

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
- 3. GAN architectures**



Drawing? => learning from examples



Recall Algo GAN

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

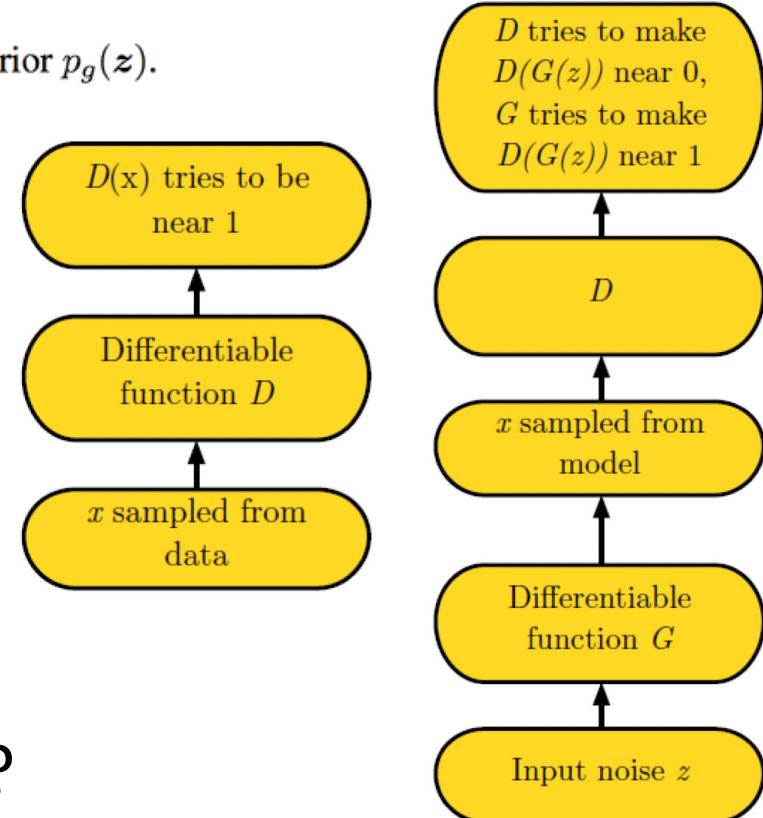
$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

end for



Functions G and D are NN

Question:

Which architectures for G and D?

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
- 3. GAN architectures**
 - 1. Basics**

Basic Archi for G and D and experiments

Models

G and D fully connected nets

or convolutional for D , (de)convolutional for G (as seen for segmentation nets)

ReLU and/or sigmoids, dropout

Datasets

MNIST, Toronto Face Database, CIFAR-10

GAN - Evaluation

- Approximate p_g by fitting a Gaussian Parzen window on the generated images.
- Cross-validate σ to maximize likelihood of validation set
- Compute the likelihood of the test set

Evaluation not trivial, can be done using generated images as inputs for deep nets => inception scores

Frechet inception distance using 2 Gaussian

(data,gen) over inception features: $FID = \|\mu - \mu_w\|_2^2 + \text{tr}(\Sigma + \Sigma_w - 2(\Sigma\Sigma_w)^{1/2})$

GAN - Qualitative results 1/2

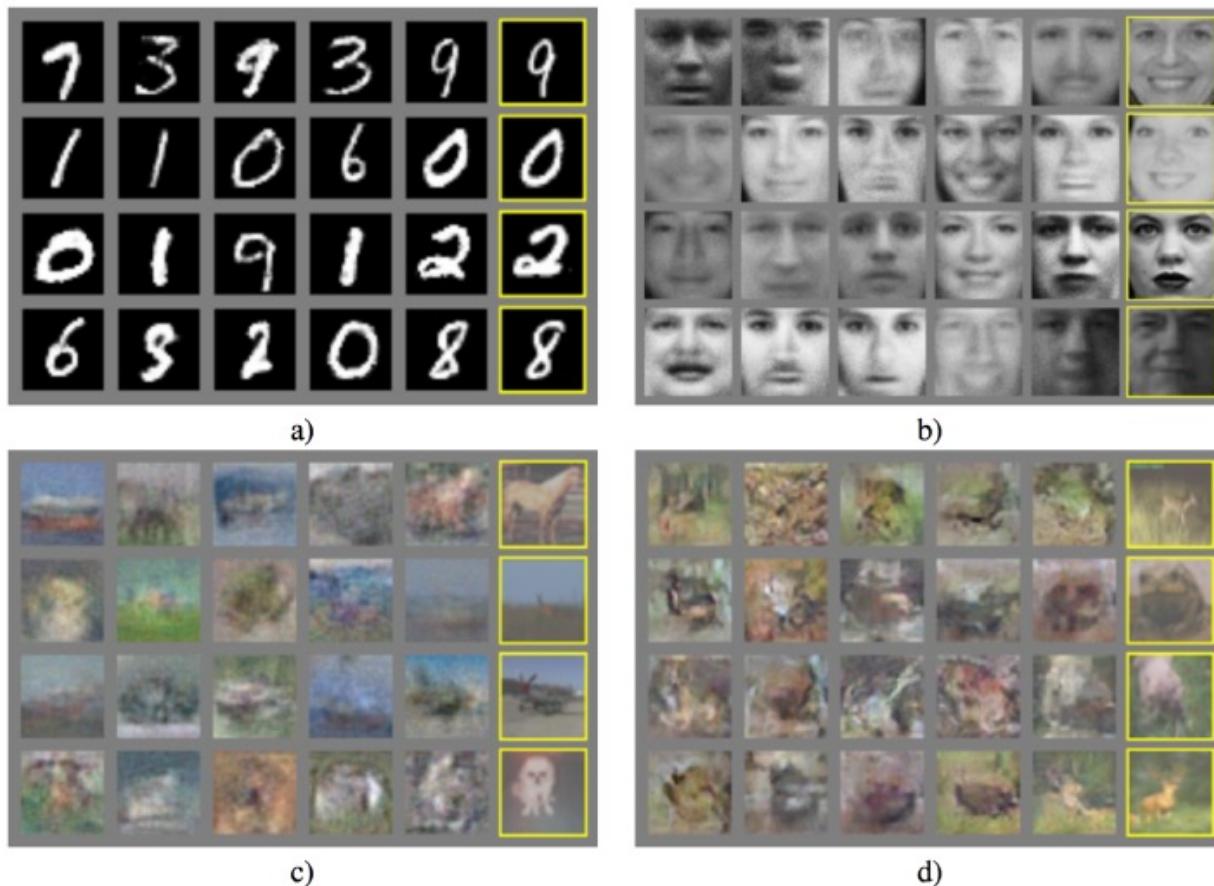


Figure: Right col nearest from dataset. a) MNIST, b) TFD, c) CIFAR-10 (fully connected), d) CIFAR-10 (convolutional D , deconvolutional G)

GAN - Qualitative results 2/2



Figure: Linear interpolation between 2 points in z space

- **Advantages:**
 - ▶ Computational advantages (no complex likelihood inference)
 - ▶ Can represent sharper distributions
- **Disadvantages:**
 - ▶ G and D must be well synchronized for the algorithm to converge correctly

GAN architectures

- How to improve result quality?
 - Spatial resolution
⇒ Cascade of GAN
 - Object quality
=> Progressive growing of spatial resolution in G

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. **LaPGAN**

Laplacian Pyramid GANs (LAPGANs)

- GANs do not work well for complex / high level / natural images.
- Idea: decompose the generation in successive tasks using Laplacian Pyramid (of GANs)

Let $d(I)$ and $u(I)$ be down-sampling and up-sampling operations.

Gaussian pyramid:

$$\mathcal{G}(I) = [I_0, I_1, \dots, I_K], I_k = d^{(k)}(I)$$

Laplacian pyramid:

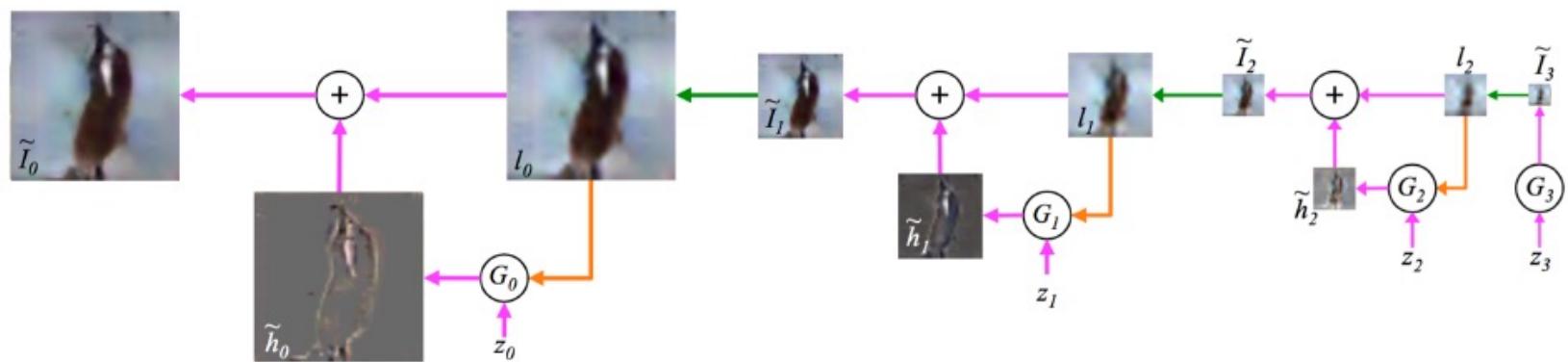
$$h_k = \mathcal{L}_k(I) = I_k - u(I_{k+1})$$

Reconstruction:

$$I_k = u(I_{k+1}) + h_k$$

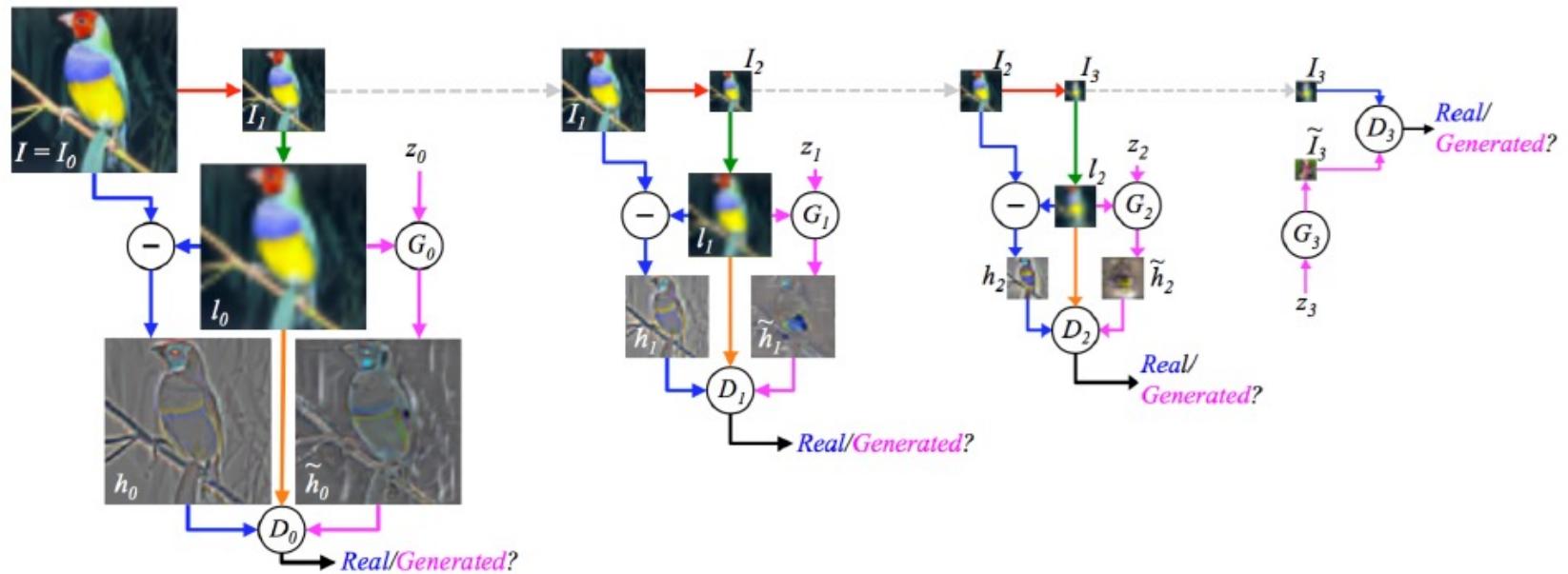
LAPGAN model - sampling

- Set of generative convnets: G_0, \dots, G_K
- Generated details: $\tilde{h}_k = G_k(z_k, u(\tilde{l}_{k+1}))$
- Reconstructed image: $\tilde{l}_k = u(\tilde{l}_{k+1}) + \tilde{h}_k$ ($\tilde{l}_{K+1} = 0$)



LAPGAN model - training

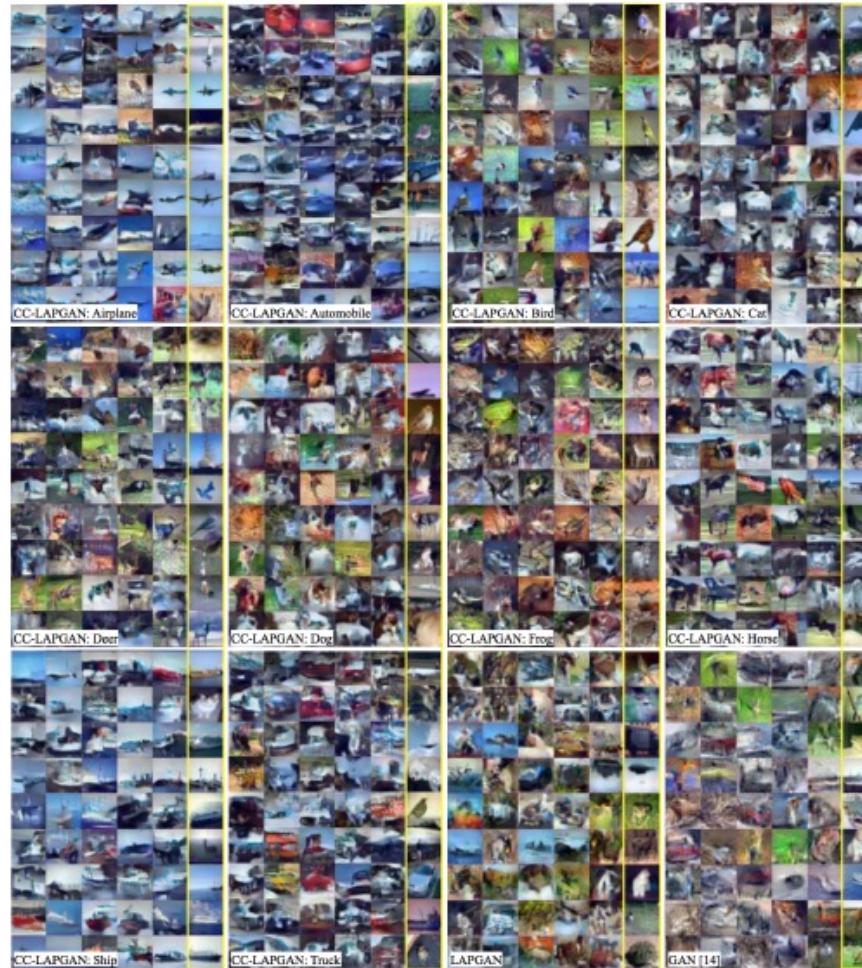
- Low-pass version of I_0 : $I_0 = u(d(I_0))$
- Either:
 - ▶ High-pass version of I_0 : $h_0 = I_0 - I_0$
 - ▶ Generate $\tilde{h}_0 = G_0(z_0, I_0)$
- Forward $D_0(I_0 + h_0$ or $\tilde{h}_0)$
- Backward D_0 and G_0
- G_K and D_K are trained as a simple GAN



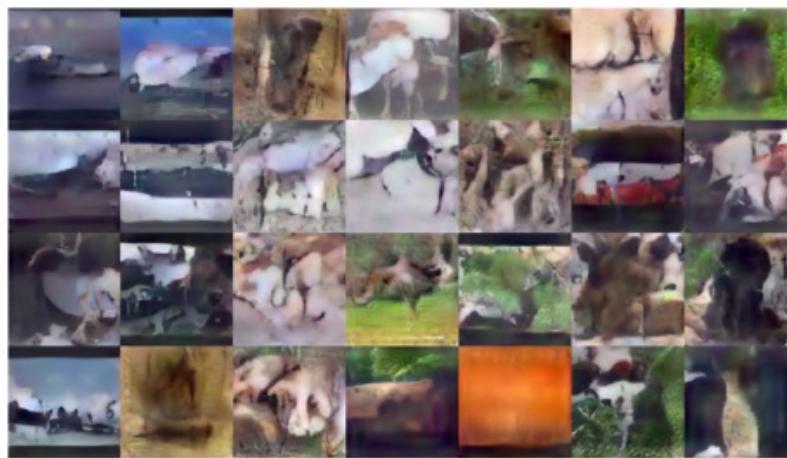
LAPGAN model - Experiments

- **Datasets:** CIFAR-10, STL
- **Initial scale:**
 - ▶ G_K and D_K have 2 hidden layers and ReLU
 - ▶ $z_K \sim U_{[-1,1]^{100}}$
 - ▶ Trained as GAN
- **Subsequent scales:**
 - ▶ G_k and D_k convnets with 3 and 2 layers
 - ▶ $z_k \sim U_{[-1,1]^{|I_k|}}$ (“color” layer)
 - ▶ Trained as CGAN

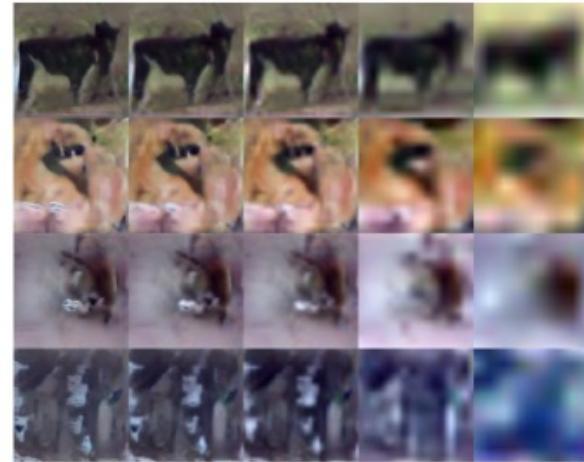
LAPGAN model - Results - CIFAR



LAPGAN model - Results - STL



(a)



(b)

Figure 4: STL samples: **(a)** Random 96x96 samples from our LAPGAN model. **(b)** Coarse-to-fine generation chain.

LAPGAN model - Results - LSUN



LAPGAN

- Good idea (**cascade**) but Generator structure too weak
- => Other approach: **progressive growing** of spatial resolution

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. **DCGAN**

Progressive growing of spatial resolution in G

Deep Convolutional GANs (DCGANs)

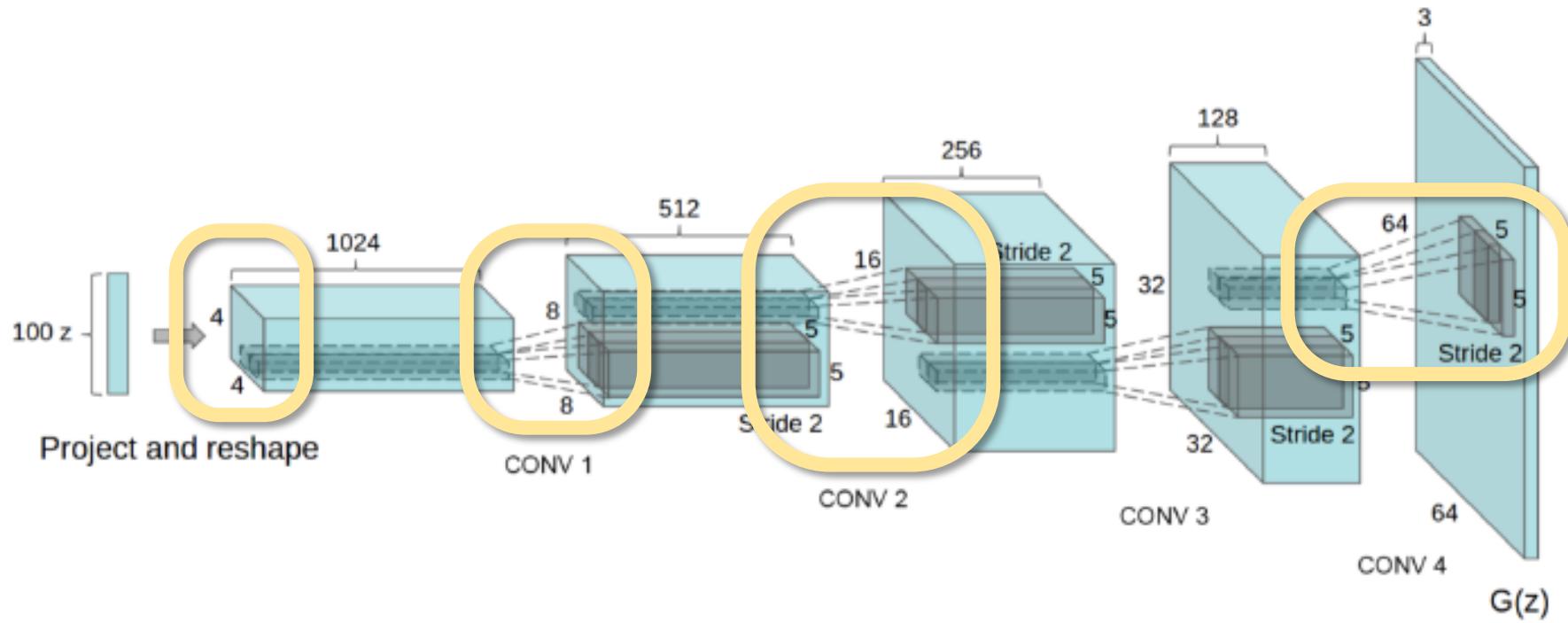
GANs are hard to scale => Identify a family of architectures that gives stable training

- Replace any pooling layers with strided convolutions
=transposed convolutions =deconvolutions
(discriminator) and fractional-strided convolutions (generator)
- Use batchnorm in both the generator and the discriminator
- Remove fully connected hidden layers for deeper architectures
- Use ReLU activation in generator for all layers except for the output, which uses Tanh
- Use LeakyReLU activation in the discriminator for all layers

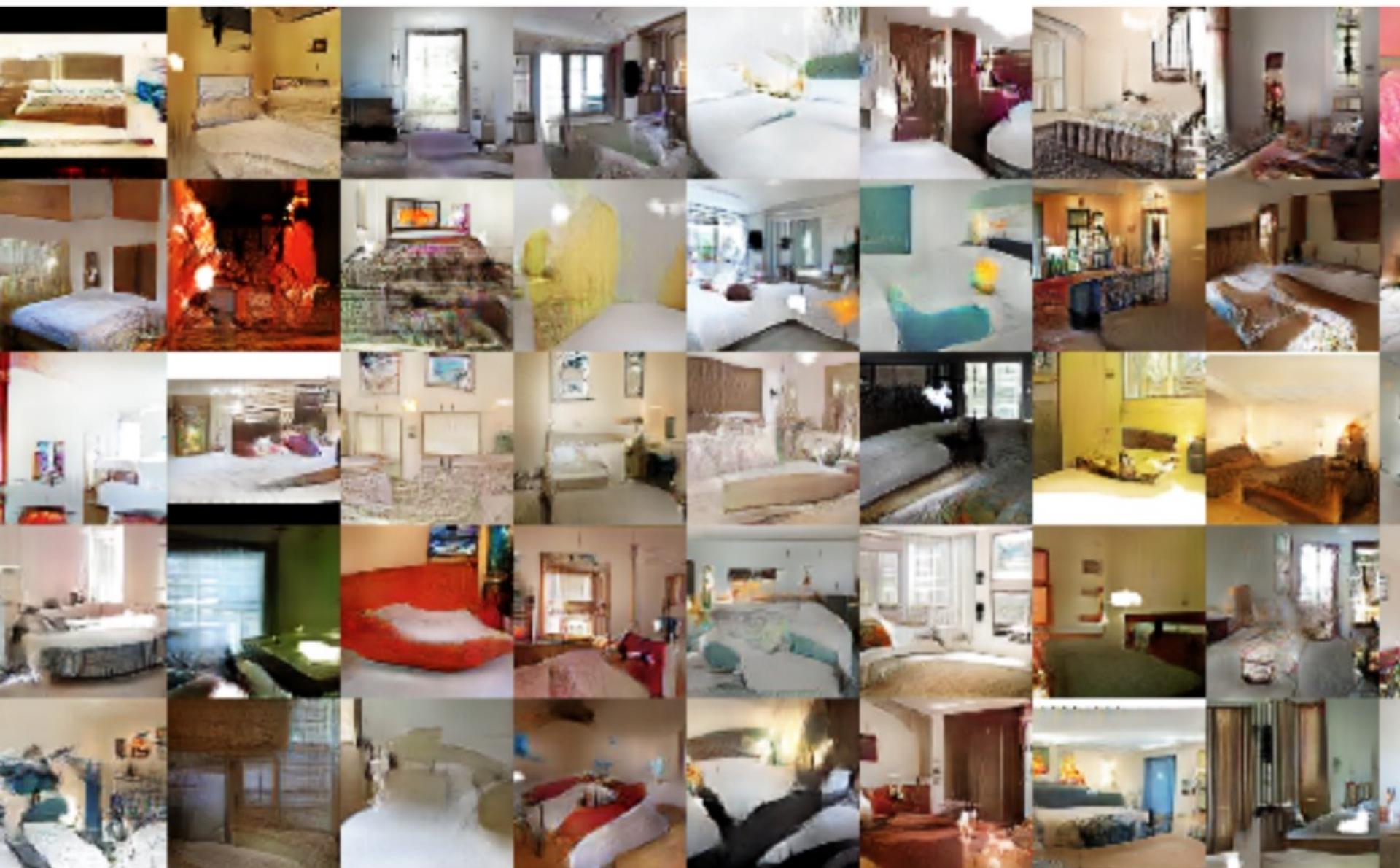
Progressive growing of spatial resolution in G: DCGAN

Upsampling step by step

Combine with convolutional layers



DCGAN - Results - generated bedrooms



DCGAN results - Faces



Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. DCGAN
 4. ProGAN

Progressive growing of GANs

[Progressive Growing of GANs for Improved Quality, Stability and Variation, Tero Karras et al. (NVIDIA); ICLR 2018]

Combine idea of LAPGAN (several output reso) and DCGAN (archi prog growing)

1. First, start with training 4x4 output images.
2. When this training has converged, add a new block to generate 8x8 output images.
3. Etc.

The transition to adding a new block is gradual, we first start with more weight on the (upsampled) output of the previous block, and then add more and more weights to the output of the current block.

All weights remain trainable during the whole process.

Discriminator = mirror image of generator

Progressive growing of GANs

[Progressive Growing of GANs for Improved Quality, Stability and Variation, Tero Karras et al. (NVIDIA); ICLR 2018]

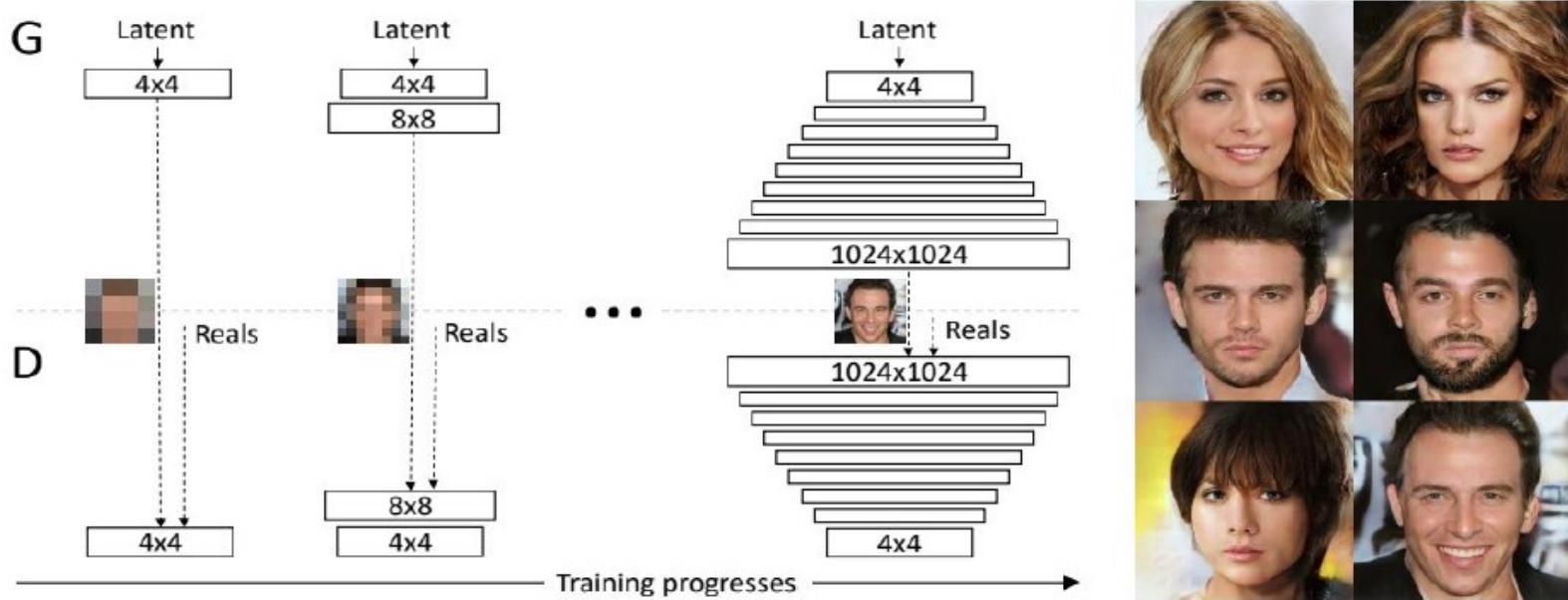


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. **All existing layers remain trainable throughout the process.** Here $N \times N$ refers to convolutional layers operating on $N \times N$ spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. One the right we show six example images generated using progressive growing at 1024×1024 .

Progressive growing of GANs

[Progressive Growing of GANs for Improved Quality, Stability and Variation, Tero Karras et al. (NVIDIA); ICLR 2018]

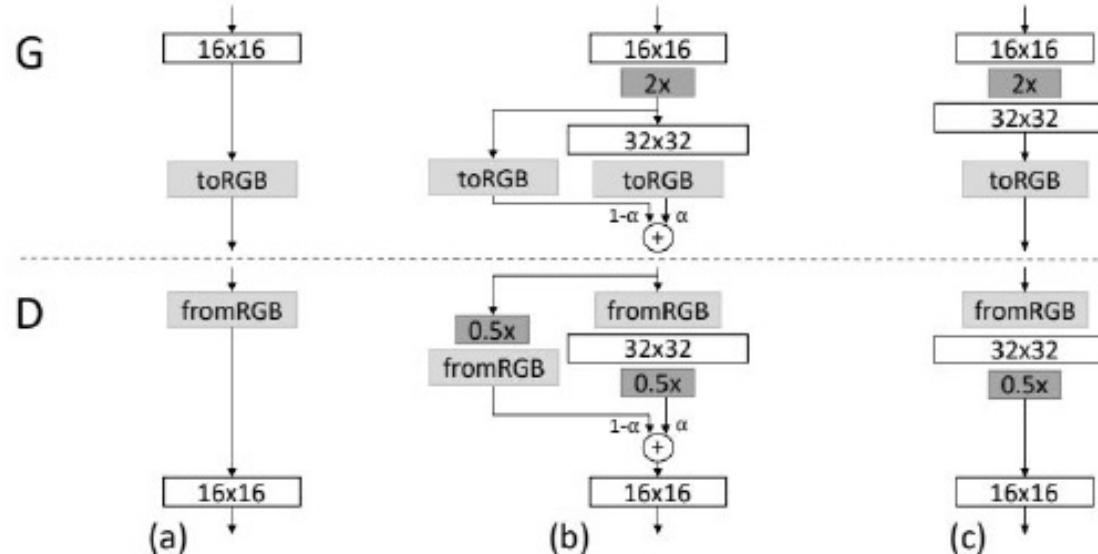


Figure 2: When doubling the resolution of the generator (G) and discriminator (D) we fade in the new layers smoothly. This example illustrates the transition from 16×16 images (a) to 32×32 images (c). During the transition (b) we treat the layers that operate on the higher resolution like a residual block, whose weight α increases linearly from 0 to 1. Here $2\times$ and $0.5\times$ refer to doubling and halving the image resolution using nearest neighbor filtering and average pooling, respectively. The toRGB represents a layer that projects feature vectors to RGB colors and fromRGB does the reverse; both use 1×1 convolutions. When training the discriminator, we feed in real images that are down-scaled to match the current resolution of the network. During a resolution transition, we interpolate between two resolutions of the real images, similarly to how the generator output combines two resolutions.

Generative models

Outline

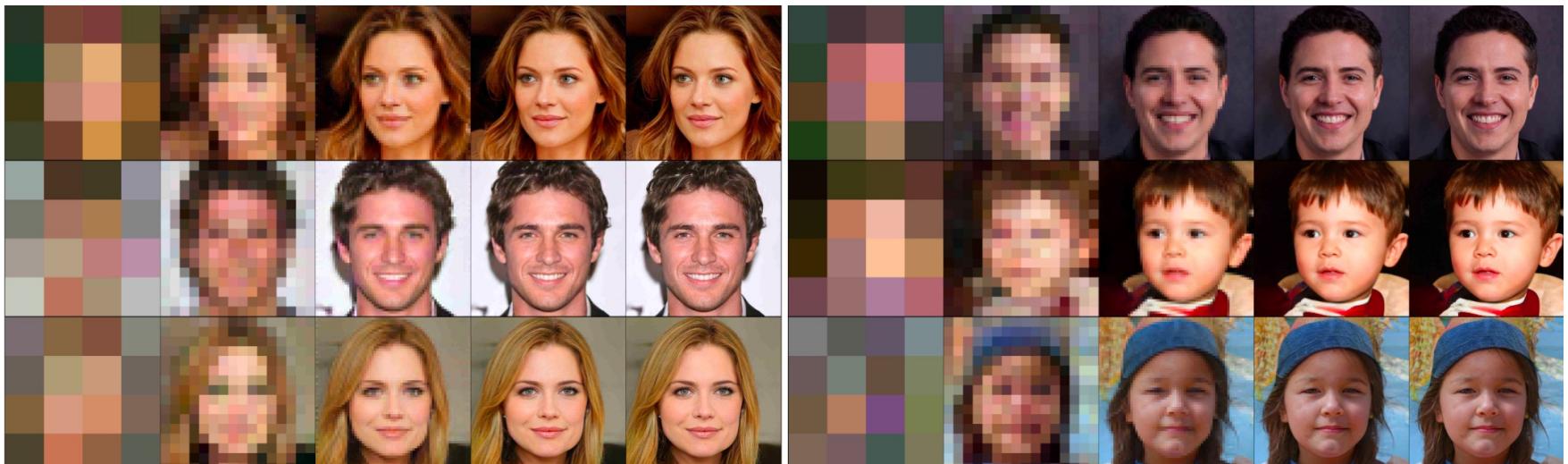
1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. DCGAN
 4. ProGAN
 5. **MSG-GAN**

MSG-GAN: Multi-Scale Gradients for Generative Adversarial Networks [CVPR 2020]

Main Idea:

- ProGAN both use progressive growing, but although this gives stability, it introduces many complicated training parameters associated with each new network.
- Training cannot be done “out of the box”, have to adjust parameters for each new dataset.

→ Train all at once without complicated adding on layers



MSG-GAN: Multi-Scale Gradients for Generative Adversarial Networks [CVPR 2020]

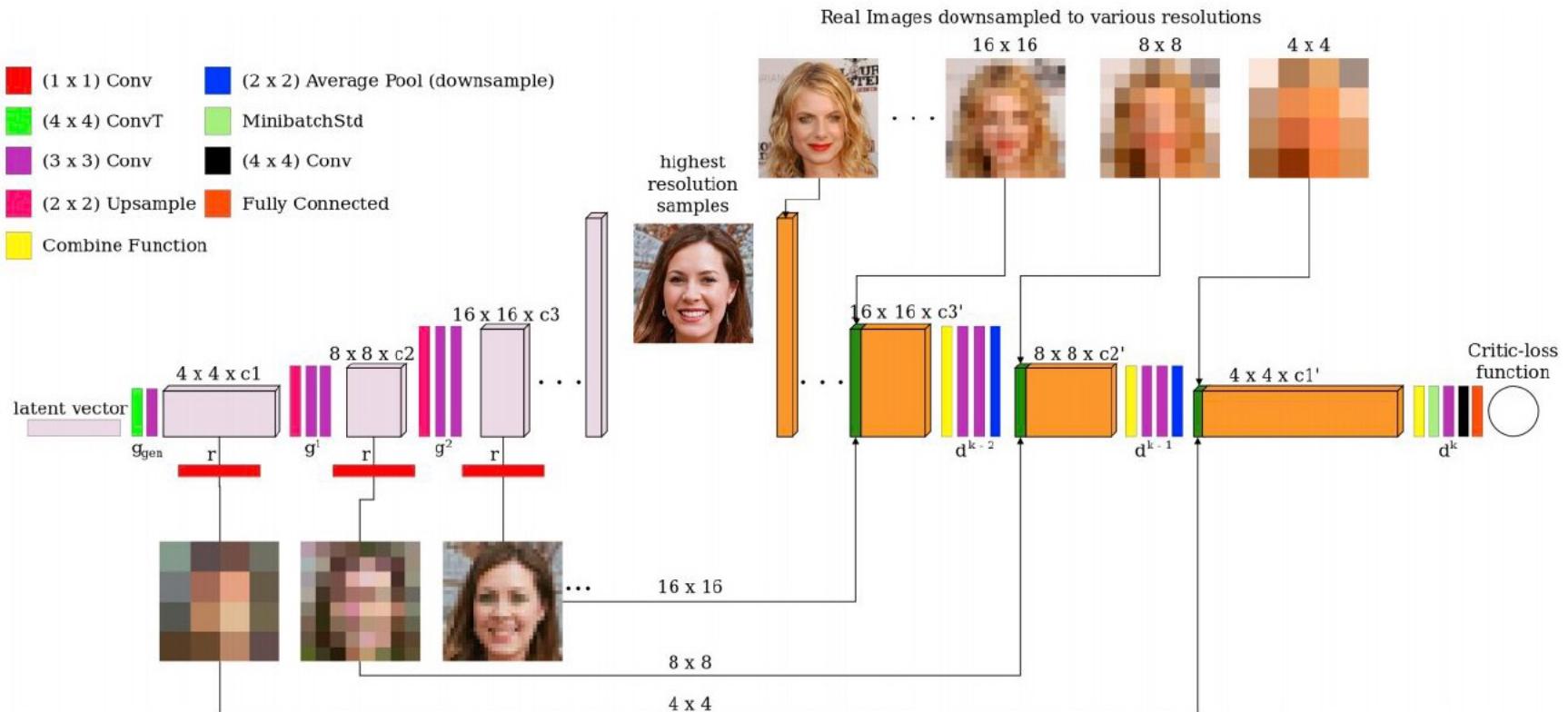


Figure 2: Architecture of MSG-GAN, shown here on the base model proposed in ProGANs [13]. Our architecture includes connections from the intermediate layers of the generator to the intermediate layers of the discriminator. Multi-scale images sent to the discriminator are concatenated with the corresponding activation volumes obtained from the main path of convolutional layers followed by a combine function (shown in yellow).

MSG-GAN: results – Random generated CelebA-HQ Faces at resolution 1024x1024



Generative models

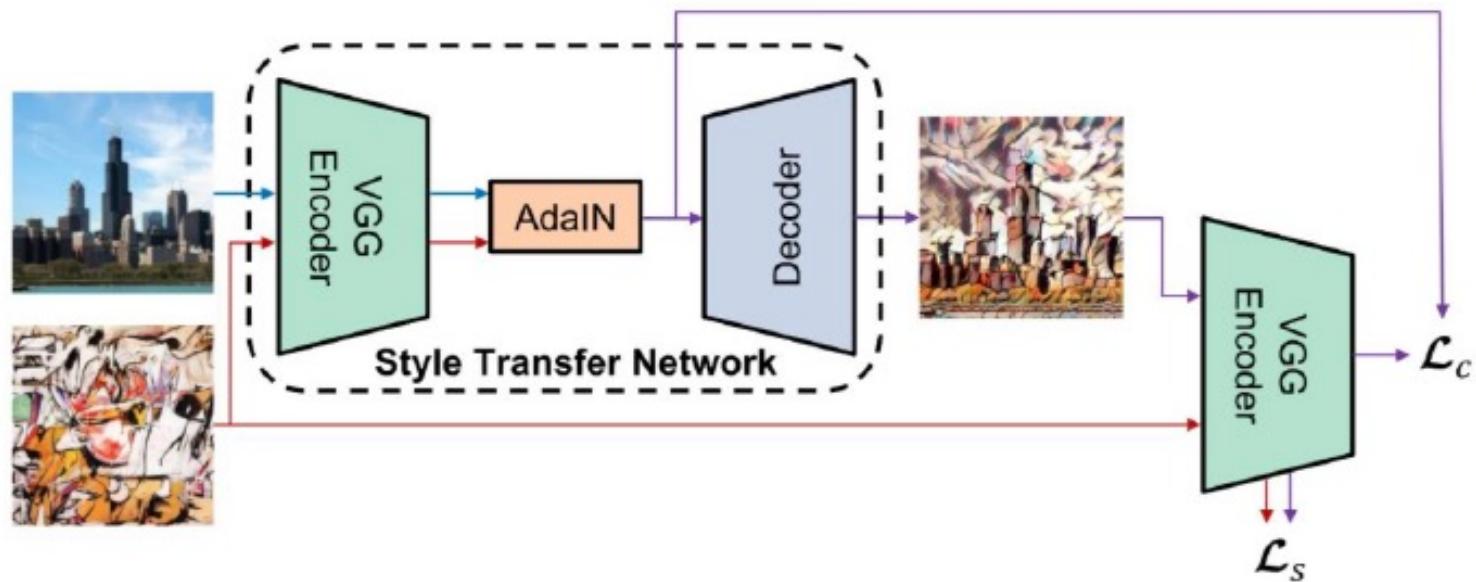
Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. DCGAN
 4. ProGAN
 5. MSG-GAN
 6. **StyleGAN**

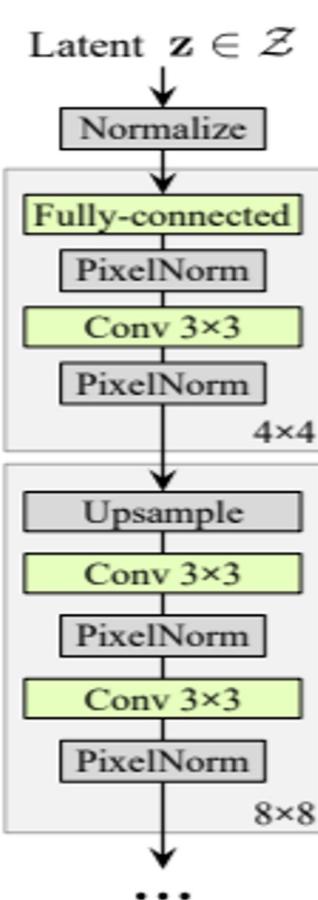
StyleGAN: A Style-Based Generator Architecture for Generative Adversarial Networks [Karras CVPR 2019]

Still progressive growing architecture but with new refinement block based on: Arbitrary Style Transfer in Real-time with Adaptive Instance Normalization (AdaIN)

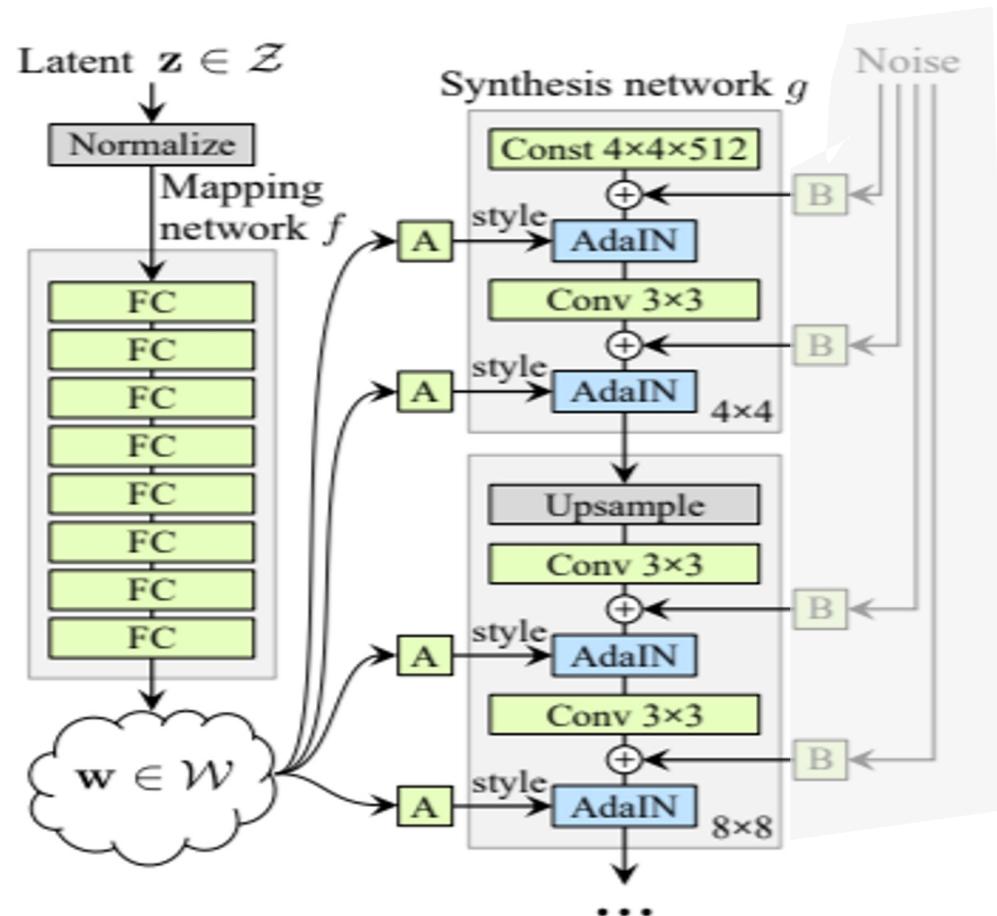
$$\text{AdaIN}(x, y) = \sigma(y) \left(\frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$



StyleGAN Network Architecture



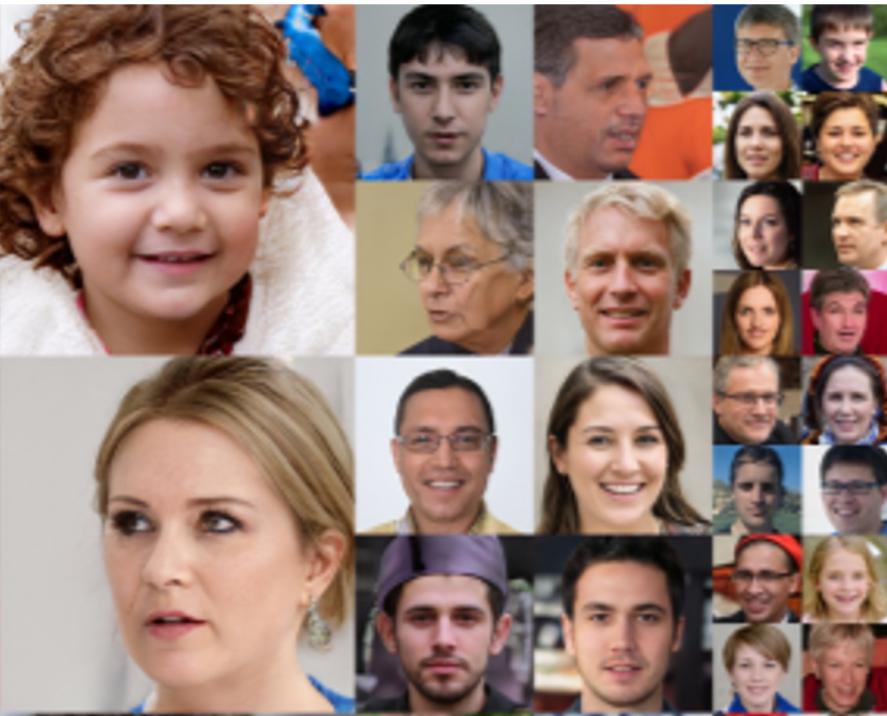
(a) Traditional



(b) Style-based generator

Building up the Model

Method	CelebA-HQ	FFHQ
A Baseline Progressive GAN [29]	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







Generative models

Outline

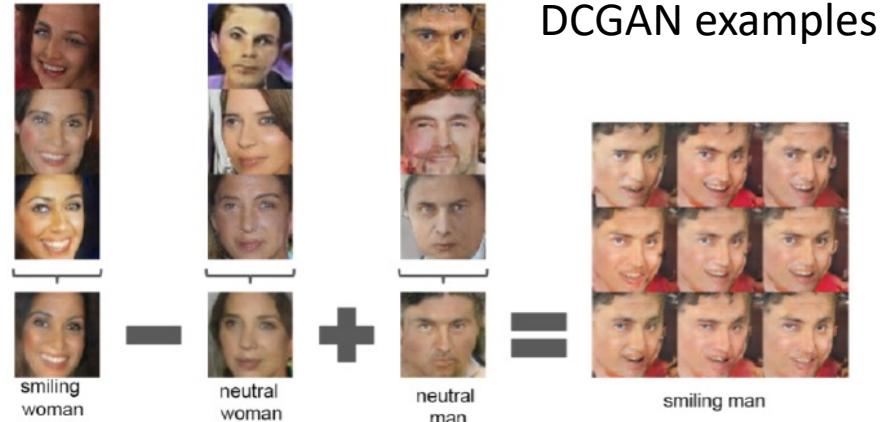
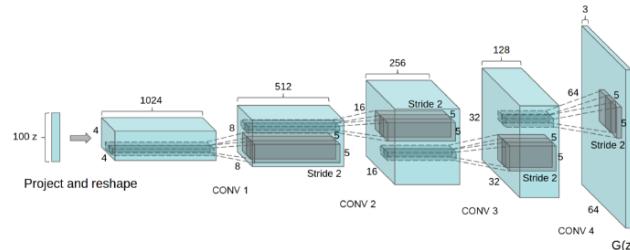
1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
- 4. Editing**

GAN editing

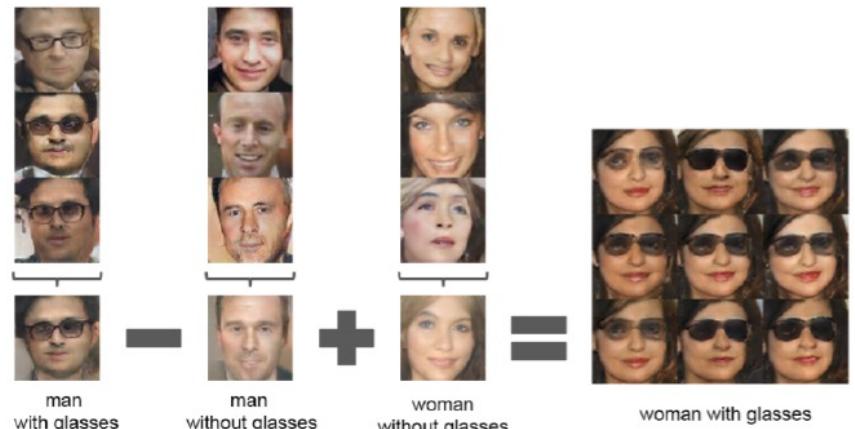
Editing by manipulation in the latent space of z

Image $I = G(z)$

Editing: changing z to z' and the new image is $I' = G(z')$

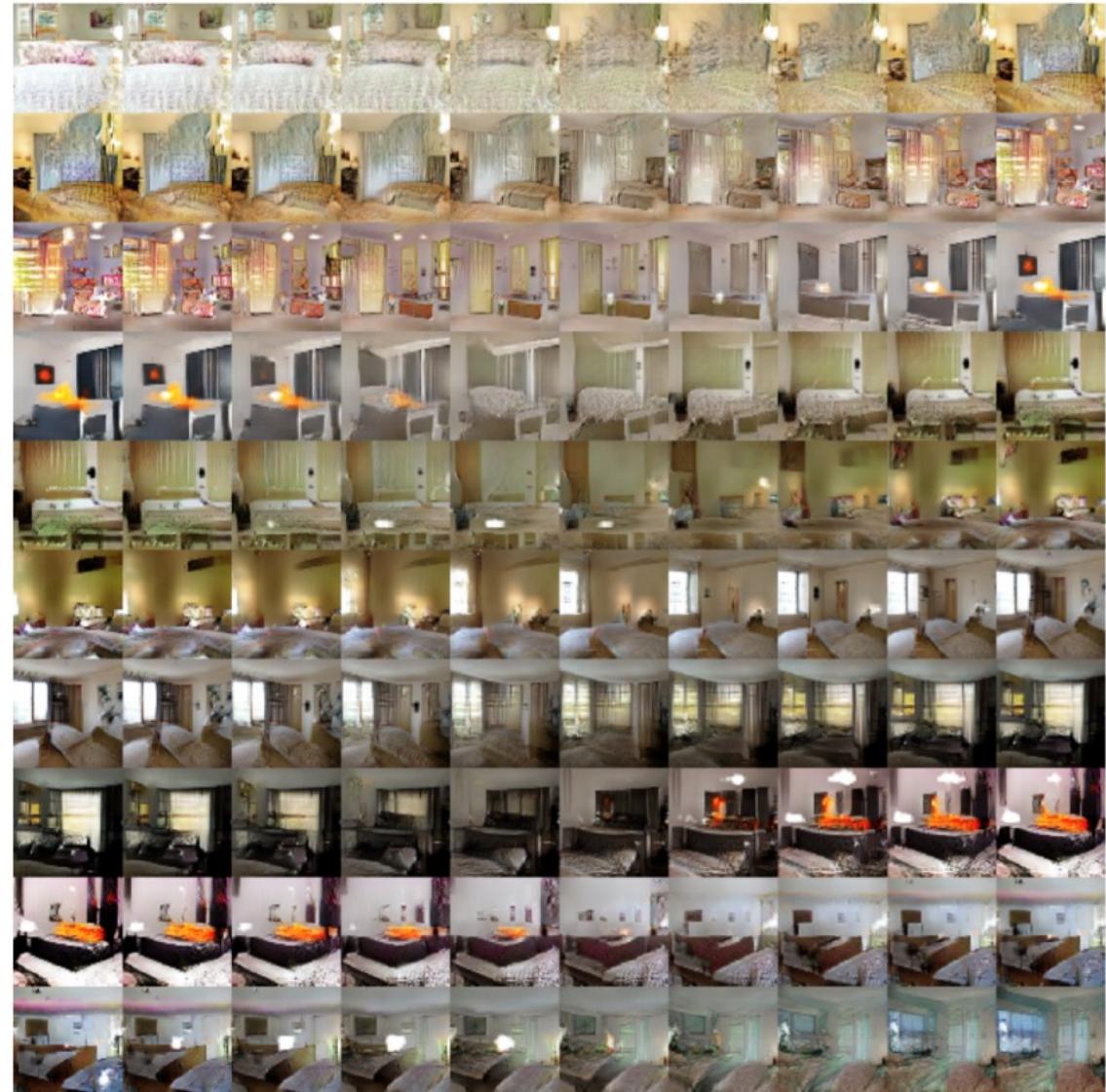


Arithmetics in latent space: vector mean, addition, subtraction



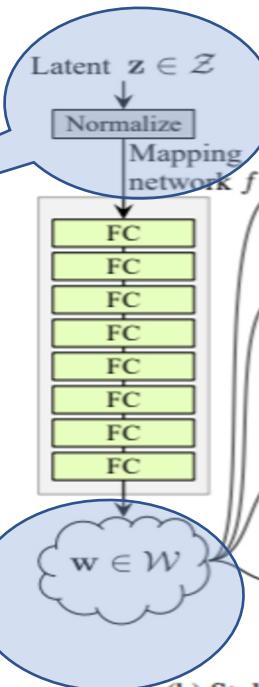
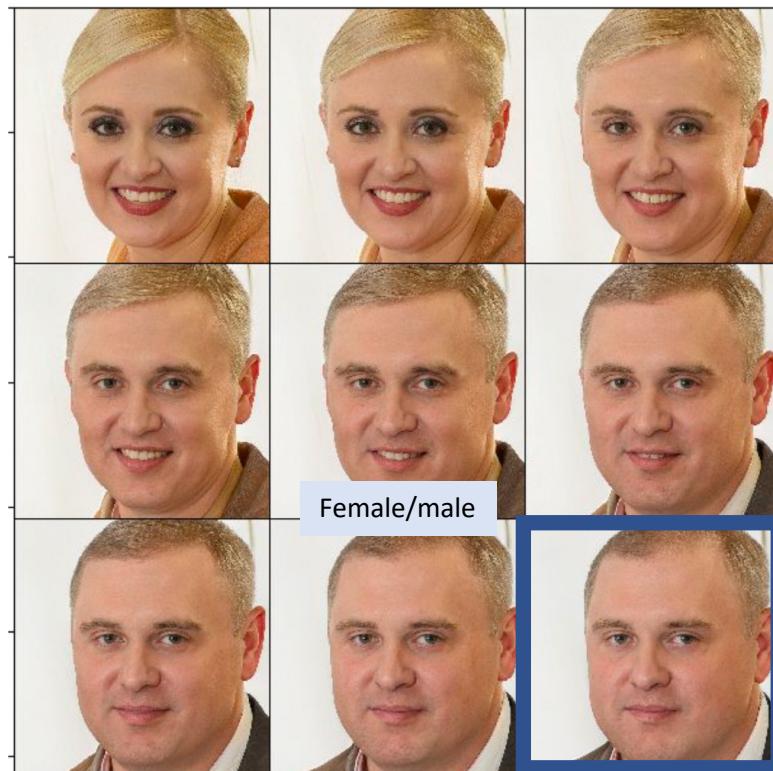
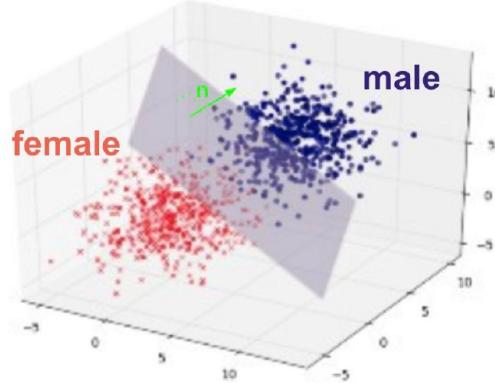
GAN editing

Linear interpolation
in latent space

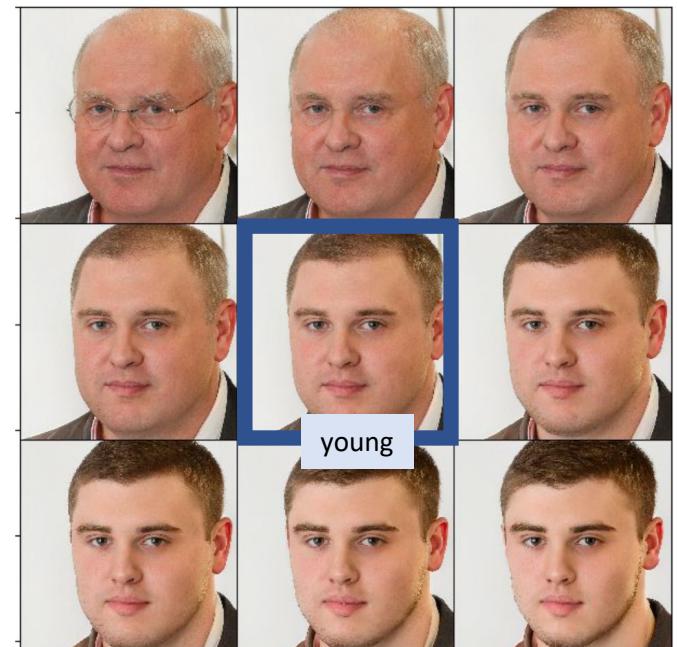


Gan Editing

Latent space analysis for GAN editing with StyleGAN

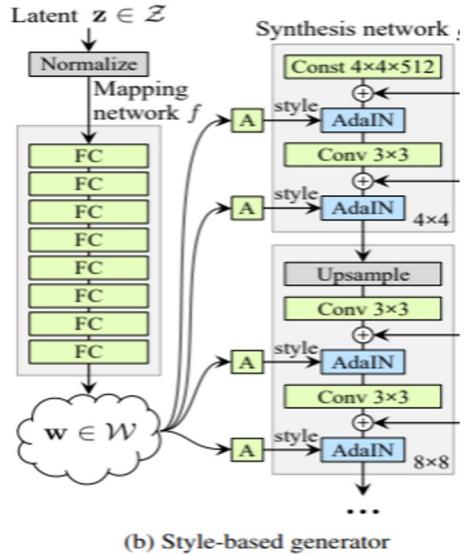


(b) Style-based generator



GAN editing

Latent space analysis for GAN editing with StyleGAN



Two sets of images were generated from their respective latent codes (sources A and B); the rest of the images were generated by copying a specified subset of styles from source B and taking the rest from source A. Copying the styles corresponding to coarse spatial resolutions ($4 \times 4 - 8 \times 8$) brings high-level aspects such as pose, general hair style, face shape, and eyeglasses from source B, while all colors (eyes, hair, lighting) and finer facial features resemble A. If we instead copy the styles of middle resolutions ($16 \times 16 - 32 \times 32$) from B, we inherit smaller scale facial features, hair style, eyes open/closed from B, while the pose, general face shape, and eyeglasses from A are preserved.

