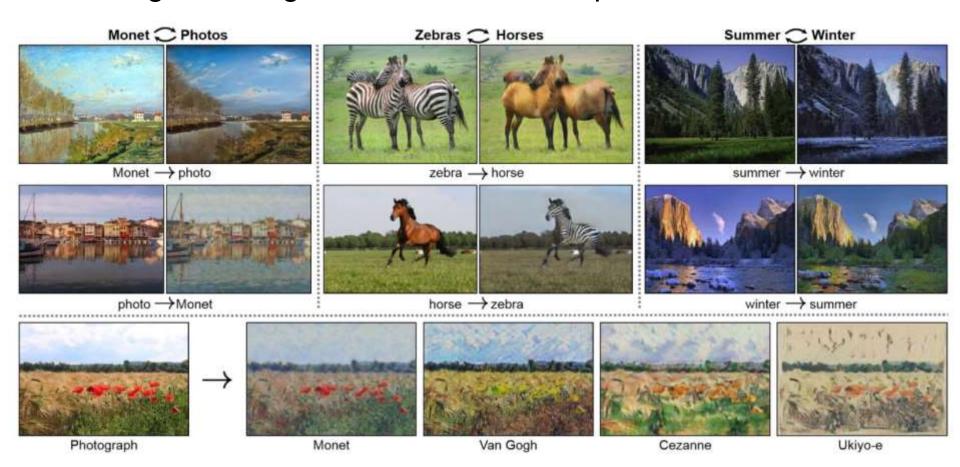
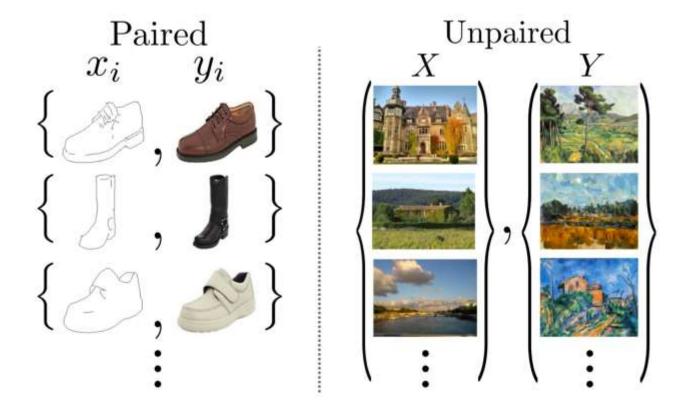
Lecture 16: Deep Generative Models V: Variants of GANs

Lan Xu SIST, ShanghaiTech Fall, 2023

Image-to-image translation without paired data

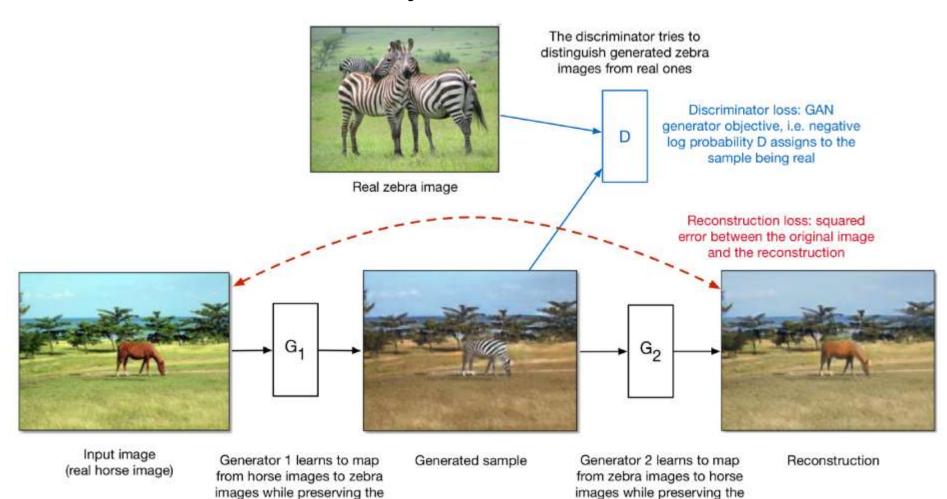


If we had paired data (same content in both styles), this would be a supervised learning problem. But this is hard to find.





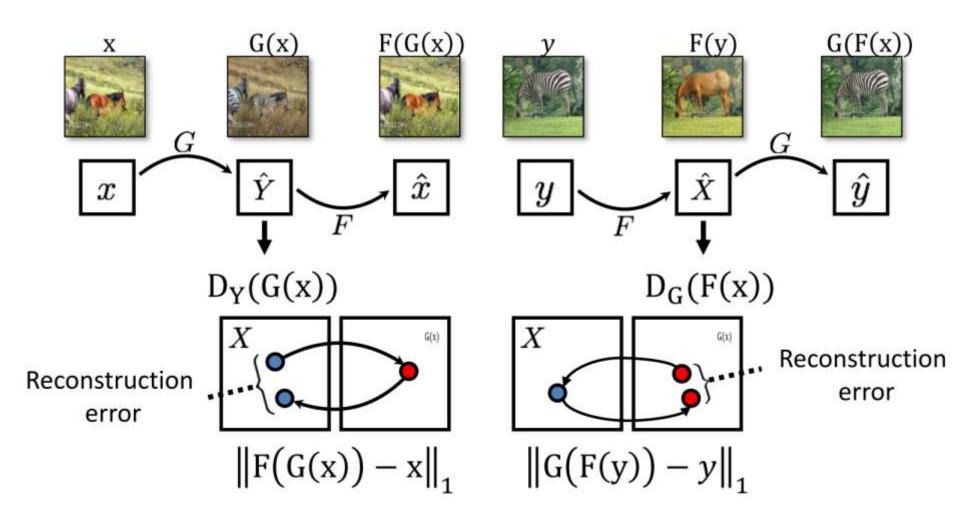
- If we had paired data (same content in both styles), this would be a supervised learning problem. But this is hard to find.
- The CycleGAN architecture learns to do it from unpaired data.
 - □ Train two different generator nets to go from style 1 to style 2, and vice versa.
 - Make sure the generated samples of style 2 are indistinguishable from real images by a discriminator net.
 - Make sure the generators are cycle-consistent: mapping from style 1 to style 2 and back again should give you almost the original image.



Total loss = discriminator loss + reconstruction loss

structure

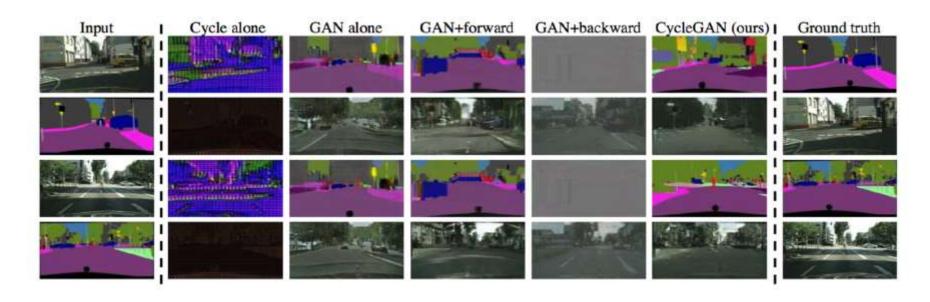
structure





Results

Style transfer between road scenes and semantic segmentations (labels of every pixel in an image by object category):



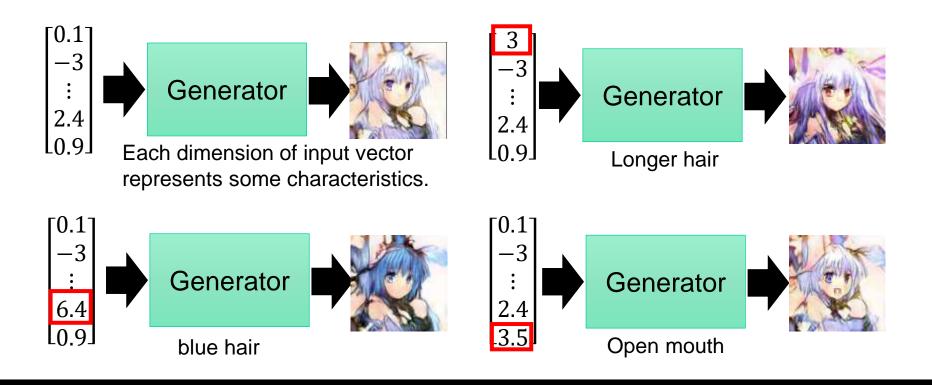
Results



More details

https://hardikbansal.github.io/CycleGANBlog/

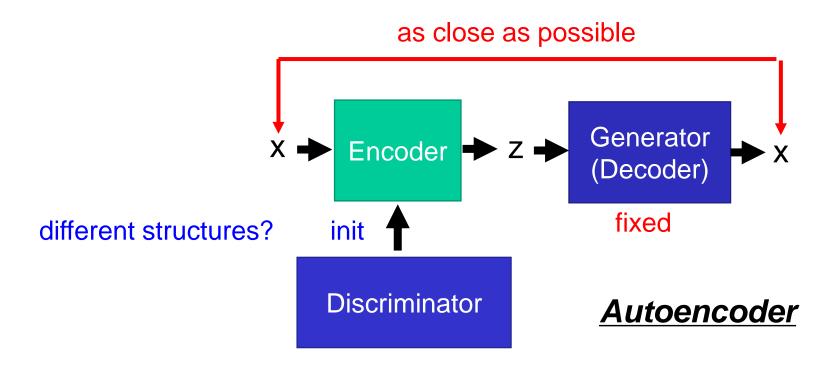
Modifying Input Code



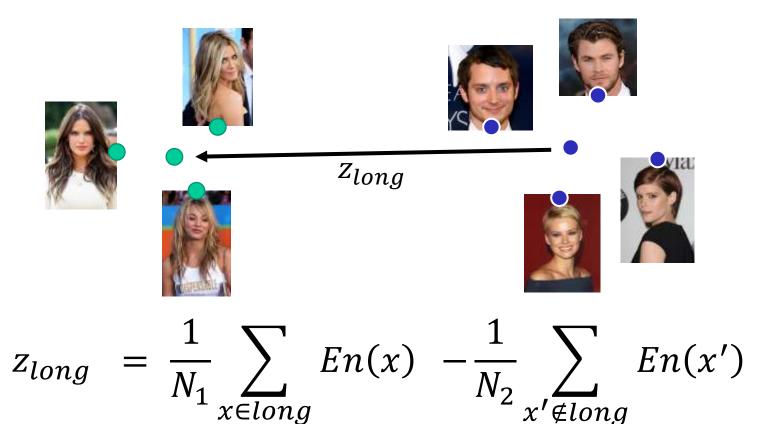
- > The input code determines the generator output.
- Understand the meaning of each dimension to control the output.

GAN + Autoencoder

- We have a generator (input z, output x)
- However, given x, how can we find z?
 - □ Learn an encoder (input x, output z)



GAN + AE: attribute representation



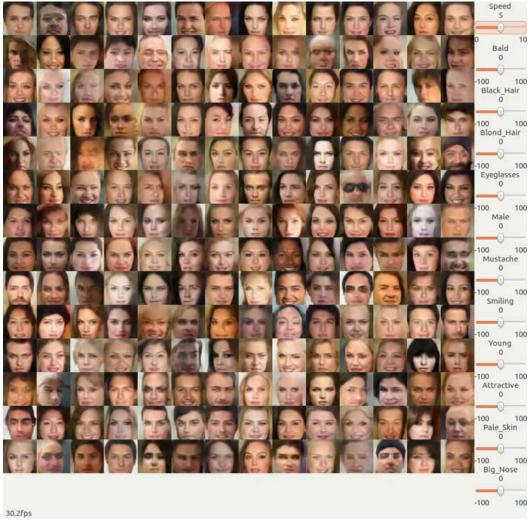
Short Hair

$$x \implies En(x) + z_{long} = z' \implies Gen(z')$$

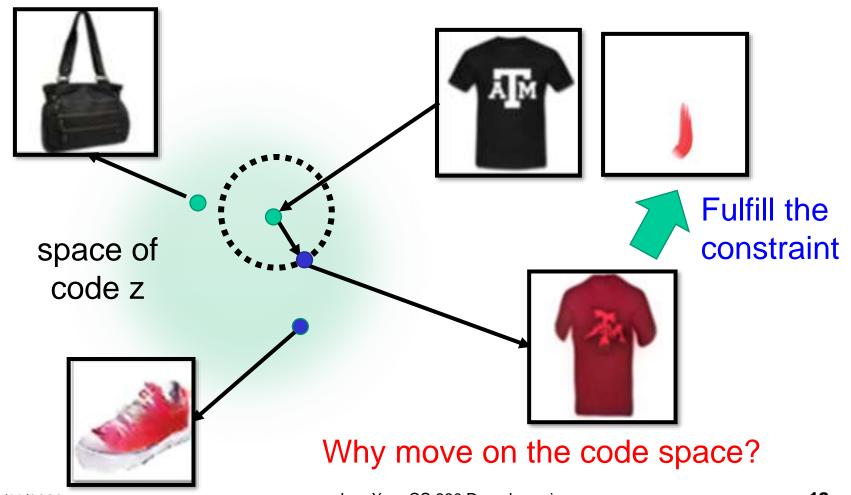
Long Hair

GAN + AE: attribute control

Photo editing via attribute control



Attribute editing: basic idea





Back to z

Method 1



$$z^* = arg \min_{z} \underline{L(G(z), x^T)}$$

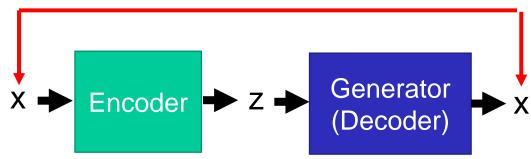
Gradient Descent

Difference between G(z) and x^T

- Pixel-wise
- By another network

as close as possible

Method 2



Method 3

Using the results from *method 2* as the initialization of *method 1*

Editing Photos





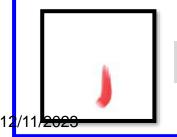
 \mathbf{z}_0 is the code of the input image

image

Using discriminator to check the image is realistic or not

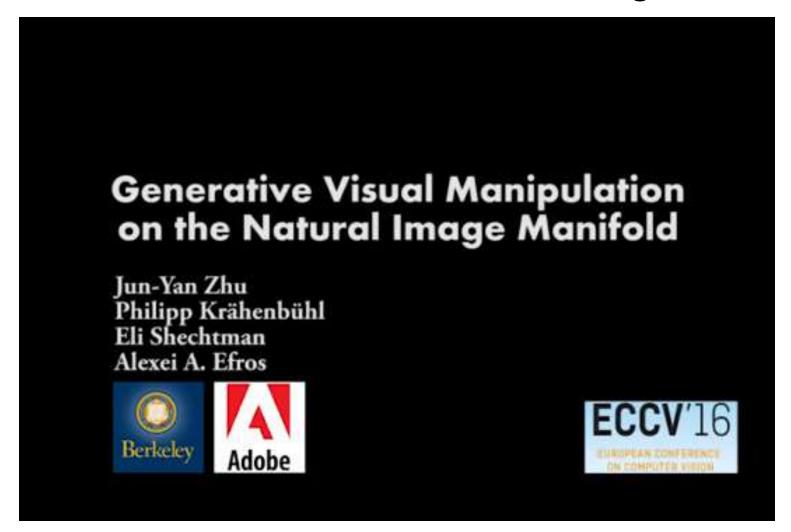
$$z^* = arg \min_{z} \underline{U(G(z))} + \lambda_1 \underline{||z - z_0||^2} - \lambda_2 \underline{D(G(z))}$$

Not too far away from the original image



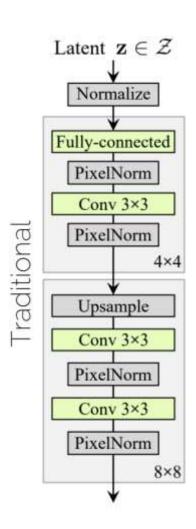
Does it fulfill the constraint of editing?

GAN + AE: attribute editing



[Zhu et al. 2016] Generative Visual Manipulation on the Natural Image Manifold

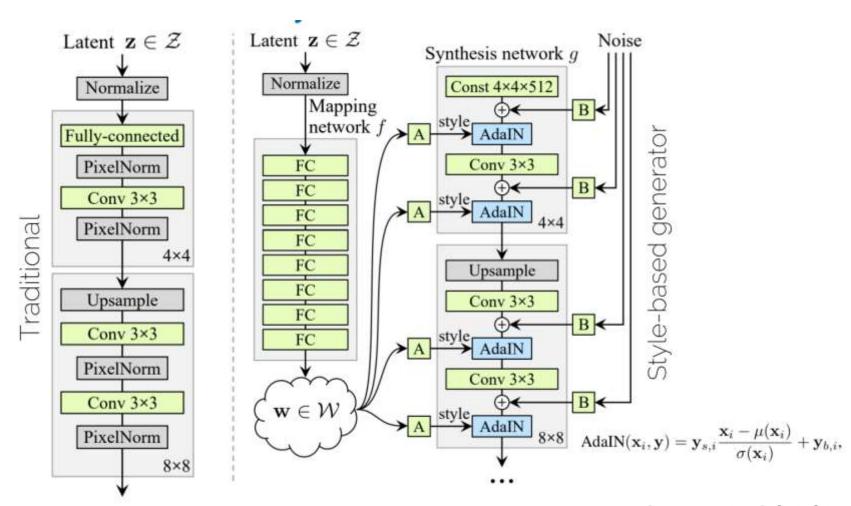




[Karras et al. 19]: StyleGAN



StyleGAN



[Karras et al. 19]: StyleGAN



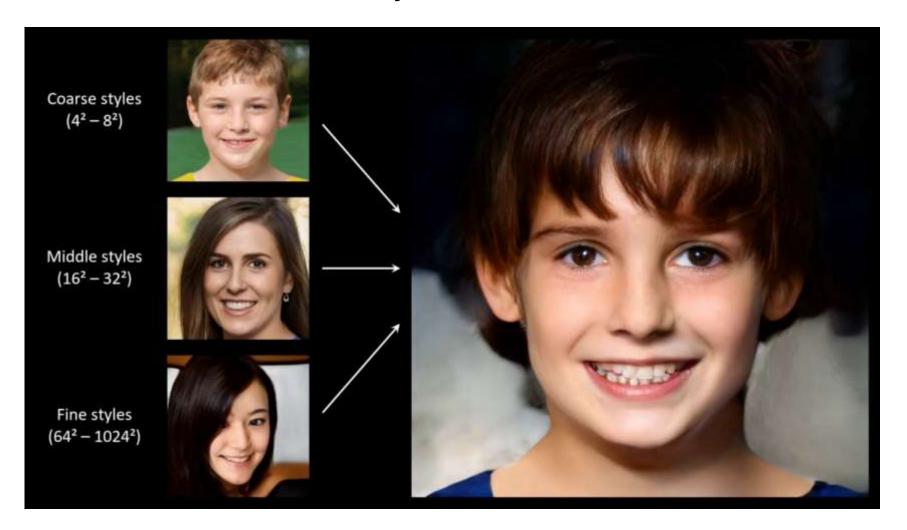
StyleGAN

- FID on 50k gen. images
- Architecture is similar to Progressive Growing GAN

Method	CelebA-HQ	FFHQ
A Baseline Progressive GAN [30]	7.79	8.04
B + Tuning (incl. bilinear up/down)	6.11	5.25
C + Add mapping and styles	5.34	4.85
D + Remove traditional input	5.07	4.88
E + Add noise inputs	5.06	4.42
F + Mixing regularization	5.17	4.40

[Karras et al. 19]: StyleGAN

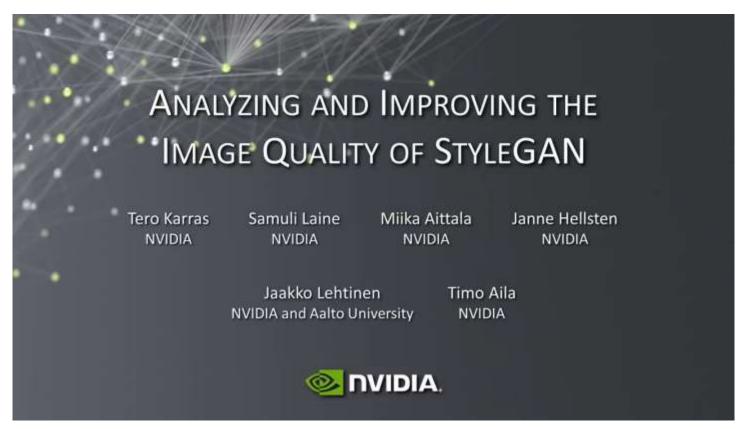
StyleGAN



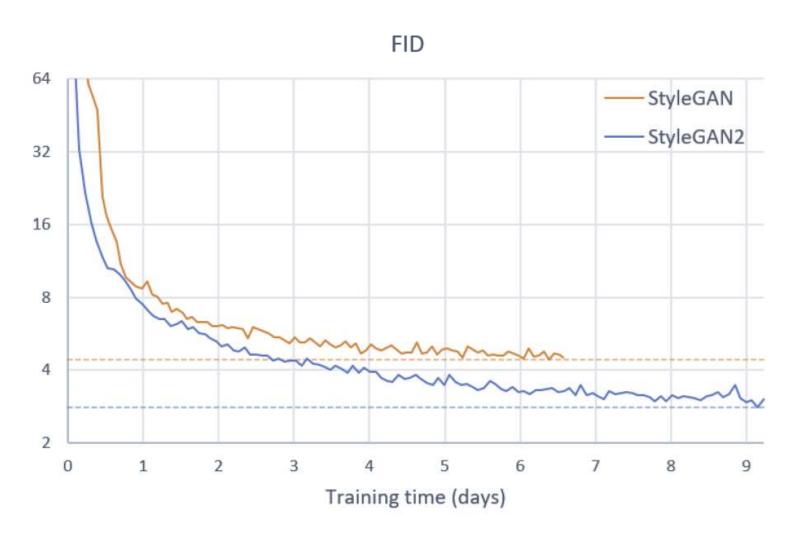
[Karras et al. 19]: StyleGAN



- Interesting analysis about design choices!
 - https://arxiv.org/pdf/1912.04958.pdf
 - https://github.com/NVlabs/stylegan2

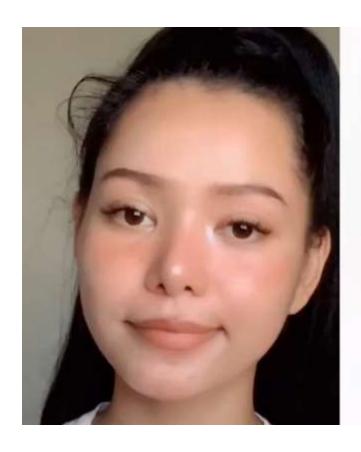






StyleGAN for style-transfer

- Keep basic and medium levels
- Fine-tune it on the fine leves



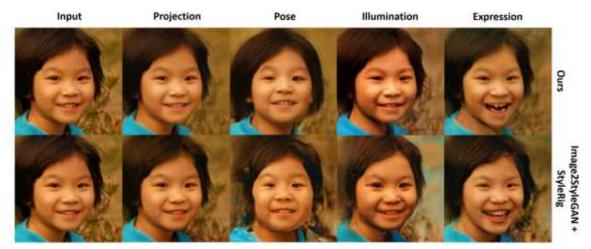


Latent space of StyleGAN

- StyleRig: Rigging StyleGAN for 3D Control over Portrait Images
 - □ Face rig-like control over StyleGAN

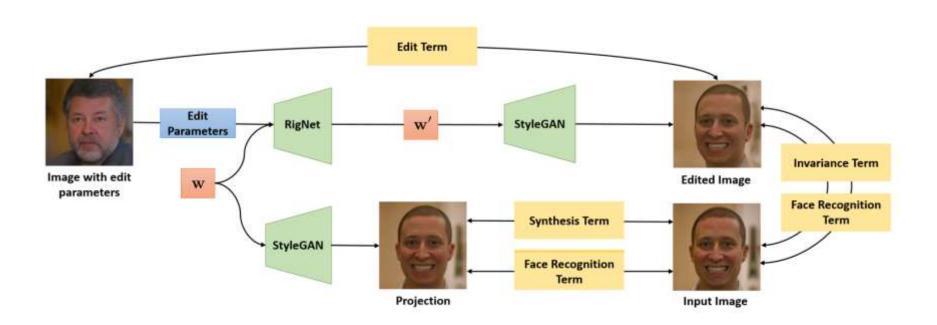


- PIE: Portrait Image Embedding for Semantic Control
 - Embedding in the latent space to edit real images



PIE: Portrait Image Embedding

Utilize pretrained StyleGAN



[TEWARI et al. 19] PIE: Portrait Image Embedding for Semantic Control

PIE: Portrait Image Embedding





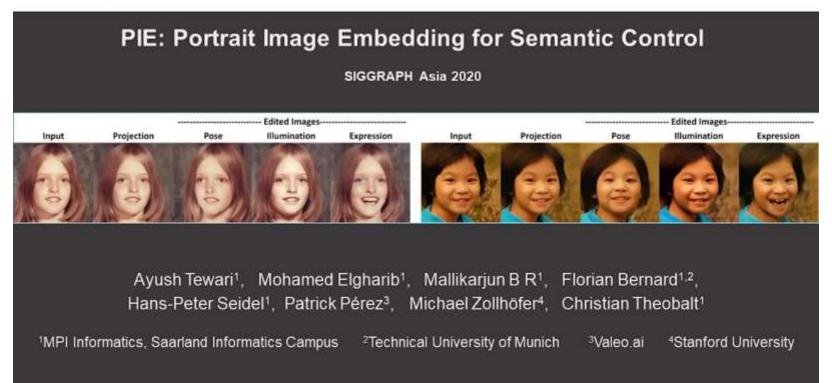












[TEWARI et al. 19] PIE: Portrait Image Embedding for Semantic Control



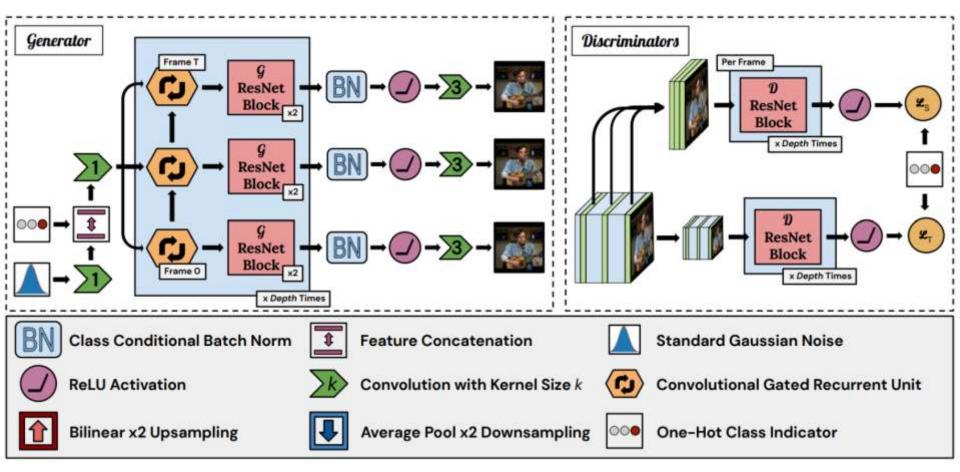
GANs on Videos

Two options

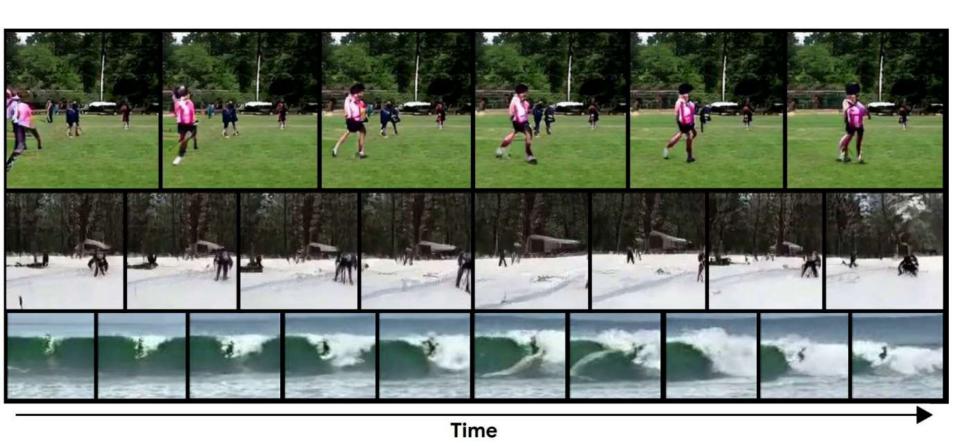
- □ Single random variable z seeds entire video (all frames)
 - Very high dimensional output
 - How to do for variable length?
 - Future frames deterministic given past
- Random variable z for each frame of the video
 - Need conditioning for future from the past
 - How to get combination of past frames + random vectors during training

General issues

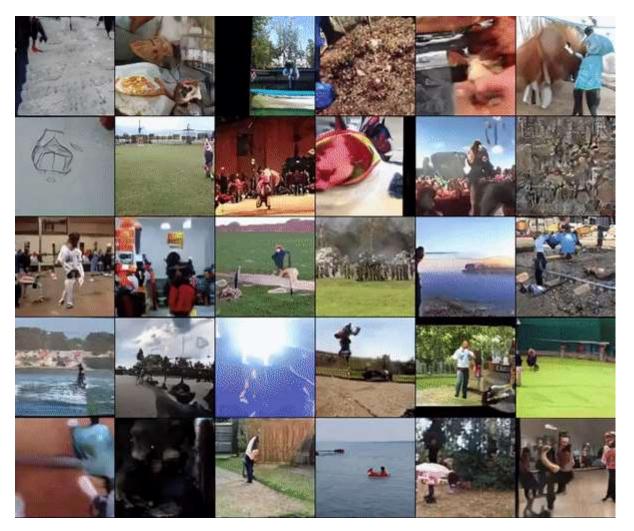
- □ Temporal coherency
- Drift over time (many models collapse to mean image)



[Clark et al. 2019] Adversarial Video Generation on Complex Datasets



[Clark et al. 2019] Adversarial Video Generation on Complex Datasets



[Clark et al. 2019] Adversarial Video Generation on Complex Datasets



- Trained on Kinetics-600 dataset
 - □ 256 x 256, 128 x 128, and 64 x 64
 - □ Lengths of up 48 frames
- This is state of the art!
- Videos from scratch still incredibly challenging

Conditional GANs on Videos

Challenge:

□ Each frame is high quality, but temporally inconsistent



и-

Video-to-Video Synthesis

Sequential Generator:

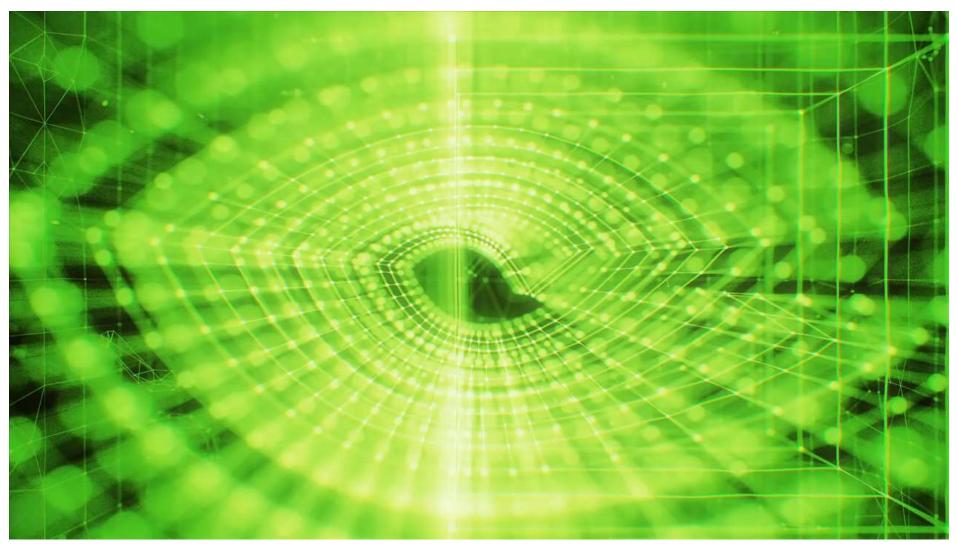
$$p(\tilde{\mathbf{x}}_1^T|\mathbf{s}_1^T) = \prod_{t=1}^T p(\tilde{\mathbf{x}}_t|\tilde{\mathbf{x}}_{t-L}^{t-1},\mathbf{s}_{t-L}^t).$$
 past L generated frames past L source frames (Set L = 2)

- lacktriangle Conditional Image Discriminator D_i (is it real image)
- lacktriangle Conditional Video Discriminator D_{v} (temp. consistency via flow)
- Full Learning Objective:

$$\min_{F} \left(\max_{D_I} \mathcal{L}_I(F, D_I) + \max_{D_V} \mathcal{L}_V(F, D_V) \right) + \lambda_W \mathcal{L}_W(F)$$

[Wang et al. 2018] Vid2Vid

Video-to-Video Synthesis



[Wang et al. 2018] Vid2Vid

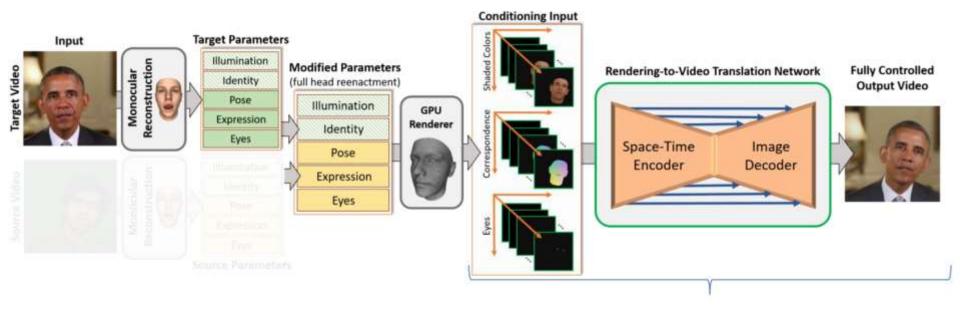


Key ideas:

- Separate discriminator for temporal parts
 - □ In this case based on optical flow
- Consider recent history of previous frames

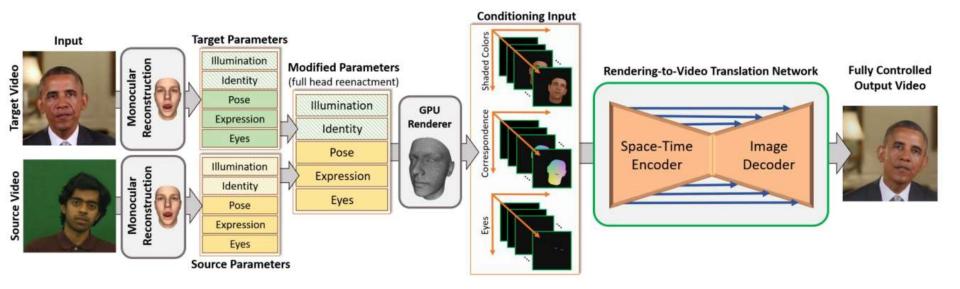
Train all of it jointly

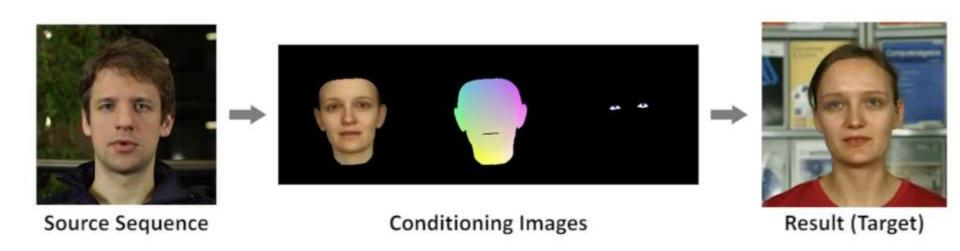
Deep Video Portraits



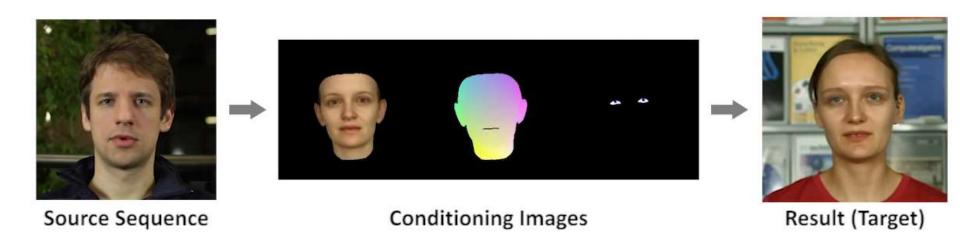
Similar to "Image-to-Image Translation" (Pix2Pix) [Isola et al.]

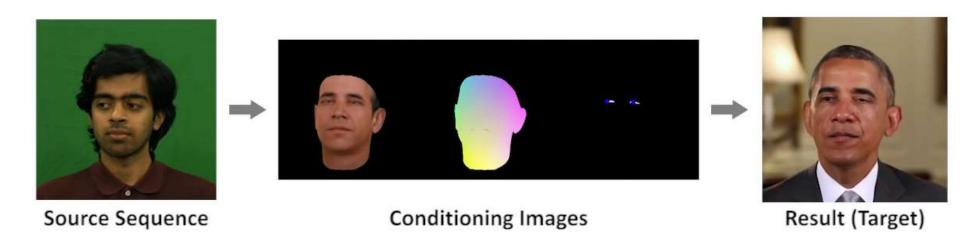
Siggraph'18 [Kim et al 2018]: Deep Portraits





Neural Network converts synthetic data to realistic video







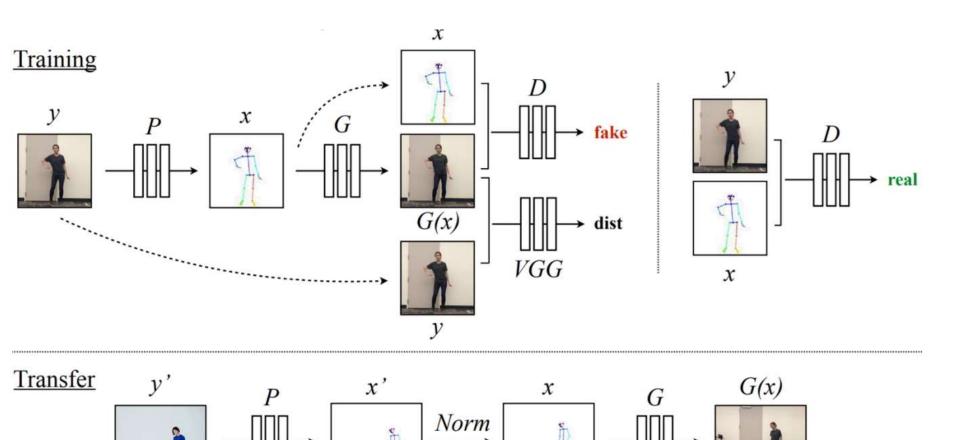






Key ideas:

- Synthetic data for tracking is great anchor / stabilizer
- Overfitting on small datasets works pretty well
- Need to stay within training set w.r.t. motions
- No real learning; essentially, optimizing the problem with SGD



[Chan et al. '18] Everybody Dance Now

Everybody Dance Now

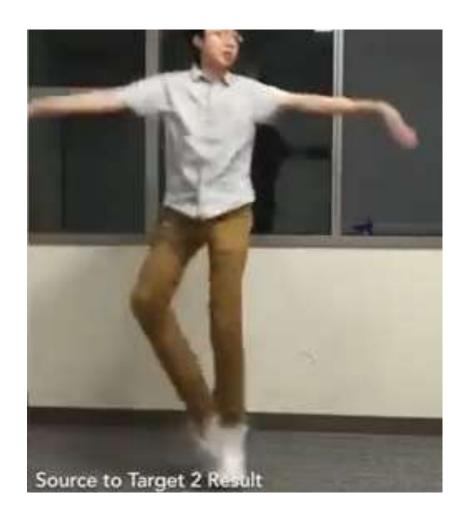
Motion Retargeting Video Subjects

Caroline Chan, Shiry Ginosar, Tinghui Zhou, Alexei A. Efros

UC Berkeley

[Chan et al. '18] Everybody Dance Now

- cGANs work with different input
- Requires consistent input i.e., accurate tracking
- Network has no explicit3D notion



[Chan et al. '18] Everybody Dance Now



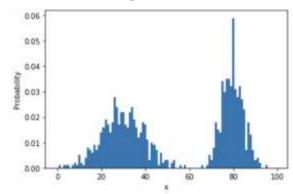
Key ideas:

- Conditioning via tracking seems promising!
- Tracking quality translates to resulting image quality
- Tracking human skeletons is less developed than faces
 - □ Temporally it's not stable... (e.g., OpenPose etc.)
- Fun fact, there were like 4 papers with a similar idea that appeared around the same time... (even more papers recently)

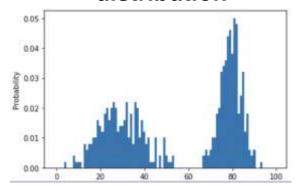
Recap of generative modeling so far

- Assume that all data comes from a data distribution that is unknown to us.
- The goal of deep generative modeling is to train neural networks that can learn this distribution from its samples (i.e., a fixed size training dataset)
- We generate new data by sampling from the learned distribution
- In practice, models are trained to maximize expected log-likelihood or minimize divergence between the learned and empirical data distributions

Training distribution



Samples from learned distribution





Prior methods

VAE

$$egin{aligned} L_{ ext{VAE}}(heta,\phi) &= -\log p_{ heta}(\mathbf{x}) + D_{ ext{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{ heta}(\mathbf{z}|\mathbf{x})) \ &= -\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} \log p_{ heta}(\mathbf{x}|\mathbf{z}) + D_{ ext{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{ heta}(\mathbf{z})) \ heta^*, \phi^* &= rg\min_{ heta,\phi} L_{ ext{VAE}} \end{aligned}$$

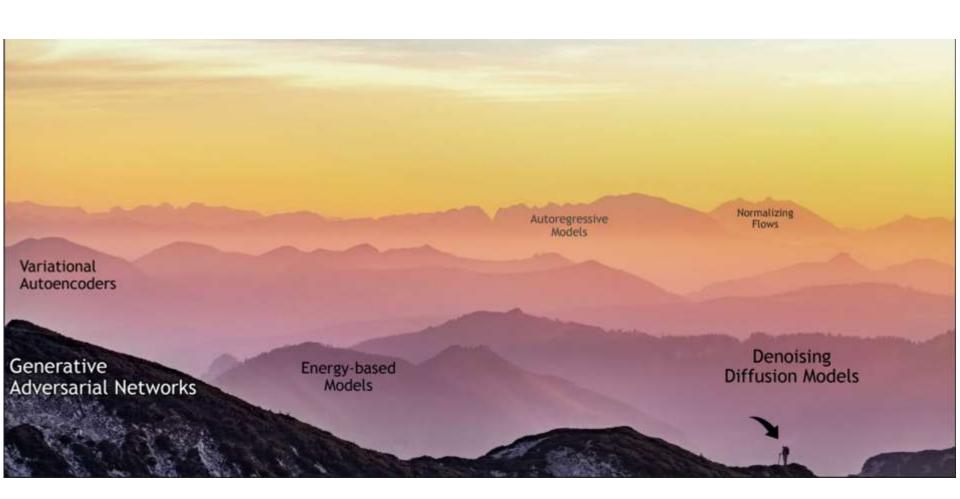
$$-L_{ ext{VAE}} = \log p_{ heta}(\mathbf{x}) - D_{ ext{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{ heta}(\mathbf{z}|\mathbf{x})) \leq \log p_{ heta}(\mathbf{x})$$

GAN

$$egin{aligned} \min_G \max_D L(D,G) &= \mathbb{E}_{x\sim p_r(x)}[\log D(x)] + \mathbb{E}_{z\sim p_z(z)}[\log(1-D(G(z)))] \ &= \mathbb{E}_{x\sim p_r(x)}[\log D(x)] + \mathbb{E}_{x\sim p_g(x)}[\log(1-D(x))] \end{aligned}$$

$$L(G, D^*) = 2 D_{JS}(p_r \| p_g) - 2 \log 2$$

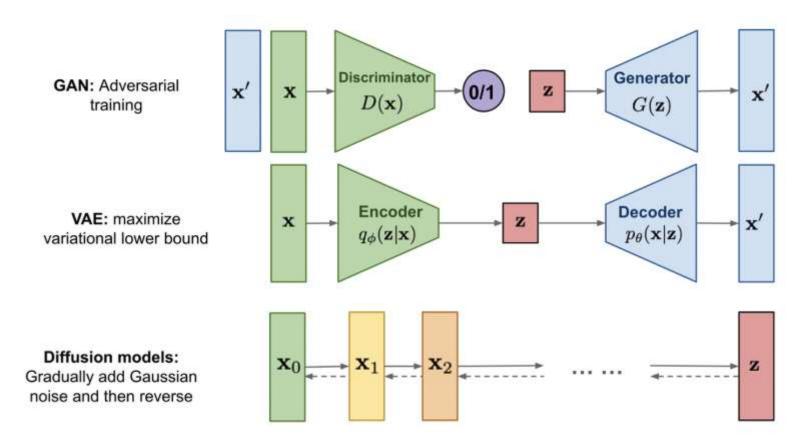
The Landscape of Generative Models





Sampling from noise

A quick paradigm



What if we add a bunch of Gaussian noise to an image?



and again...



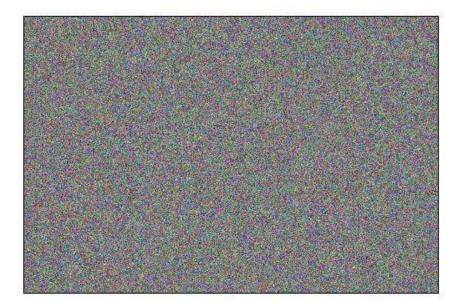
and again...



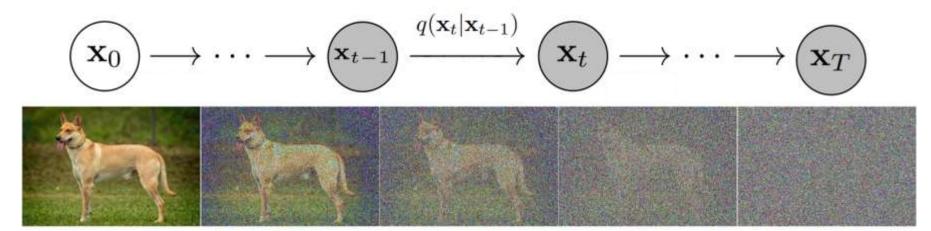
and again...



... until it resembles pure noise

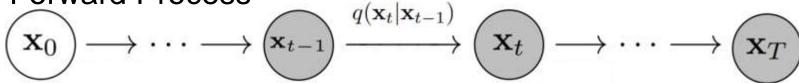


Forward Process



- Take a datapoint x_0 and gradually add very small amounts of Gaussian noise to it
- Let x_t be the datapoint after t iterations
- This is called the forward diffusion process
- Repeat this process for T steps over time, more and more features of the original input are destroyed until you get something resembling pure noise

Forward Process







More formally, we update each image over time as

$$x_t = \sqrt{1-eta_t} x_{t-1} + \sqrt{eta_t} \epsilon ~~ ext{where}~~\epsilon \sim \mathcal{N}(0,\,I)$$

where

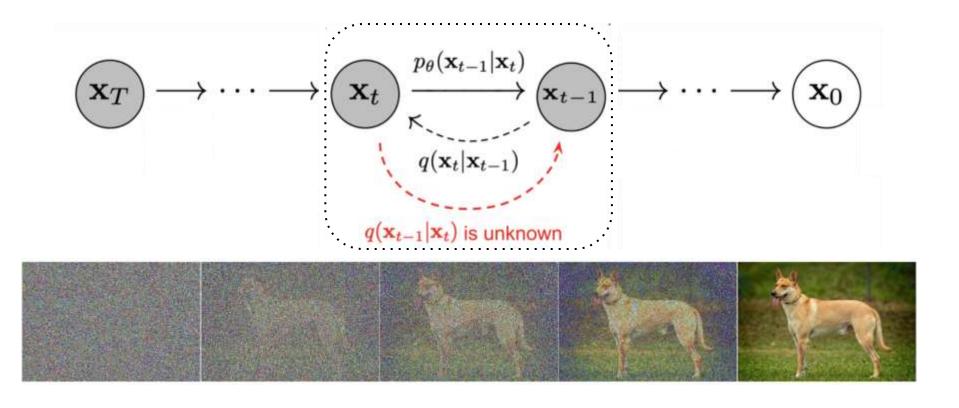
$$\{eta_t \in (0,1)\}_{t=1}^T$$

is called the **noise schedule** (basically a hyperparameter describing how much noise to add at a given timestep).

The update above can equivalently be written as a sampling process from the following Gaussian distribution:

$$q(x_t \mid x_{t-1}) \sim \mathcal{N}(\sqrt{1-eta_t}x_{t-1},\,eta_t I)$$

Can we go in the other direction?



58

Reverse Process

Fixed Forward Process

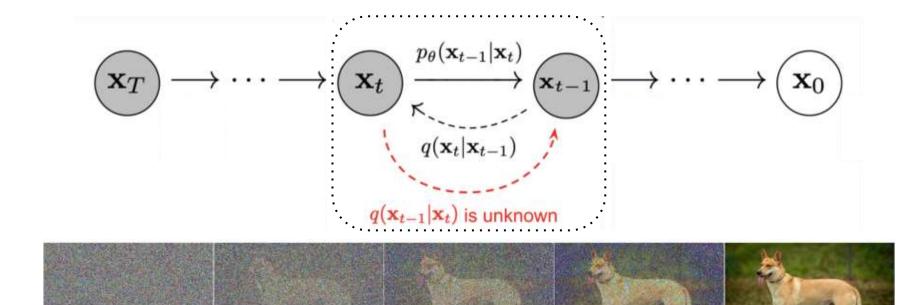


Learned Reverse Process

- The goal of a diffusion model is to learn the reverse denoising process to iteratively undo the forward process
- In this way, the reverse process appears as if it is generating new data from random noise!

Reverse Process

We are given $q(x_t | x_{t-1})$. How do we find $q(x_{t-1} | x_t)$?





- Finding the exact distribution is hard
- Remember Bayes rule?

$$f(\theta \mid x) = \frac{f(\theta, x)}{f(x)} = \frac{f(\theta) f(x \mid \theta)}{f(x)}$$

$$q(x_{t-1} \mid x_t) = q(x_t \mid x_{t-1}) \frac{q(x_{t-1})}{q(x_t)}$$

$$q(x_t) = \int q(x_t \mid x_{t-1}) q(x_{t-1}) dx$$

- This is computationally intractable
 - Need to integrate over the whole data distribution to find q(x_t) and q(x_{t-1})
 - Where else have we seen this dilemma?
- We still need the posterior distribution to carry out the reverse process. Can we approximate this somehow?



- What does the final reverse process look like?
- In practice, we choose our noise schedule such that the forward process steps are very small. Thus, we approximate the reverse posterior distributions as Gaussians and learn their parameters (i.e., the mean and variance) via neural networks

$$p_{ heta}(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1};oldsymbol{\mu}_{ heta}(\mathbf{x}_t,t),oldsymbol{\Sigma}_{ heta}(\mathbf{x}_t,t))$$

Cool, we now have an idea of what the model looks like. How do we train it?



A preliminary objective

We want to maximize the log-likelihood of the data generated by a reverse process.

Remember that VAEs tried to do something similar but they maximized a lower bound on the likelihood instead because the actual likelihood is computationally intractable

$$-L_{ ext{VAE}} = \log p_{ heta}(\mathbf{x}) - D_{ ext{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \| p_{ heta}(\mathbf{z}|\mathbf{x})) \leq \log p_{ heta}(\mathbf{x})$$

We can apply the same trick to diffusion!

$$-L = \mathbb{E}_{q(\mathbf{x}_{0:T})}igg[\lograc{p_{ heta}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}\,|\,\mathbf{x}_0)}igg] \leq \log p_{ heta}(\mathbf{x}_0)$$



A simplified objective

Recall the forward process:

$$x_t = \sqrt{1-eta_t} x_{t-1} + \sqrt{eta_t} \epsilon ~~ ext{where}~~\epsilon \sim \mathcal{N}(0,\,I)$$

Skipping the derivation here, but you can show that

our estimate
for x_{t-1}

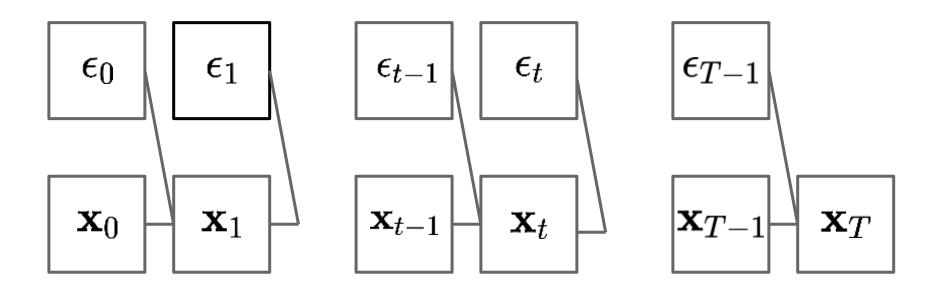
where

$$L \, = \, \mathbb{E}_{\mathbf{x}_0,\epsilon} \Bigg[w_t igg \| rac{1}{\sqrt{lpha_t}} igg(\mathbf{x}_t - rac{eta_t}{\sqrt{1-ar{lpha}_t}} \epsilon igg) - \mu_ heta(\mathbf{x}_t,\,t) igg \|_2^2 \Bigg]$$

$$lpha_t = 1 - eta_t \;\; ext{and} \;\; ar{lpha}_t = \prod_{i=1}^t lpha_t$$



- A simplified objective
- Forward process re-visualized





An even more simplified objective!

Instead of predicting an estimate for x_{t-1} , it's actually better to predict the noise ϵ that was added to x_{t-1} to produce x_t .

In other words,

$$x_t = \sqrt{1-eta_t} \, x_{t-1} + \sqrt{eta_t} \, \epsilon$$
 where $\epsilon \sim \mathcal{N}(0,\,I)$ Don't predict this Predict this instead

Thus, our final loss becomes

$$egin{aligned} L &= \mathbb{E}_{\mathbf{x}_0,\,\epsilon} \Big[\|\epsilon - \epsilon_{ heta}(\mathbf{x}_t,\,t)\|_2^2 \Big] \ &= \mathbb{E}_{\mathbf{x}_0,\,\epsilon} \Big[\Big\|\epsilon - \epsilon_{ heta} \Big(\sqrt{ar{lpha}}_t \mathbf{x}_0 + \sqrt{1 - ar{lpha}_t} \epsilon,\,t \Big) \Big\|_2^2 \Big] \end{aligned}$$



Training

Algorithm 1 Training

```
1: repeat
```

- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1,\ldots,T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on

$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$

6: until converged



Sampling

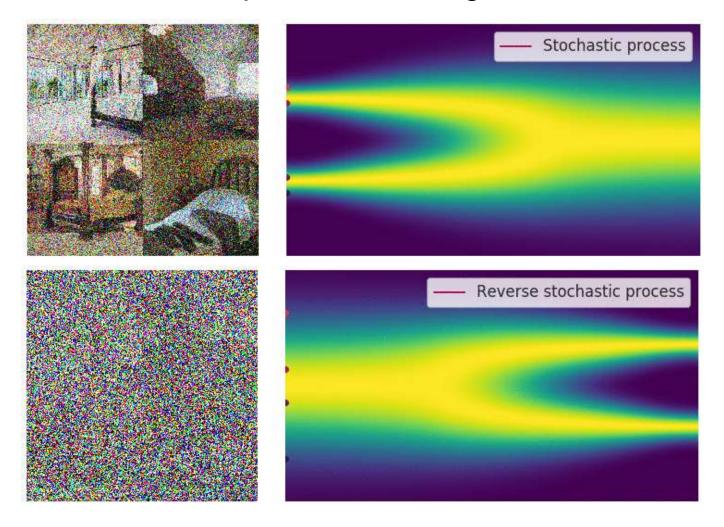
Algorithm 2 Sampling

- 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for** t = T, ..., 1 **do**
- 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if t > 1, else $\mathbf{z} = \mathbf{0}$

4:
$$\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$$

- 5: end for
- 6: return \mathbf{x}_0

Forward/Reverse process for Image Generation





Summary

- Variants of GANs
 - CycleGAN: Image-to-image translation with unpaired data
 - ☐ GANs on videos
- Diffusion, very quick start
- Keep working on the projects!