

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Autumn 2023

Generative Adversarial Networks (GANs)

Modeling Probability Distributions on Images

Suppose we want to train a model of the probability distribution of natural images using cross-entropy loss.

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim \text{pop}} - \ln p_{\Phi}(y)$$

Images are continuous structured objects — a continuous value at every pixel.

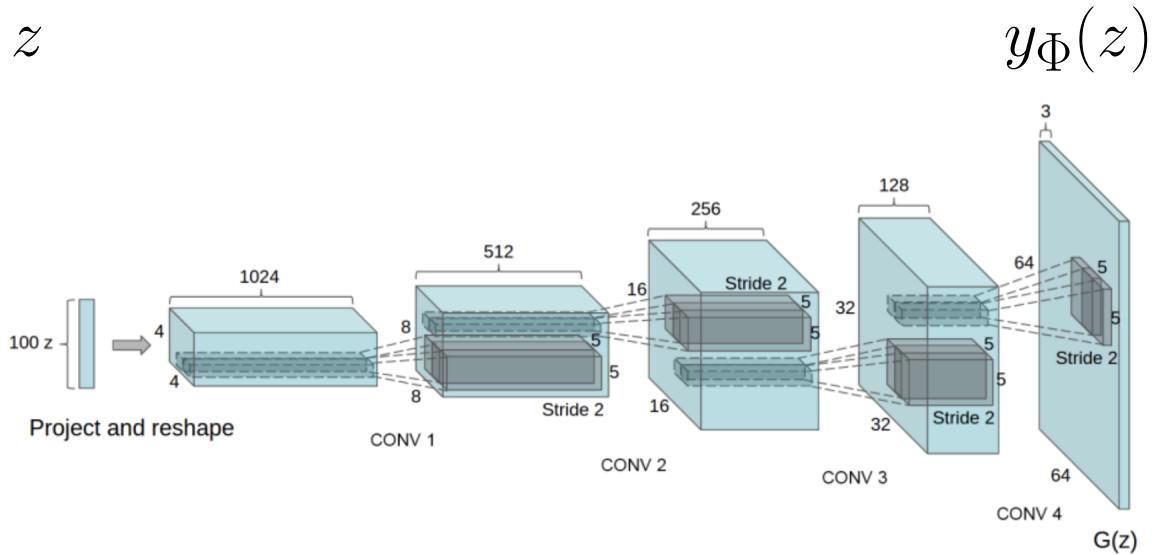
It is difficult to build probability models for images (or other continuous structured values) that both accurately model the distribution and also allow us to calculate $p_{\Phi}(y)$.

Generative Adversarial Networks (GANs)

GANs represent $p_\Phi(y)$ implicitly by constructing an image generator and abandon the ability to compute $p_\Phi(y)$.

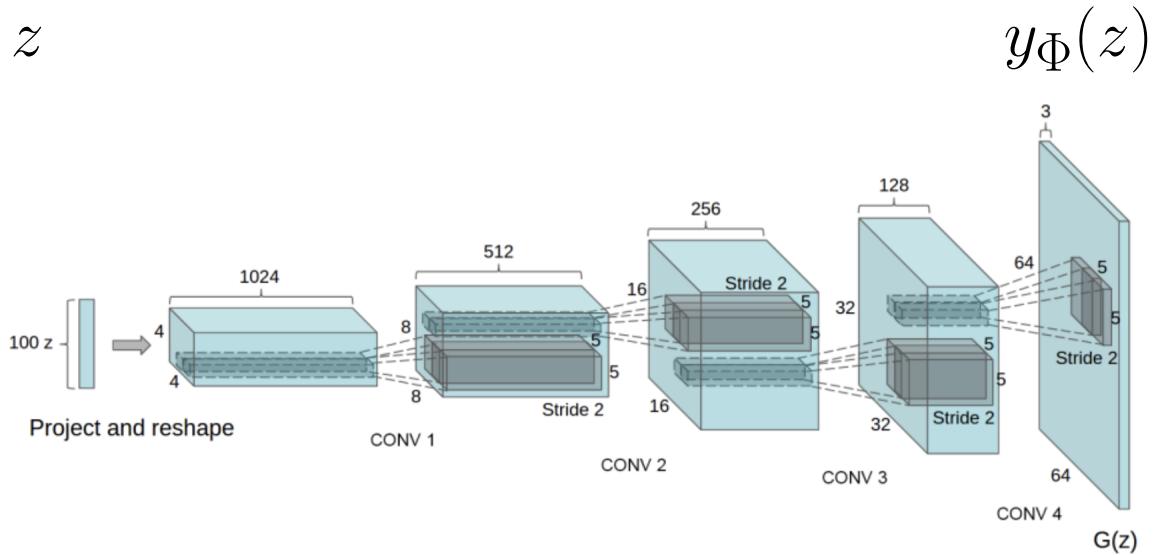
The cross-entropy loss is replaced by an adversarial discriminator which tries to distinguish between generated images and real images.

Representing a Distribution with a Generator



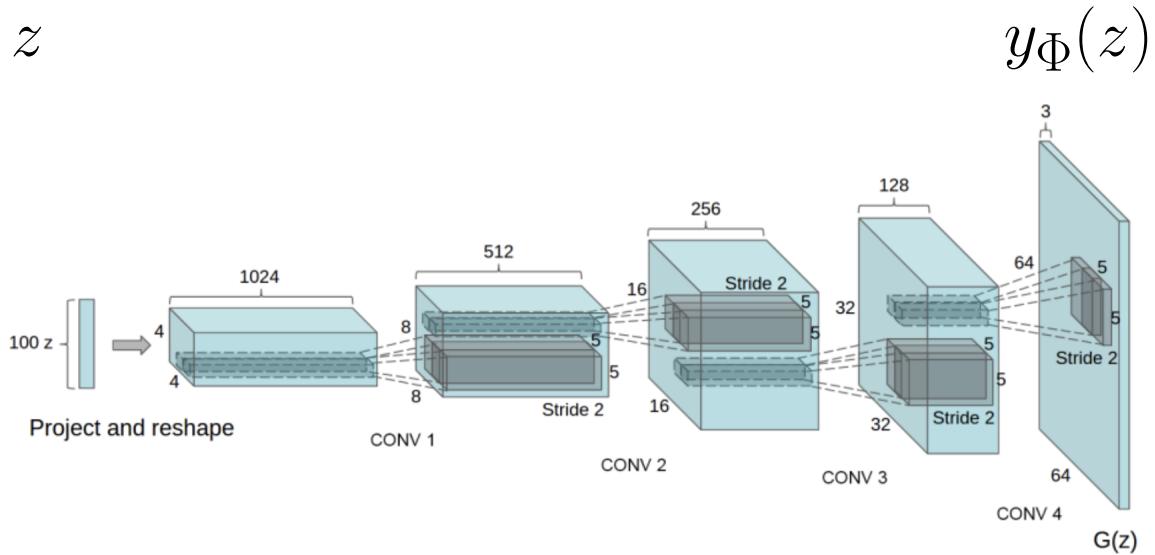
The random input z defines a probability density on images $y_\Phi(z)$. We will write this as $p_\Phi(y)$ for the image y .

Representing a Distribution with a Generator



We want $p_\Phi(y)$ to model a natural image distribution such as the distribution over human faces.

Representing a Distribution with a Generator



We can sample from $p_\Phi(y)$ by sampling z . But we cannot compute $p_\Phi(y)$ for y sampled from the population.

Increasing Spatial Dimension

Reducing spatial dimension with strided convolution:

For $x, y, j, \Delta x, \Delta y, i$

$$L_{\ell+1}[\textcolor{red}{x}, \textcolor{red}{y}, j] += W[\Delta x, \Delta y, i, j] L_\ell[\textcolor{red}{s} * x + \Delta x, \textcolor{red}{s} * y + \Delta y, i]$$

Increasing spatial dimension with PyTorch ConvTranspose2d:

For $x, y, j, \Delta x, \Delta y, i$

$$L_{\ell+1}[\textcolor{red}{s} * x + \Delta x, \textcolor{red}{s} * y + \Delta y, i] += W[\Delta x, \Delta y, i, j] L_\ell[\textcolor{red}{x}, \textcolor{red}{y}, j]$$

Generative Adversarial Networks (GANs)

Let y range over images. We have a generator p_Φ . For $i \in \{-1, 1\}$ we define a probability distribution over pairs $\langle i, y \rangle$ by

$$\begin{aligned}\tilde{p}_\Phi(i = 1) &= 1/2 \\ \tilde{p}_\Phi(y|i = 1) &= \text{pop}(y) \\ \tilde{p}_\Phi(y|i = -1) &= p_\Phi(y)\end{aligned}$$

We also have a discriminator $P_{\text{disc}}(i|y)$ that tries to determine the source i given the image y .

The generator tries to fool the discriminator.

$$\text{gen}^* = \underset{\text{gen}}{\operatorname{argmax}} \underset{\text{disc}}{\min} E_{\langle i, y \rangle \sim \tilde{p}_{\text{gen}}} - \ln P_{\text{disc}}(i|y)$$

GANs

The generator tries to fool the discriminator.

$$\text{gen}^* = \underset{\text{gen}}{\operatorname{argmax}} \min_{\text{disc}} E_{\langle i, y \rangle \sim \tilde{p}_{\text{gen}}} - \ln P_{\text{disc}}(i|y)$$

Assuming universality (next slide) of both the generator p_{gen} and the discriminator P_{disc} we have $p_{\text{gen}^*} = \text{pop}$.

Note that this involves only discrete cross-entropy.

The Universality Assumption

DNNs are universally expressive (can model any function) and trainable (the desired function can be found by SGD).

Universality assumption is clearly false but is useful.

The success of GANs (to the extent that they have been successful) is a tribute to the utility of the universality assumption.

Jensen-Shannon Divergence

$$\text{gen}^* = \underset{\text{gen}}{\operatorname{argmax}} \min_{\text{disc}} E \langle i, y \rangle_{\sim \tilde{p}_{\text{gen}}} - \ln P_{\text{disc}}(i|y)$$

$$= \underset{\text{gen}}{\operatorname{argmin}} KL \left(\text{pop}, \frac{\text{pop} + p_{\text{gen}}}{2} \right) + KL \left(p_{\text{gen}}, \frac{\text{pop} + p_{\text{gen}}}{2} \right)$$

$$\text{gen}^* = \underset{\text{gen}}{\operatorname{argmin}} \text{JSD}(\text{pop}, p_{\text{gen}})$$

GAN Mode Collapse

A major concern is “mode collapse” where the learned distribution omits a significant fraction of the population distribution.

There is no quantitative performance measure that provides a meaningful guarantee against mode collapse.

In practice GANS are evaluated on FID score.

The Fréchet Inception Score (FID)

Consider two distributions P and Q on the same set (perhaps distributions or densities on images).

Generative image models are (still) evaluated using a certain measure of the “distance” between the population distribution and the generation distribution.

For GANs we cannot compute the probability of a generated image so we cannot use cross entropy or KL-divergence.

The Fréchet Inception Score (FID)

The Fréchet distance $F(P, Q)$ can be measured purely by sampling.

$$F(p, q) = \inf_{\mu} \left(E_{(x,y) \sim \mu} ||x - y||^2 \right)^{\frac{1}{2}}$$

This is also known as the 2-Wasserstein distance.

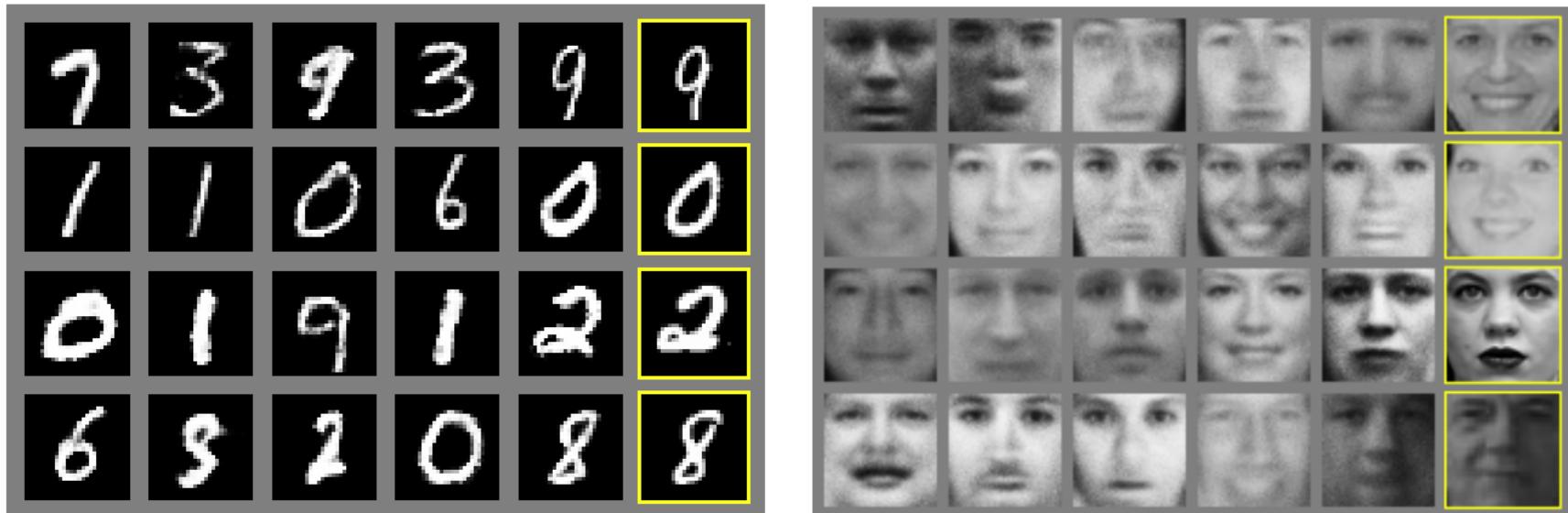
The Fréchet Inception Score (FID)

But rather than measure the L_2 distance between images we measure the L_2 distance between the “inception feature vectors” $I(x)$ and $I(y)$.

For an image x the feature vector $I(x)$ is computed from a certain layer in the inception image classification network.

Generative Adversarial Nets

Goodfellow et al., June 2014



The rightmost column (yellow boarders) gives the nearest neighbor in the training data to the adjacent column.

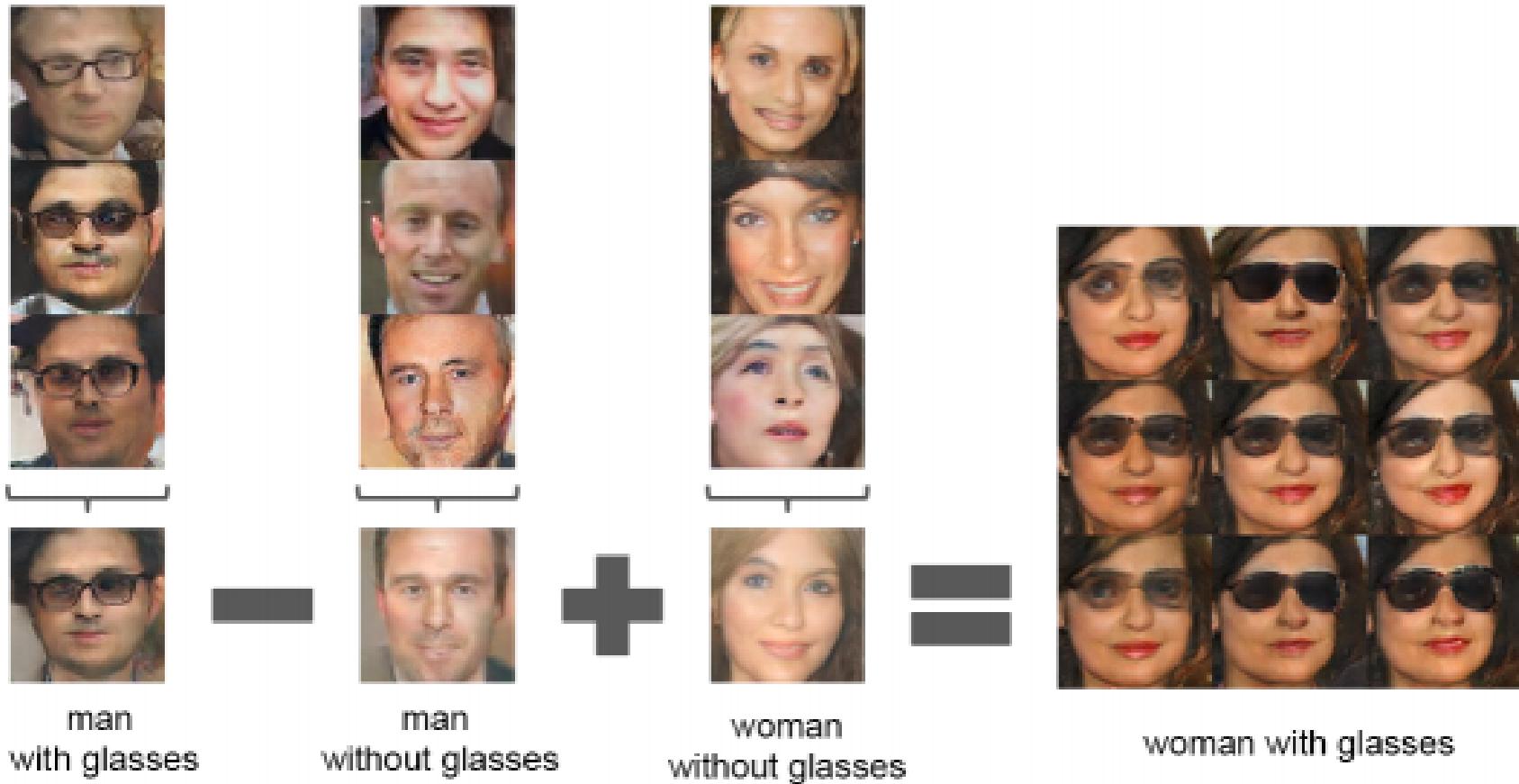
Unsupervised Representation Learning ... (DC GANS)

Radford et al., Nov. 2015



Unsupervised Representation Learning ... (DC GANS)

Radford et al., Nov. 2015



Interpolated Faces

[Ayan Chakrabarti, January 2017]



Conditional GANS

In the conditional case we have a population distribution over pairs $\langle x, y \rangle$.

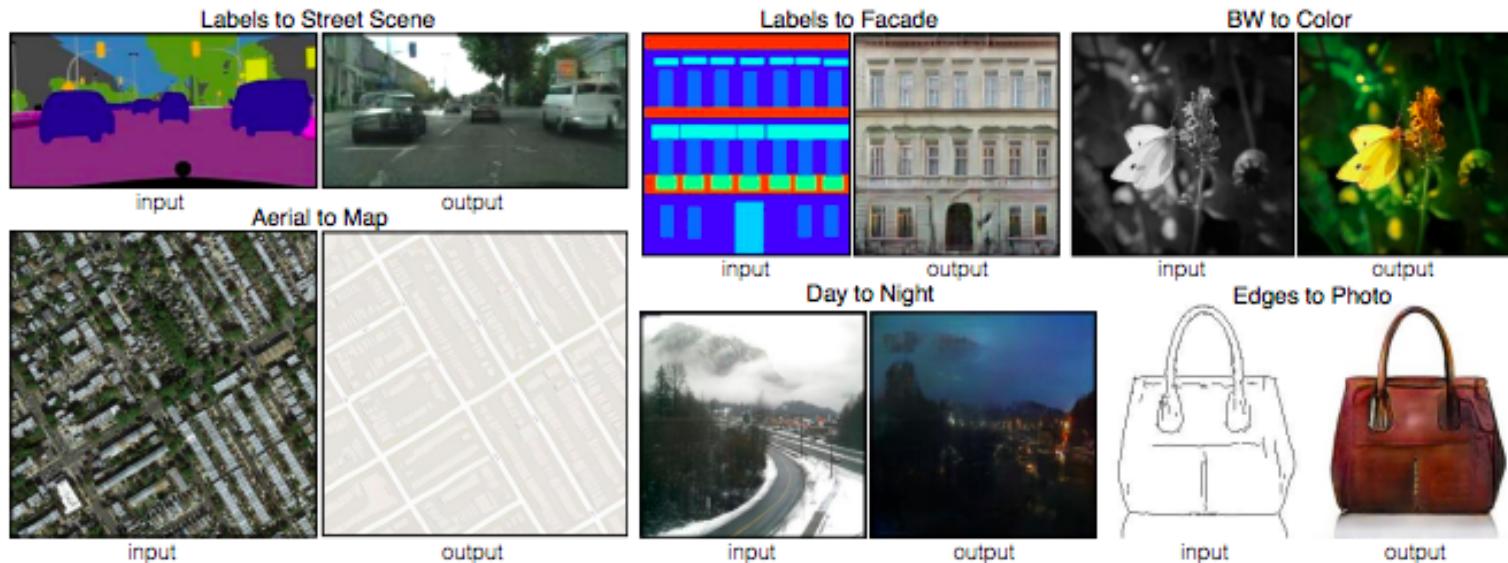
For conditional GANs we have a generator $p_{\text{gen}}(y|x)$ and a discriminator $P_{\text{disc}}(i|x, y)$ where $i = 1$ if y is the real “label” for x and -1 if y is generated from x .

$$\text{gen}^* = \underset{\text{gen}}{\operatorname{argmax}} \underset{\text{disc}}{\min} E_{\langle x, y, i \rangle \sim \tilde{p}_{\text{gen}}} - \ln P_{\text{disc}}(i|x, y)$$

Image-to-Image Translation (Pix2Pix)

Isola et al., Nov. 2016

We assume a corpus of “image translation pairs” such as images paired with semantic segmentations.



UNets

Pix2Pix uses a U-Net.

U-Net: Convolutional Networks for Biomedical Image Segmentation, Ronneberger, Fischer and Brox, May 2015.

U-Nets are fundamental to current diffusion models such as DALL·E.

Adversarial Discrimination as an Additional Loss

$$\text{gen}^* = \underset{\text{gen}}{\operatorname{argmin}} \ E_{(x,y) \sim \text{pop}} \ ||y - y_{\text{gen}}(x)||_1 + \lambda \mathcal{L}_{\text{Discr}}(\text{gen})$$

$$\mathcal{L}_{\text{Discr}}(\text{gen}) = \max_{\text{disc}} \ E_{x,y,i \sim \tilde{p}_{\text{gen}}} \ \ln P_{\text{disc}}(i|y, x)$$

Discrimination as an Additional Loss

$$\text{L1 : } \text{gen}^* = \operatorname{argmin}_{\text{gen}} E_{(x,y) \sim \text{pop}} \|y - y_{\text{gen}}(x)\|_1$$

$$\text{cGAN : } \text{gen}^* = \operatorname{argmin}_{\text{gen}} \mathcal{L}_{\text{Discr}}(\text{gen})$$

$$\text{L1 + cGAN : } \text{gen}^* = \operatorname{argmin}_{\text{gen}} E_{(x,y) \sim \text{pop}} \|y - y_{\text{gen}}(x)\|_1 + \lambda \mathcal{L}_{\text{Discr}}(\text{gen})$$

Image-to-Image Translation (Pix2Pix)

Isola et al., Nov. 2016



Arial Photo to Map and Back

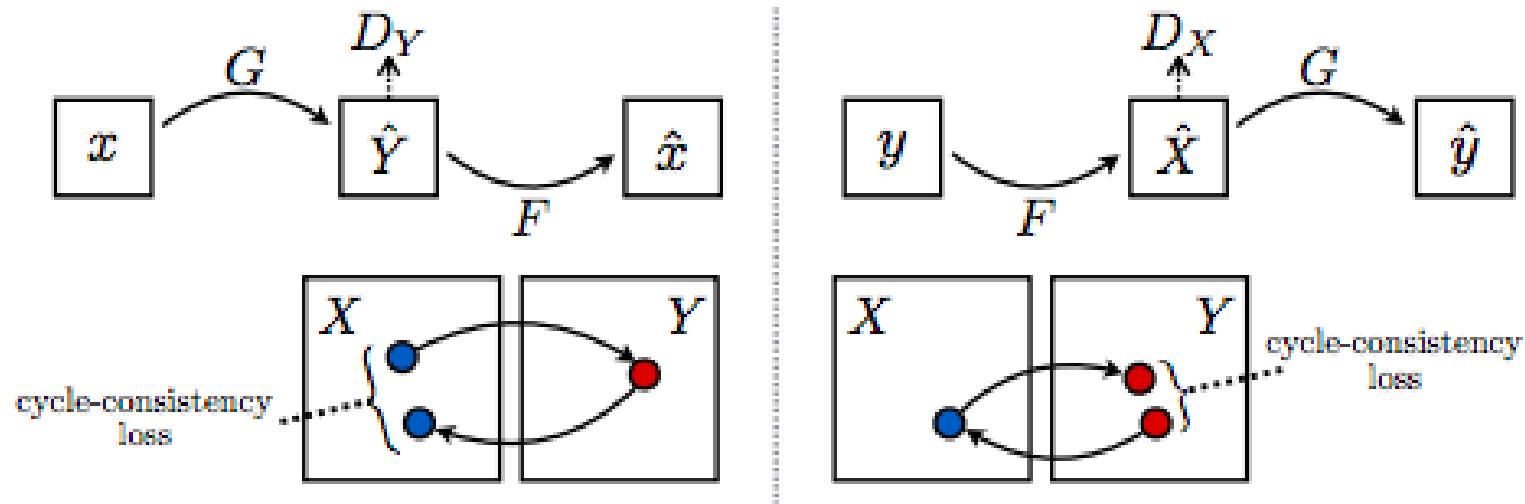


Unpaired Image-to-Image Translation (Cycle GANs)

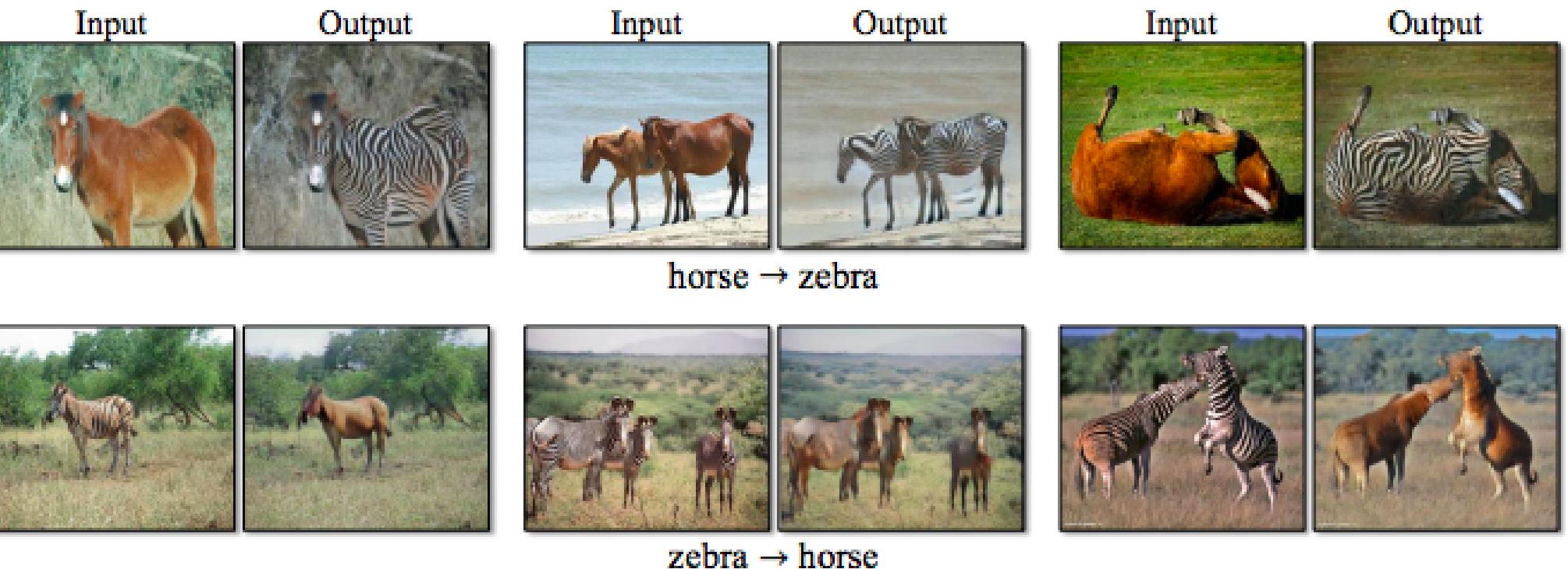
Zhu et al., March 2017

We have two corpora of images, say images of zebras and unrelated images of horses, or photographs and unrelated paintings by Monet.

We want to construct translations between the two classes.



Cycle Gans



Cycle Gans



Horse → Zebra

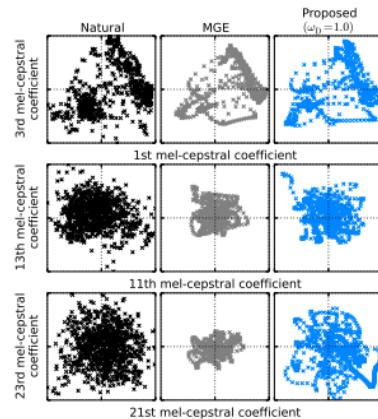
Unsupervised Machine Translation (UMT)

Lample et al, Oct. 2017, also Artetxe et al., Oct. 2017

In unsupervised machine translation the cycle loss is called **back-translation**.

Feature Alignment by Discrimination

Text to Speech (Saito et al. Sept. 2017)

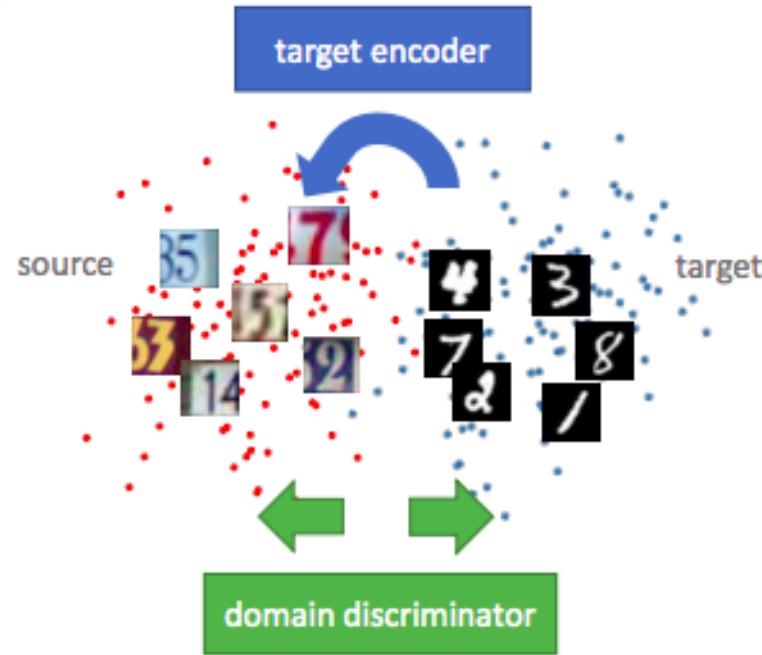


Minimum Generation Error (MGE) uses **perceptual distortion** — a distance between the feature vector of the generated sound wave and the feature vector of the original.

Perceptual Naturalness can be enforced by a feature discrimination loss.

Adversarial Discriminative Domain Adaptation

Tzeng et al. Feb. 2017



A feature discrimination loss can be used to align source and target features.

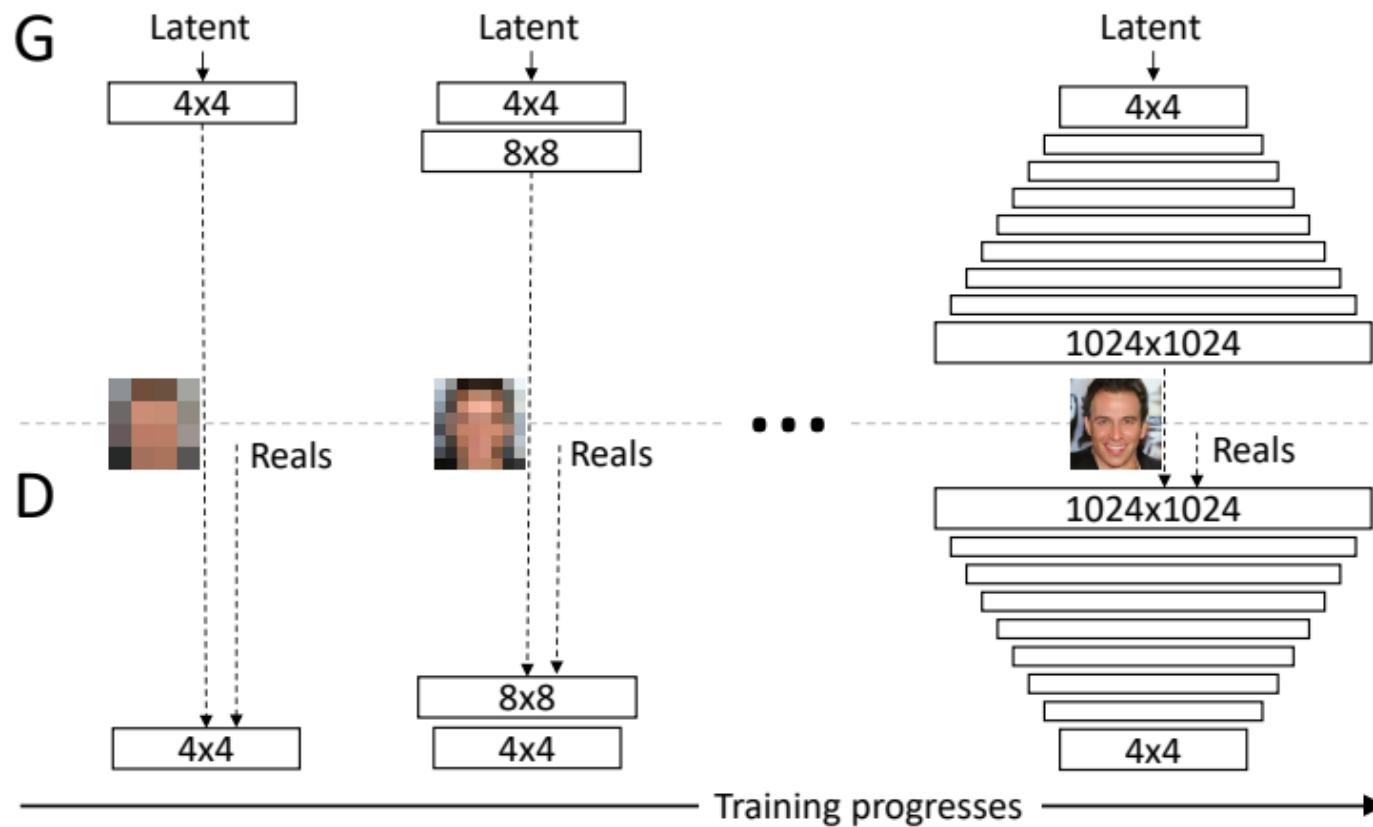
Progressive GANs

Progressive Growing of GANs, Karras et al., Oct. 2017

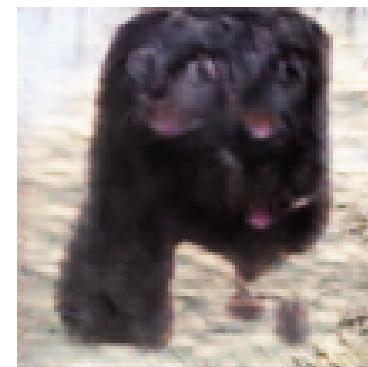


Figure 5: 1024×1024 images generated using the CELEBA-HQ dataset. See Appendix F for a larger set of results, and the accompanying video for latent space interpolations.

Progressive GANs



Early GANs on ImageNet



BigGans

Large Scale GAN Training, Brock et al., Sept. 2018



Figure 1: Class-conditional samples generated by our model.

This is a class-conditional GAN — it is conditioned on the imangenet class label.

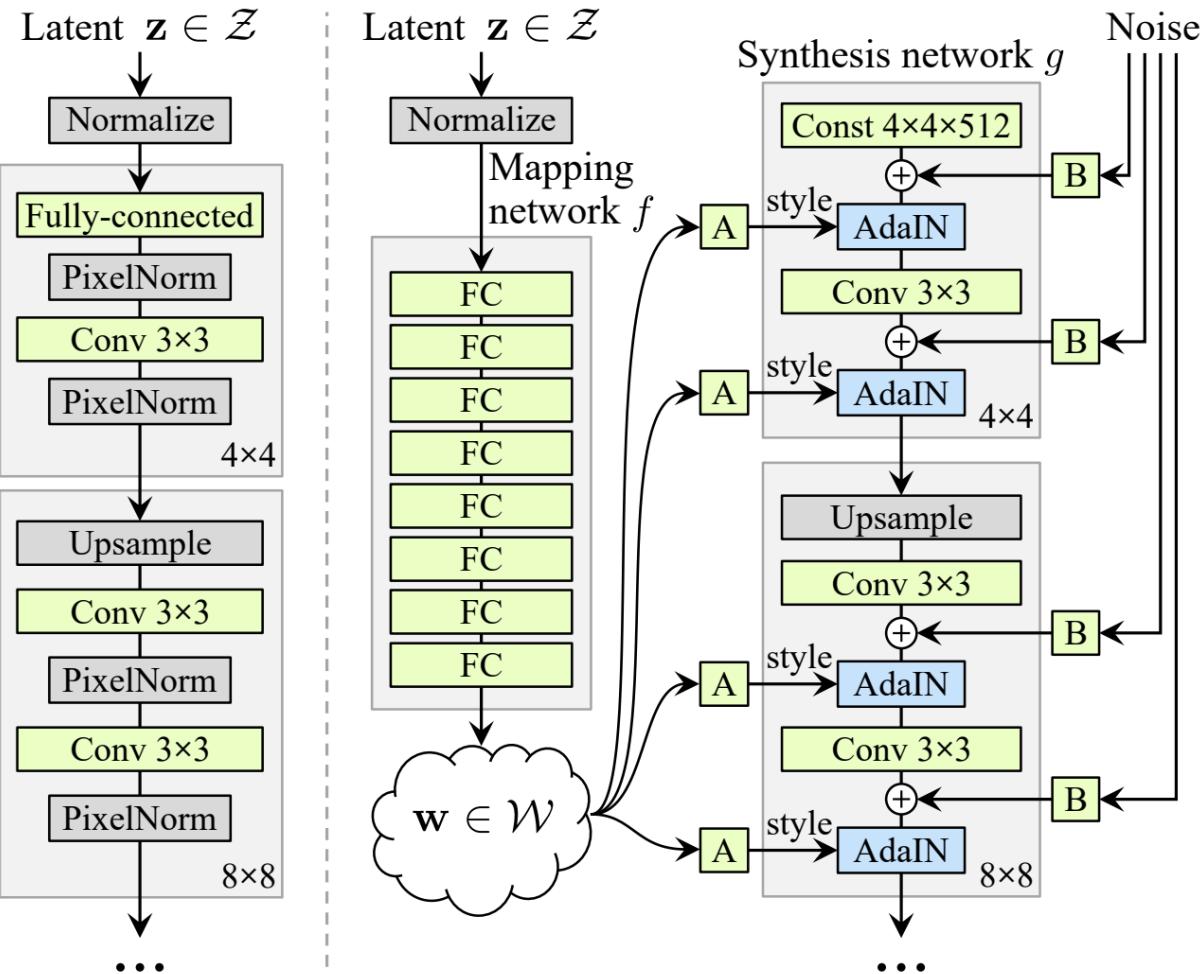
This generates 512 X 512 images without using progressive training.

StyleGAN

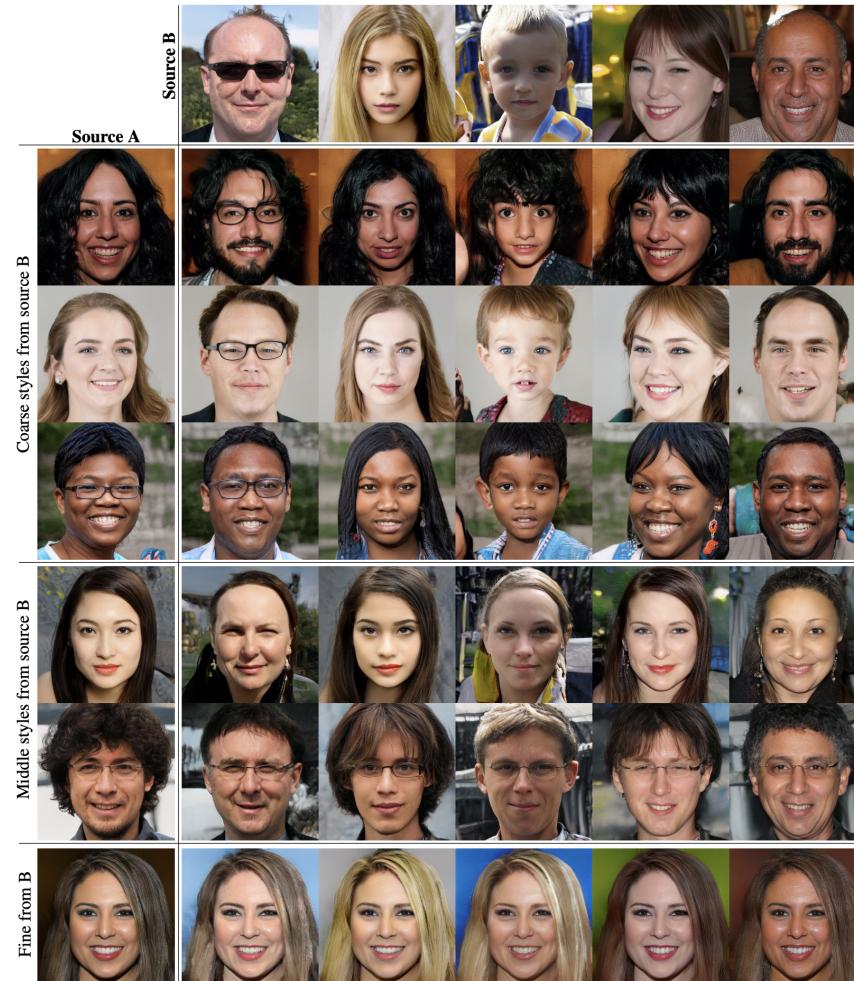
A Style-Based Generator Architecture for Generative Adversarial Networks, Karras et al., Dec. 2018



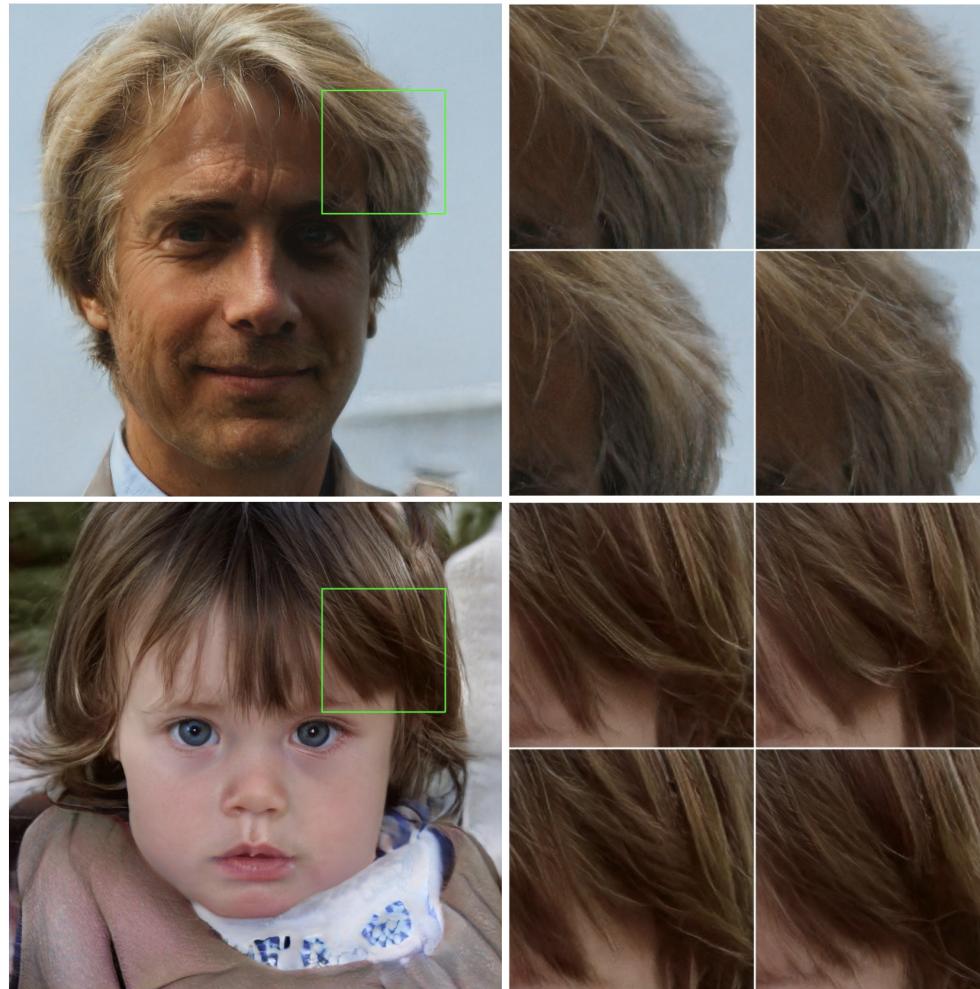
StyleGAN: Architecture



StyleGAN: Style Transfer



StyleGAN: Noise Variation



StyleGAN2 and StyleGAN3

StyleGan2 appeared in December of 2019 with significant improvements.

It was demonstrated to work on many classes of images, not just faces.

StyleGAN3 appeared in June 2021.

Projecting Images into Latent Space

Given an image, can we find a noise vector (a latent vector) that generate a close approximation of the given image. Can we invert the generator?

This appears to be possible with StyleGAN2 but not with the original even though StyleGAN2 produces better images.

By measuring the match between an image y and $g(g^{-1}(y))$ we can determine whether y was generated by StyleGAN2.

June 2023: StyleGAN Seems to “Understand” Images

StyleGAN knows Normal, Depth, Albedo, and More

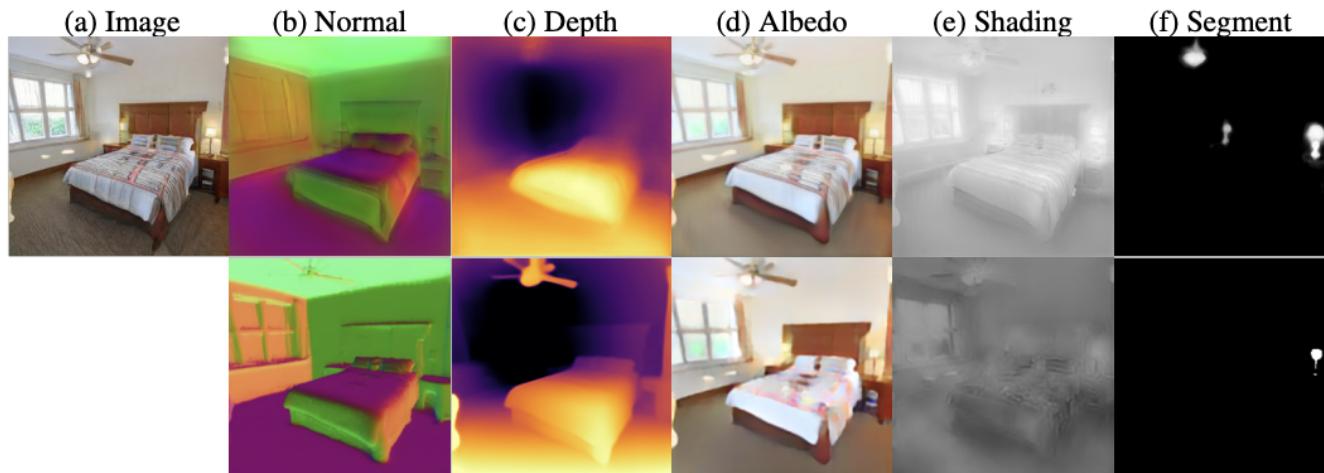
Anand Bhattacharjee

Daniel McKee

Derek Hoiem

D.A. Forsyth

University of Illinois Urbana Champaign



END