of Generators



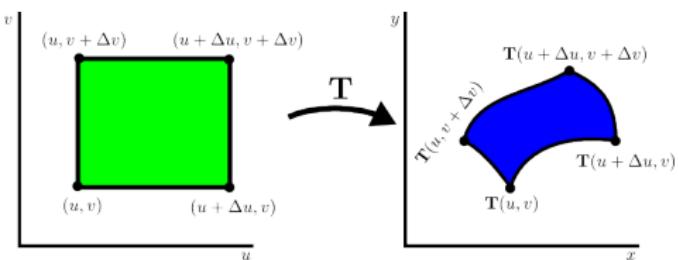
source: ian goodfellow

Density Estimation

- ▶ We've come close to learning density, but not exactly.
- ▶ Density is **very** difficult to learn.
- ▶ Having true density *should* mean better ability to perform our downstream tasks (generation, density estimation, latent variable inference, etc)

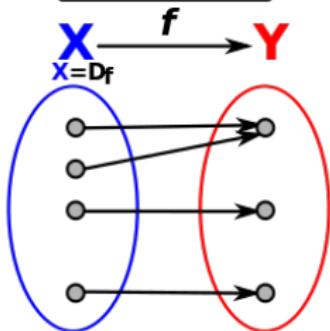
Change Of Variables

- ▶ What if we just map to a different space?
- ▶ We'll convert one probability distribution to another tractable distribution (such as a Gaussian)

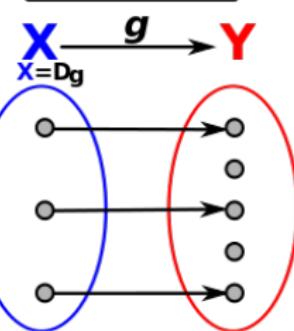


Bijection (One-to-One and On-To)

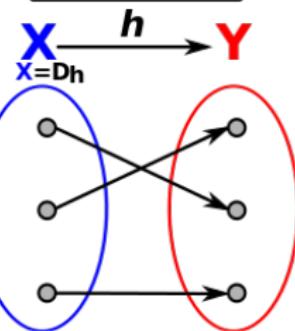
SURJECTION



INJECTION



BIJECTION



Change Of Variables

- ▶ $\int p(x)dx = \int q(z)dz = 1$ (Definition of PDF)

Change Of Variables

- ▶ $\int p(x)dx = \int q(z)dz = 1$ (Definition of PDF)
- ▶ $p(x) = q(z) \left| \frac{dz}{dx} \right|$

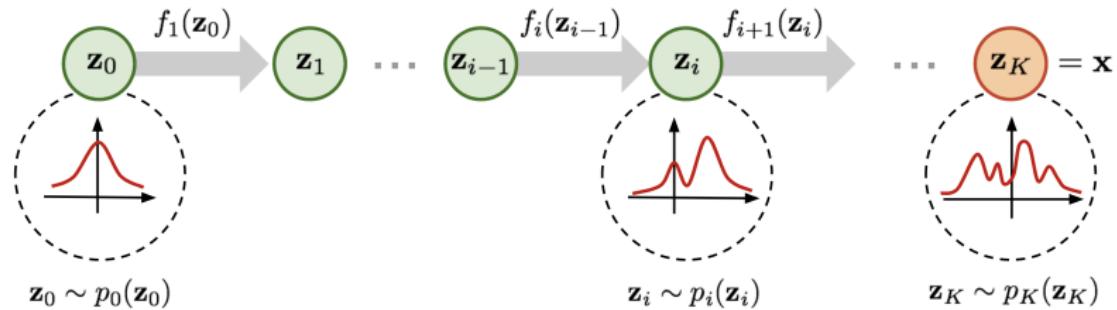
Change Of Variables

- ▶ $\int p(x)dx = \int q(z)dz = 1$ (Definition of PDF)
- ▶ $p(x) = q(z) \left| \frac{dz}{dx} \right|$
- ▶ $= q(f^{-1}(x)) \left| \frac{df^{-1}(x)}{dx} \right|$

Change Of Variables

- ▶ $\int p(x)dx = \int q(z)dz = 1$ (Definition of PDF)
- ▶ $p(x) = q(z) \left| \frac{dz}{dx} \right|$
- ▶ $= q(f^{-1}(x)) \left| \frac{df^{-1}(x)}{dx} \right|$
- ▶ $p(x) = q(f^{-1}(x)) \left| \det \frac{df^{-1}(x)}{dx} \right|$ (using the Jacobian determinant)

Normalizing Flow



Source: lilianweng.github.io

Normalizing Flow

- ▶ We need our function $f(x)$ to be easily invertible
- ▶ The Jacobian determinant needs to be *easy* to compute.
- ▶ BUT we get the exact log-likelihood of $p(x)$ (tractable) and we can train our model on the NLL

$$\mathcal{L}(D) = -\frac{1}{|D|} \sum_{x \in D} \log p(x)$$

NICE: Non-linear Independent Components Estimation

- ▶ Key Idea: What if split the transformation into two blocks? $x \mapsto (x_0, x_1)$

NICE: Non-linear Independent Components Estimation

- ▶ Key Idea: What if split the transformation into two blocks? $x \mapsto (x_0, x_1)$
- ▶ Use an easy to invert equation.

$f(x)$	$f^{-1}(y)$
$y_0 = x_0$	$x_0 = y_0$
$y_1 = x_1 + m(x_0)$	$x_1 = y_1 - m(y_0)$

NICE: Non-linear Independent Components Estimation

- ▶ Key Idea: What if split the transformation into two blocks? $x \mapsto (x_0, x_1)$
- ▶ Use an easy to invert equation.

$f(x)$	$f^{-1}(y)$
$y_0 = x_0$	$x_0 = y_0$
$y_1 = x_1 + m(x_0)$	$x_1 = y_1 - m(y_0)$

- ▶ Jacobian is trivial and triangular.

NICE: Non-linear Independent Components Estimation

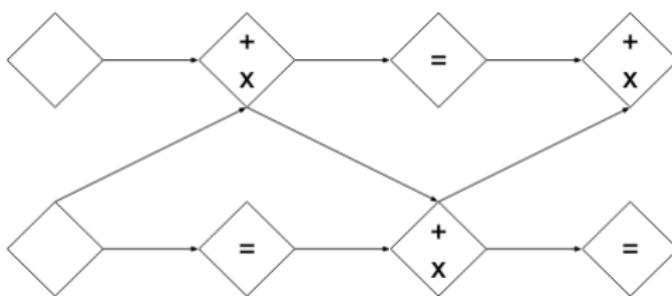
- ▶ Key Idea: What if split the transformation into two blocks? $x \mapsto (x_0, x_1)$
- ▶ Use an easy to invert equation.

$f(x)$	$f^{-1}(y)$
$y_0 = x_0$	$x_0 = y_0$
$y_1 = x_1 + m(x_0)$	$x_1 = y_1 - m(y_0)$

- ▶ Jacobian is trivial and triangular.
- ▶ 2nd Key Idea: Permute x split:
 $x_0 \mapsto (x_{0,0}, x_{0,1}), x_1 \mapsto (x_{1,1}, x_{1,0}), x_2 \mapsto \dots$

Throwing It Together

- ▶ Sample from a tractable distribution $z \sim f(x)$.
 - ▶ Make $m(x)$ a neural net.
 - ▶ Use at least 3 layers (permuting x).
 - ▶ Add a scaling term at the last layer $e^s \odot f(x)$.



NICE: Results

7	7	7	7	9	8	8	1	1	0	9
7	7	7	9	9	8	8	8	1	1	9
7	9	9	8	8	8	8	8	1	1	7
7	9	8	8	8	8	8	8	9	7	7
7	9	8	0	8	8	8	9	7	7	7
7	9	0	1	1	8	8	9	7	7	7
7	9	1	1	8	8	8	9	7	7	7
7	7	1	8	8	8	8	8	9	9	9
7	7	7	9	8	8	8	8	8	8	9
7	7	7	9	8	8	8	0	8	8	9

Why Is This Better/Worse?

- ▶ We've learned a richer family of priors compared to VAE.
- ▶ More difficult to train because we need invertible functions.
- ▶ Able to perform complex/accurate downstream tasks because we've better approximated the distribution. (e.g. Inpainting, conceptual compression)

GLOW

- ▶ NICE (and RealNVP) do not produce sharp images.
- ▶ NICE could not compete with common GAN tasks like face generation.
- ▶ Need a new model to get sharp features and produce high quality images.

GLOW: What Changed?

- ▶ Add an Actnorm layer.
- ▶ Add an invertible convolution layer to learn checkerboarding.
- ▶ Make it bigger!

GLOW: What Changed?

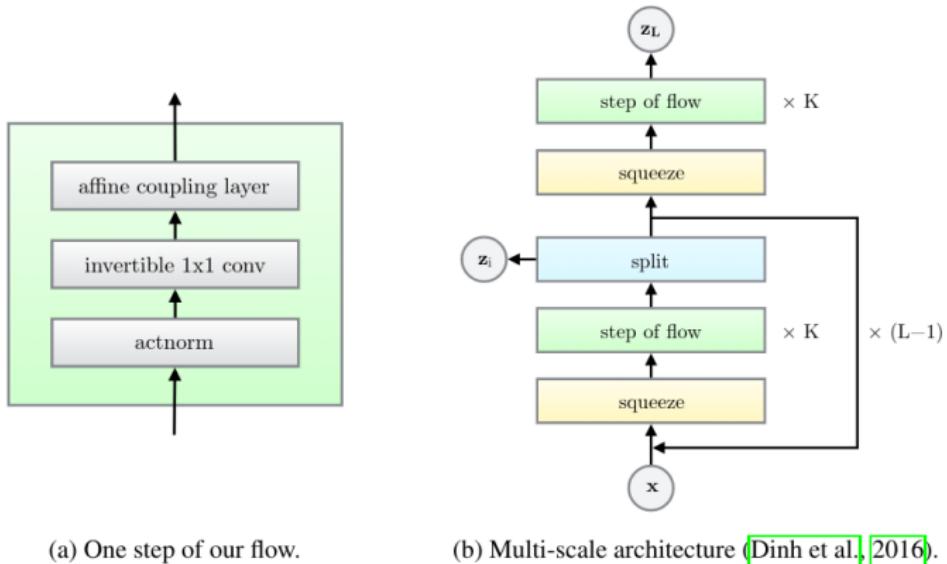


Figure 2: We propose a generative flow where each step (left) consists of an *actnorm* step, followed by an invertible 1×1 convolution, followed by an affine transformation (Dinh et al., 2014). This flow is combined with a multi-scale architecture (right). See Section 3 and Table I.

GLOW: What Changed?

Table 1: The three main components of our proposed flow, their reverses, and their log-determinants. Here, \mathbf{x} signifies the input of the layer, and \mathbf{y} signifies its output. Both \mathbf{x} and \mathbf{y} are tensors of shape $[h \times w \times c]$ with spatial dimensions (h, w) and channel dimension c . With (i, j) we denote spatial indices into tensors \mathbf{x} and \mathbf{y} . The function $\text{NN}()$ is a nonlinear mapping, such as a (shallow) convolutional neural network like in ResNets (He et al., 2016) and RealNVP (Dinh et al., 2016).

Description	Function	Reverse Function	Log-determinant
Actnorm. See Section 3.1	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{s} \odot \mathbf{x}_{i,j} + \mathbf{b}$	$\forall i, j : \mathbf{x}_{i,j} = (\mathbf{y}_{i,j} - \mathbf{b}) / \mathbf{s}$	$h \cdot w \cdot \sum(\log \mathbf{s})$
Invertible 1×1 convolution. $\mathbf{W} : [c \times c]$. See Section 3.2	$\forall i, j : \mathbf{y}_{i,j} = \mathbf{W} \mathbf{x}_{i,j}$	$\forall i, j : \mathbf{x}_{i,j} = \mathbf{W}^{-1} \mathbf{y}_{i,j}$	$h \cdot w \cdot \log \det(\mathbf{W}) $ or $h \cdot w \cdot \sum(\log \mathbf{s})$ (see eq. (10))
Affine coupling layer. See Section 3.3 and (Dinh et al., 2014)	$\mathbf{x}_a, \mathbf{x}_b = \text{split}(\mathbf{x})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{x}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{y}_a = \mathbf{s} \odot \mathbf{x}_a + \mathbf{t}$ $\mathbf{y}_b = \mathbf{x}_b$ $\mathbf{y} = \text{concat}(\mathbf{y}_a, \mathbf{y}_b)$	$\mathbf{y}_a, \mathbf{y}_b = \text{split}(\mathbf{y})$ $(\log \mathbf{s}, \mathbf{t}) = \text{NN}(\mathbf{y}_b)$ $\mathbf{s} = \exp(\log \mathbf{s})$ $\mathbf{x}_a = (\mathbf{y}_a - \mathbf{t}) / \mathbf{s}$ $\mathbf{x}_b = \mathbf{y}_b$ $\mathbf{x} = \text{concat}(\mathbf{x}_a, \mathbf{x}_b)$	$\sum(\log(\mathbf{s}))$

GLOW: Training

- ▶ Set $K = 32$ and $L = 6$ for $(32 * 6 + 32)$ 224 flow steps (NICE had 5!)
- ▶ New type of layers (NICE only had the additive coupling)
- ▶ Can train on CelebA-HQ (256x256) images.

GLOW: Results

Table 2: Best results in bits per dimension of our model compared to RealNVP.

Model	CIFAR-10	ImageNet 32x32	ImageNet 64x64	LSUN (bedroom)	LSUN (tower)	LSUN (church outdoor)
RealNVP	3.49	4.28	3.98	2.72	2.81	3.08
Glow	3.35	4.09	3.81	2.38	2.46	2.67



Figure 4: Random samples from the model, with temperature 0.7

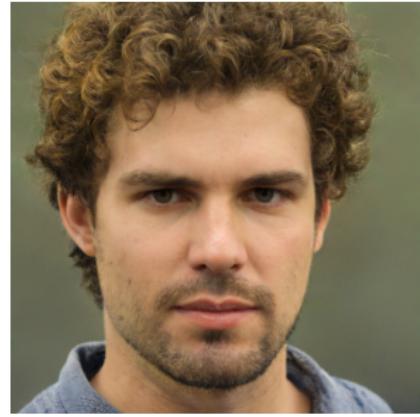
GLOW vs StyleGAN vs StyleGAN2



GLOW



StyleGAN



StyleGAN2

GLOW Downstream Tasks: Interpolation



Figure 5: Linear interpolation in latent space between real images

GLOW Downstream Tasks: Feature Manipulation



(a) Smiling



(b) Pale Skin



(c) Blond Hair



(d) Narrow Eyes



(e) Young

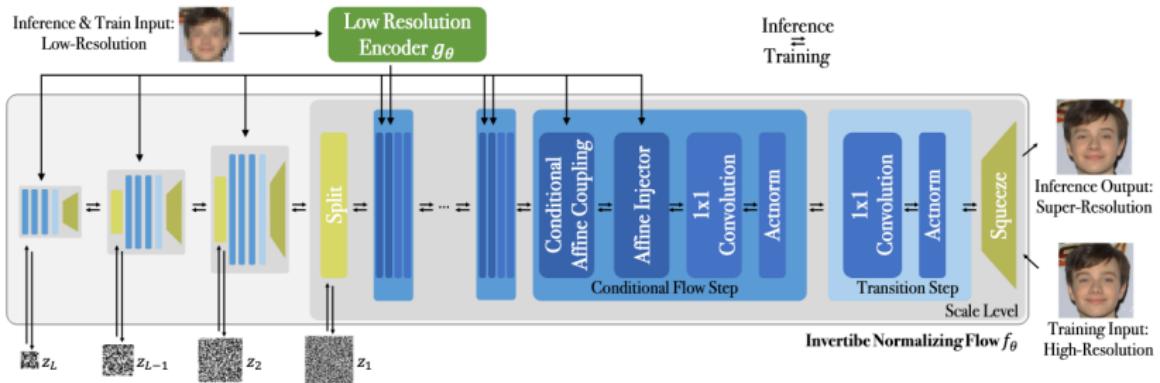


(f) Male

What About Using Flows for Super Resolution?

- ▶ Use flows as a replacement for GAN super resolution.
- ▶ Training within certain scale uses a normalizing flow.
- ▶ When scaling up and down use injective and surjective functions.
- ▶ Output is a distribution instead of a single image.

SRFlow



SRFlow: Comparison

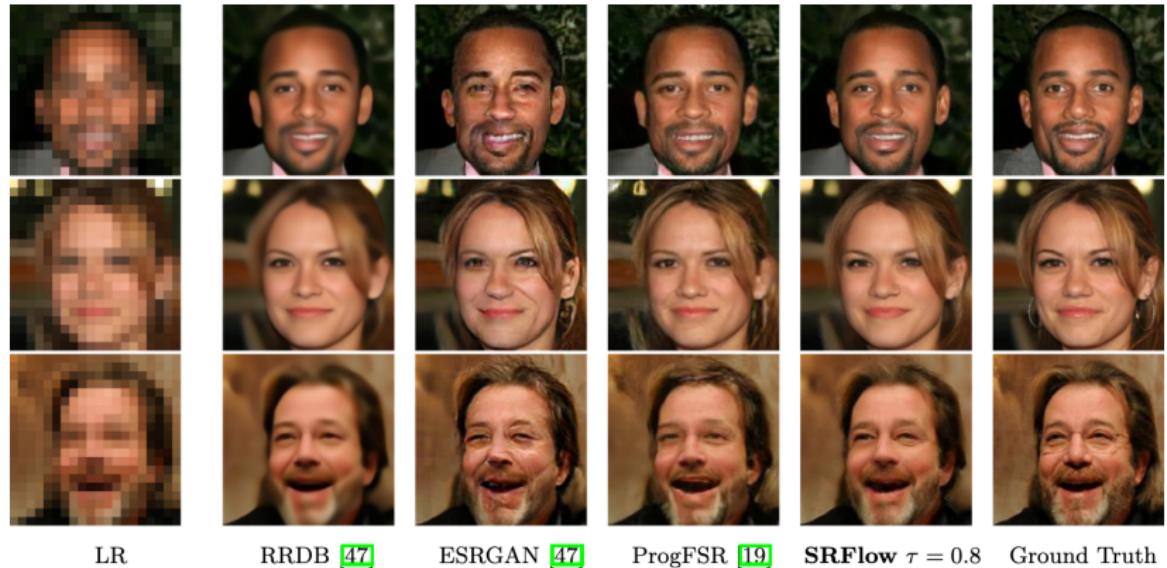
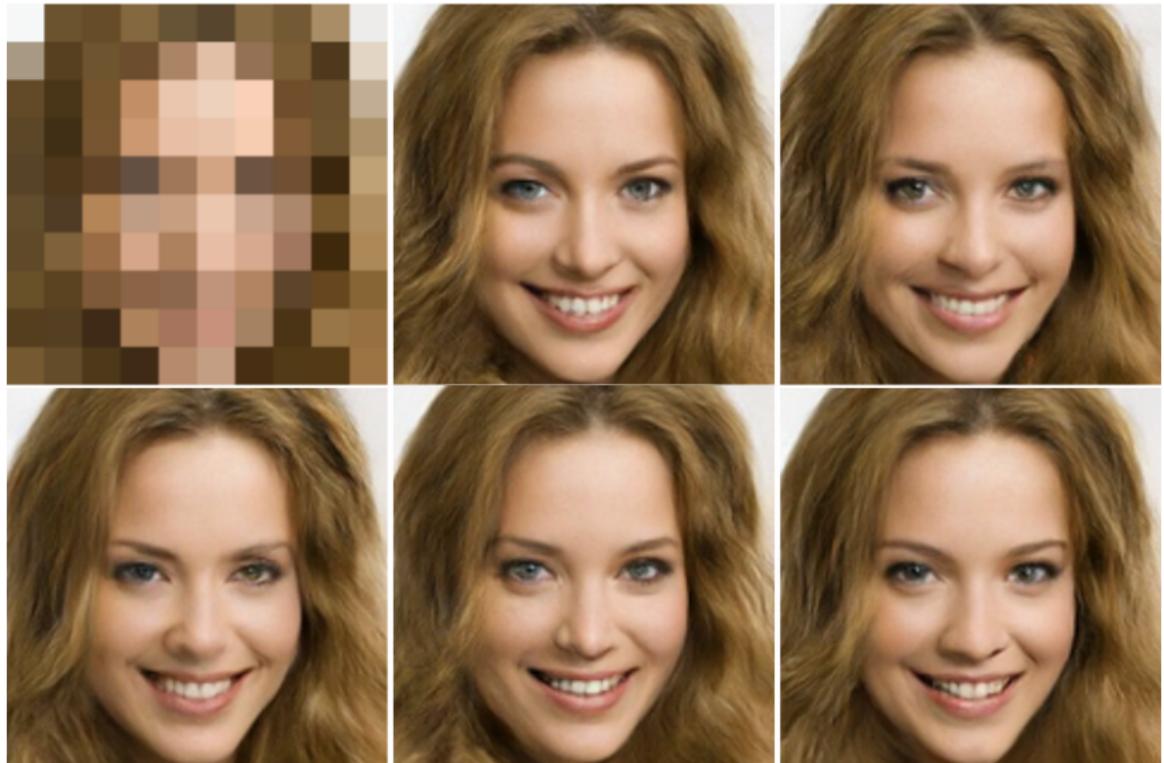


Fig. 7. Comparison of our SRFlow with state-of-the-art for $8\times$ face SR on CelebA.

SRFlow: Stochastic Output



SRFlow: Stochastic Output



Recap: Generative Models

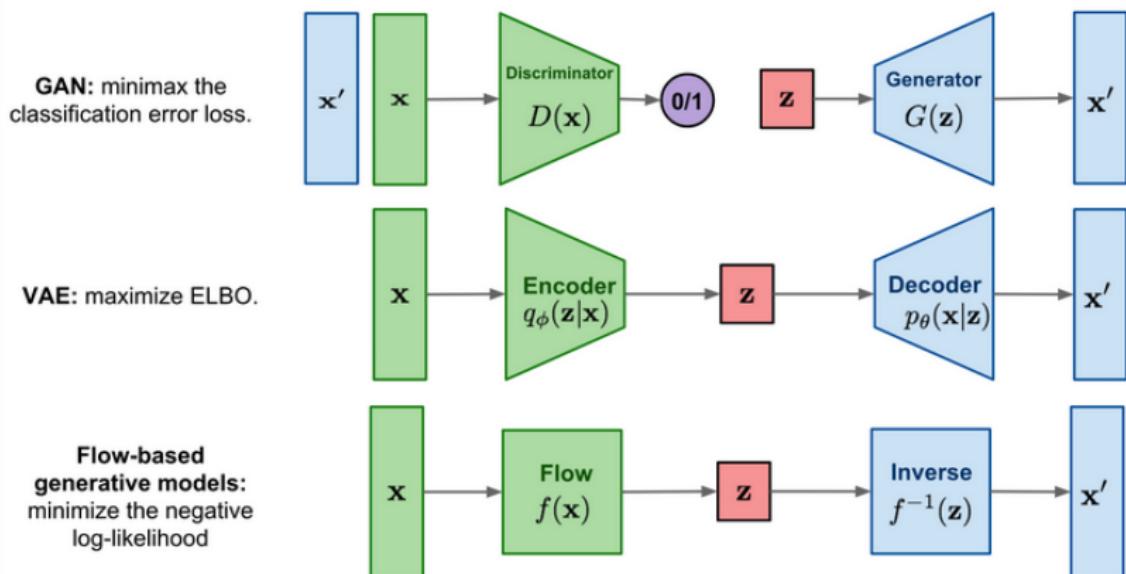


Fig. 1. Comparison of three categories of generative models.