

Generative Adversarial Networks

David Coufal

Institute of Computer Science
The Czech Academy of Sciences
david.coufal@cs.cas.cz

Vision for Robotics - FEL CTU

November 29, 2021

Neural networks

- a **neural network** is a complex composite function built from individual **layers of neurons**
neurons represent **simple computation units**
- **neurons are parametrized**, so the whole network is a **highly parametrized function**
- adjustment of parameters is called **network learning** via **back propagation** of some **loss function** error
- **shallow networks** - one hidden layer of neurons
- **deep networks** - multiple layers
(up to 200 layers, millions of parameters)

Standard neural networks

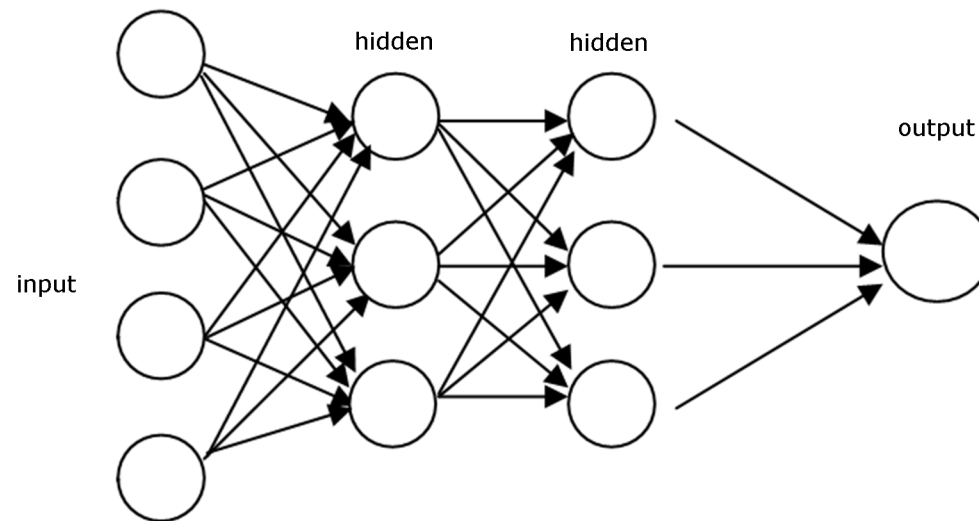
- standard **neuron** $h : \mathbb{R}^d \rightarrow \mathbb{R}$ has form

$$h(\mathbf{x}) = \text{act}(\mathbf{w}\mathbf{x} + \mathbf{b})$$

- $\text{act}(z) = \frac{1}{1+e^{-\beta z}}$ (sigmoid)

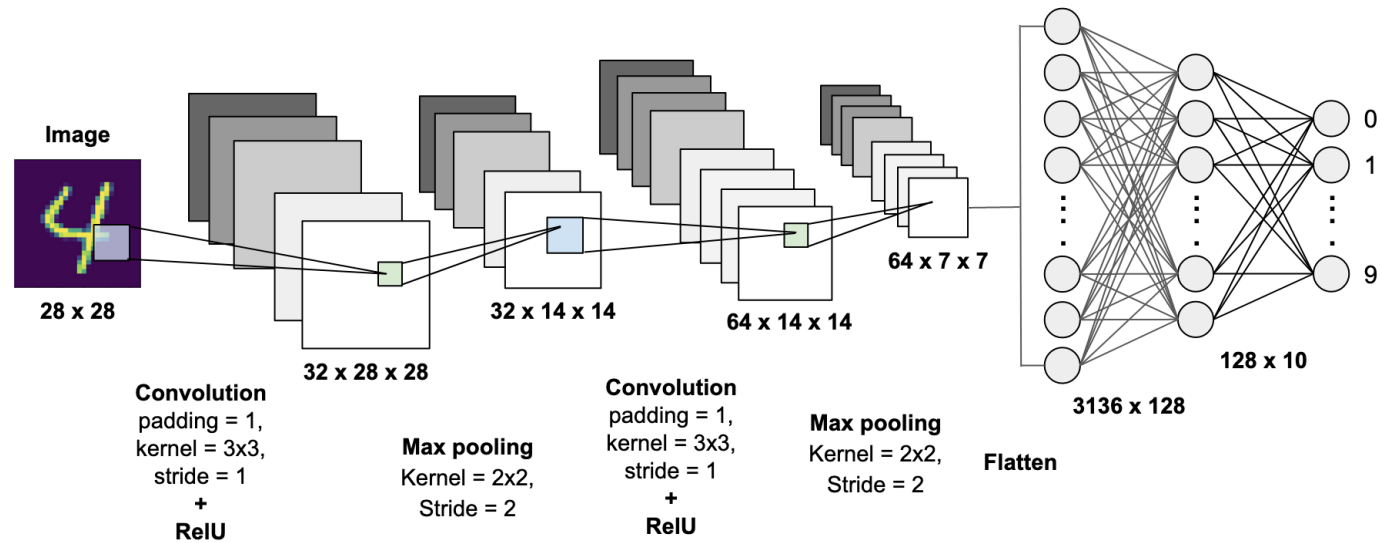
- $\text{act}(z) = \max(0, z)$ (ReLU)

- $\mathbf{w}, \mathbf{b} \in \mathbb{R}^d$ - **parameters**



Convolutional neural networks

- **convolution filters** moving over the input

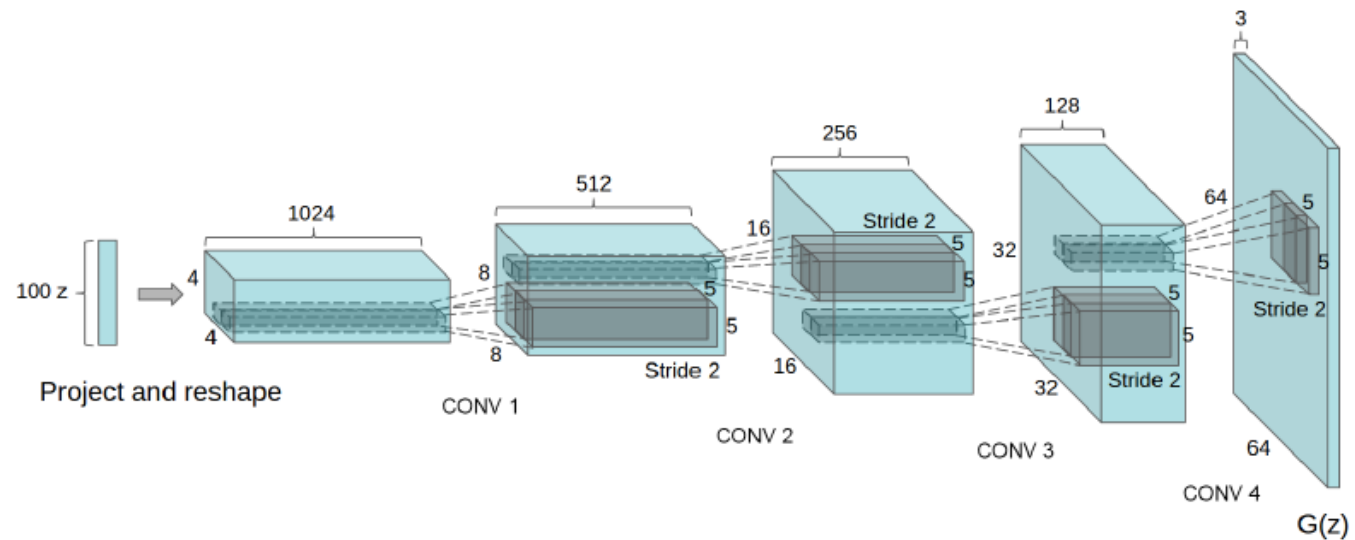


source: <https://towardsdatascience.com/mnist-handwritten-digits-classification-using-a-convolutional-neural-network-cnn-af5fafbc35e9>

- **down-sampling** and **up-sampling** operations, pooling

Convolutional neural networks - upsampling

- **transposed convolutions** - increase in spatial dimensions



- standard convolutions with manipulated inputs

Well recognized DL tasks

- **classification**
ImageNet Large Scale Visual Recognition Challenge AlexNet
CNN network won the contest in 2012
- **recurrent neural networks** (RNNs)
LSTM, GRU - units, NLP tasks, Google Translator
(but CNNs are used in DeepL.com)
- **reinforcement learning** DeepMind (UK, Google 2014)
AlhaGo vs. Lee Sedol (4:1, 2016), AlphaGo Zero vs. AlphaGo
(100:0, 2017) AlphaZero vs. Stockfish (28:72:0, 2018),
Dota 2 tournaments, [AlphaFold \(2021\)](#)
- **generative programming**
Ian Goodfellow et al. (2014) - [Generative Adversial Networks](#)
<https://arxiv.org/abs/1406.2661>

Elementary concepts

- **random variable** $X \sim P_X$, $(\Omega, \mathcal{A}, P_X)$
 - Ω - space of elementary events $X \in \Omega$
 - \mathcal{A} - sigma algebra of measurable events
 - P_X - distribution of X
- **distribution of X**
 - set function on \mathcal{A} , $P_X : \mathcal{A} \rightarrow [0, 1]$
 - obeys Kolmogorov's laws of probability
 - typically $\Omega \in \mathbb{R}^d$ and $\mathcal{A} = \mathcal{B}(\mathbb{R}^d)$
- **data** $D = \{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$ **comes from distribution P_D**
i.e., we assume that there exists a random variable D
such that $D \sim P_D$ (sometimes we use P_{data} instead of P_D)
- **How to specify P_D on the basis of D ?**

Elementary concepts

- if Ω is **countable**, P_D can be given **by enumeration**, i.e., $P_D(\omega_i) = p_i$, for $i = 1, \dots, n$ (finite) or $i \in \mathbb{N}$ (countable)
- if $\Omega = \mathbb{R}^d$, specification of cdf is possible, but inconvenient in higher dimensions, so the most common approach is **to specify a density** $p_D : \mathbb{R}^d \rightarrow [0, \infty)$ of P_D and one has

$$P_D(A) = \int_A p_D(\mathbf{x}) d\mathbf{x} \quad \text{for } A \in \mathcal{B}(\mathbb{R}^d)$$

- cannot handle distributions which do not have densities **complex formulas in high dimensions for dependent data**
- **How to get the density from empirical data?**

Elementary concepts

- if $p_D \in \{p_\theta, \theta \in \Theta\}$ (a parametric set of densities) task reduces to estimate θ^* from data D and $p_D = p_{\theta^*}$
maximum likelihood estimation
- in a non-parametric context, kernel density estimation is the standard choice

$$p_D^*(x) = \frac{1}{nh^d} \sum_{k=1}^n K\left(\frac{x - x_i}{h}\right)$$

- $K : \mathbb{R}^d \rightarrow \mathbb{R}$, a kernel (bump) function, $h > 0$ is the bandwidth
practically applicable for d up to 5
- How to sample from a given distribution/density?

Distance of probability distributions

- **space of probability distributions** on $\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d)$:
 $\mathcal{P} = \{P : \text{probability distribution on } (\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))\}$
 \mathcal{P} is metrizable, e.g., using Lévy-Prokhorov metric
 $\pi : \mathcal{P}^2 \rightarrow [0, \infty)$, complicated formulas
- another "metric" is the **Kullback-Leibler divergence**
let $P, Q \in \mathcal{P}$, $P \ll Q$ (if $Q(x) = 0$, then $P(x) = 0$)

$$\begin{aligned} KL(P||Q) &= \int \frac{dP}{dQ} dP \\ &= \int \log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

- **properties:**
 $KL(P||P) = 0$, $KL(P||Q) \geq 0$, $KL(P||Q) \neq KL(Q||P)$
- tight relation to **theory of information** (relative entropy),
theory of large deviations

Kullback-Leibler divergence

- (Wikipedia entry ...) In applications, P typically represents the "true" distribution of data, observations, or a precisely calculated theoretical distribution, while Q typically represents a theory, model, description, or approximation of P . In order to find a distribution Q that is closest to P , we can minimize KL divergence.
- Jensen-Shannon divergence - symmetrized KL divergence

$$JSD(P||Q) = \frac{1}{2}KL(P||M) + \frac{1}{2}KL(Q||M)$$

where $M = \frac{1}{2}(P + Q)$

$$JSD(P||P) = 0, 0 \leq JSD(P||Q) \leq 1, JSD(P||Q) = JSD(Q||P)$$

Reverse information projection (M-projection)

- let $P \in \mathcal{P}$ is fixed, and $\mathcal{Q} \subset \mathcal{P}$ (subset of prob. distributions)

$$Q_{KL}^* = \arg \min_{Q \in \mathcal{Q}} KL(P||Q)$$

or for JS

$$Q_{JSD}^* = \arg \min_{Q \in \mathcal{Q}} JSD(P||Q)$$

Q^* is the closest distribution from subset of \mathcal{Q} to P

- easy to state, generally hard to solve (i.e., to find Q^*)

Specification of $\mathcal{Q} \subset \mathcal{P}$

- via **parametrized densities** $\mathcal{Q} = \{p_\theta, \theta \in \Theta\}$
- via **parametrized transformations**
 X has some simple distribution which is easy to sample from and is transformed to a complex one using a deterministic function G
e.g., let $X \sim N(0, 1)$ then $X^2 \sim \chi^2(1)$ and $G(z) = z^2$
- \mathcal{Q} is given by set of parametrized functions $G_\theta, \theta \in \Theta$
(**neural networks parametrized via their weights**)
- **easy sampling** from $G_\theta(X)$, sample $x \sim X$ (easy)
and then pass x through $G_\theta(X)$, i.e., compute $G_\theta(x)$
- **How to solve the information projection problem?**

Maximum likelihood estimation

- **task**

given the set of data $\{\mathbf{x}_i \sim P_D\}_{i=1}^n$, describe distribution P_D

- **MLE estimate** $P_D \in P_\theta = \{P_\theta, \theta \in \Theta\}$

assume that P_θ has density, i.e., $dP_\theta = p_\theta(\mathbf{x}) d\mathbf{x}$

assume that \mathbf{x}_i i.i.d.

search for optimal $\theta_{\text{mle}} \in \Theta$ and then set $P_D = P_{\theta_{\text{mle}}}$

$$\theta_{\text{mle}} = \operatorname{argmax}_\theta \mathbb{E}_{\mathbf{x} \sim P_D} \log p_\theta(\mathbf{x})$$

$$\text{estimate } \theta_{\text{mle}}^* = \operatorname{argmax}_\theta \frac{1}{n} \sum_{i=1}^n \log p_\theta(\mathbf{x}_i)$$

- **optimization in terms of KL-divergence**

$$\theta_{\text{mle}} = \operatorname{argmin}_\theta KL(P_D(\mathbf{x}) || P_\theta(\mathbf{x}))$$

$$= \operatorname{argmin}_\theta \int p_D(\mathbf{x}) \frac{p_D(\mathbf{x})}{p_\theta(\mathbf{x})} d\mathbf{x}$$

MLE in terms of KL-divergence

- best approximation of P_D using P_θ
 - \hat{P}_D proxy for P_D , $\hat{P}_D(d\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}(d\mathbf{x})$ (Dirac m.)
 - P_θ - model distribution with density $p_{\text{model}}(\mathbf{x}|\theta)$

- maximization MLE = minimization of $KL(P_D||P_\theta)$

$$\begin{aligned} KL(P_D||P_\theta) &= \int \log \frac{dP_D}{dP_\theta} dP_D = \int \log \frac{p_D(\mathbf{x})}{p_\theta(\mathbf{x})} dP_D \\ &= \int \log p_D(\mathbf{x}) dP_D - \int \log p_\theta(\mathbf{x}) dP_D \\ &\approx -H[P_D] - \int p_\theta(\mathbf{x}) d\hat{P}_D \quad (P_D \approx \hat{P}_D) \\ &\propto - \int \log p_\theta(\mathbf{x}) d\hat{P}_D \quad (\text{integration over Dirac}) \\ &\propto - \underbrace{\frac{1}{n} \sum_{i=1}^n \log p_\theta(\mathbf{x}_i)}_{=\text{MLE}} \end{aligned}$$

Generative modeling

- **purpose**

given data from an unknown distribution $\mathbf{x} \sim p(\mathbf{x})$
model $p(\mathbf{x})$ using a differentiable mapping G so that

$$p(\mathbf{x}) \sim G_{\theta_g}(p(\mathbf{z})) = G(p(\mathbf{z}); \theta_g)$$

where $p(\mathbf{z})$ is a selected, simple prior, e.g. mv Gaussian

- **maximum likelihood estimation** direct setting of density under i.i.d. assumption, **KL divergence minimization**

Generative modeling

- solution to the information projection problem
KL-divergence minimalization
via playing an adversarial game between
generator and discriminator



Partial criteria

- an ideal discriminator

$D : \mathbf{x} \in \mathbb{R}^d \rightarrow (0, 1)$, i.e., $\log D : \mathbf{x} \rightarrow (-\infty, 0)$

we would like $D_{\theta_d}(\mathbf{x}^{real}) \rightarrow 1$, $D_{\theta_d}(\mathbf{x}^{fake}) \rightarrow 0$

i.e., maximize w.r.t. θ_d

$$\log(D_{\theta_d}(\mathbf{x}^{real})) + \log((1 - D_{\theta_d}(\mathbf{x}^{fake})))$$

- an ideal generator

generator wants to fool discriminator,

i.e., it generates \mathbf{x}^{fake} so that $D_{\theta_d}(\mathbf{x}^{fake}) \rightarrow 1$

tune weights θ_g of the generator to minimize

$$\log(1 - D_{\theta_d}(\mathbf{x}^{fake})) = \log(1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})))$$

w.r.t θ_g for θ_d fixed

Compound criterion

- compound criterion

$$V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D_{\theta_d}(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_z(\mathbf{x})} [\log(1 - D_{\theta_d}(G_{\theta_g}(z)))]$$

- minimax optimization - set θ_d, θ_g using

$$\min_{\theta_g} \max_{\theta_d} V(D_{\theta_d}, G_{\theta_g})$$

- alternate optimization

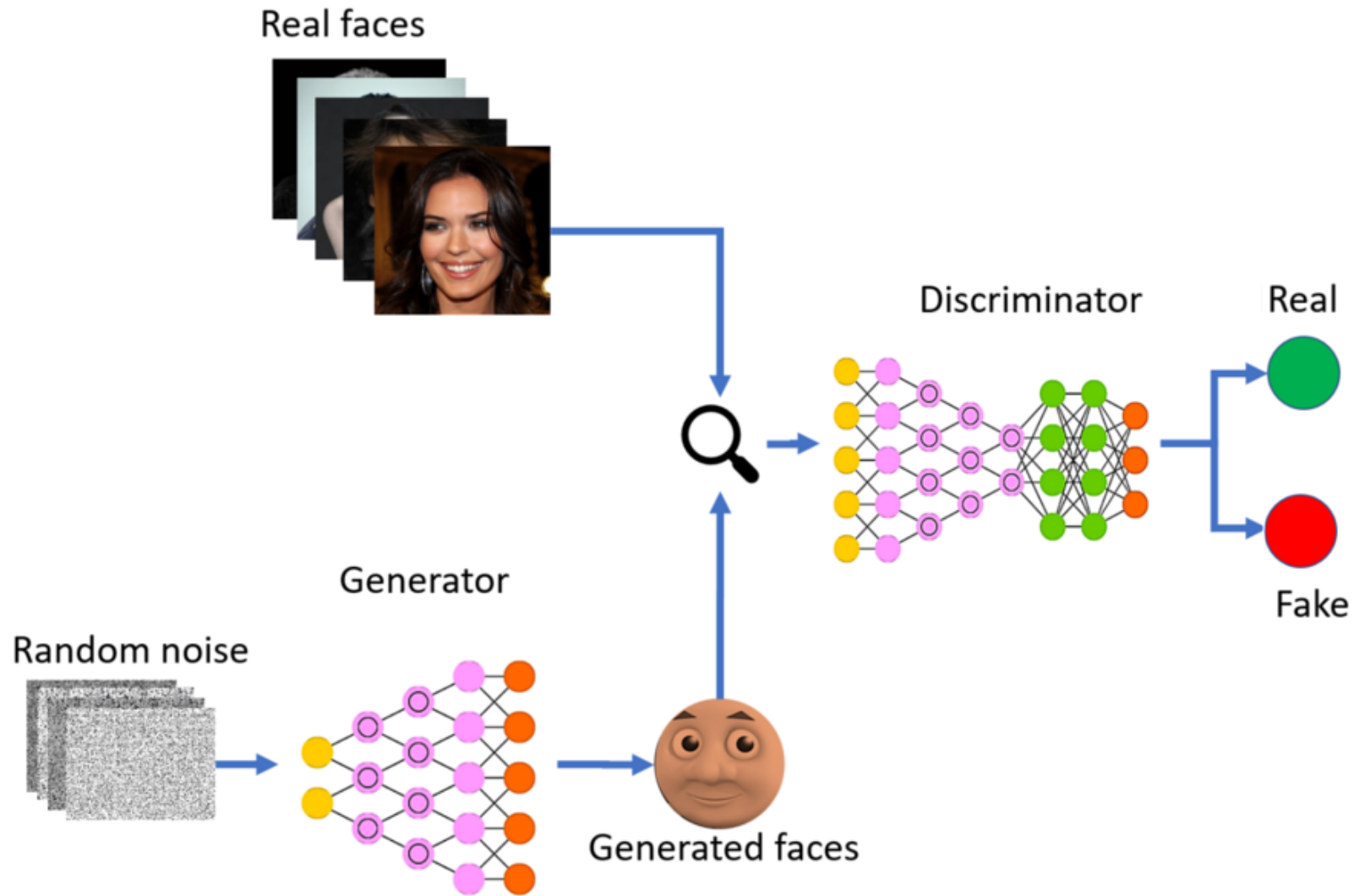
- for fixed generator G_{θ_g} maximize $V(D_{\theta_d}, \cdot)$
- for fixed discriminator D_{θ_d} minimize $V(\cdot, G_{\theta_g})$

Theoretical analysis

- **Proposition.** Optimizing $\min_G \max_D V(D, G)$ corresponds to minimizing $JSD(p_{\text{data}} || p_g)$, which is minimal ($=0$) if and only if $p_{\text{data}} = p_g$

source: <https://arxiv.org/abs/1406.2661>

A GAN concept



source: <https://medium.com/sigmoid/a-brief-introduction-to-gans>

Learning algorithm

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

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

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

end for

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

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

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

MNIST dataset

- 60000 - 28x28 greyscale **images of handwritten digits**
<http://yann.lecun.com/exdb/mnist/>



MNIST dataset

- 60000 - 28x28 greyscale images of handwritten digits
GAN architecture: D,G - perceptron networks



MNIST dataset

- 60000 - 28x28 greyscale images of handwritten digits GAN architecture: D,G - convolution networks

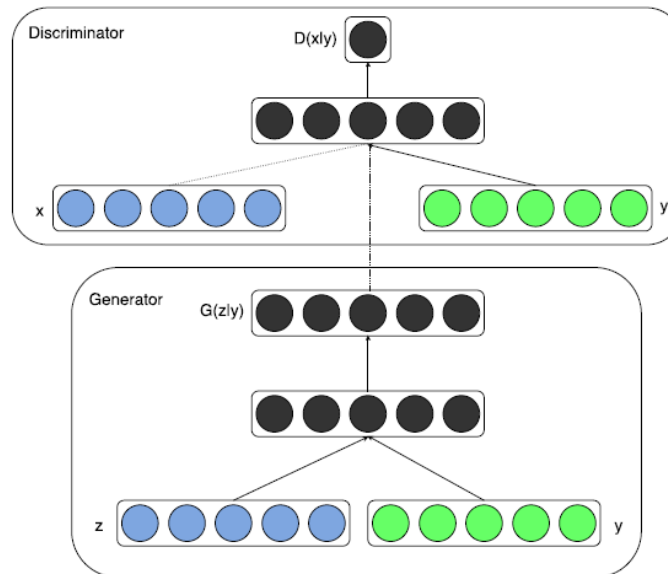


cGAN - 2014

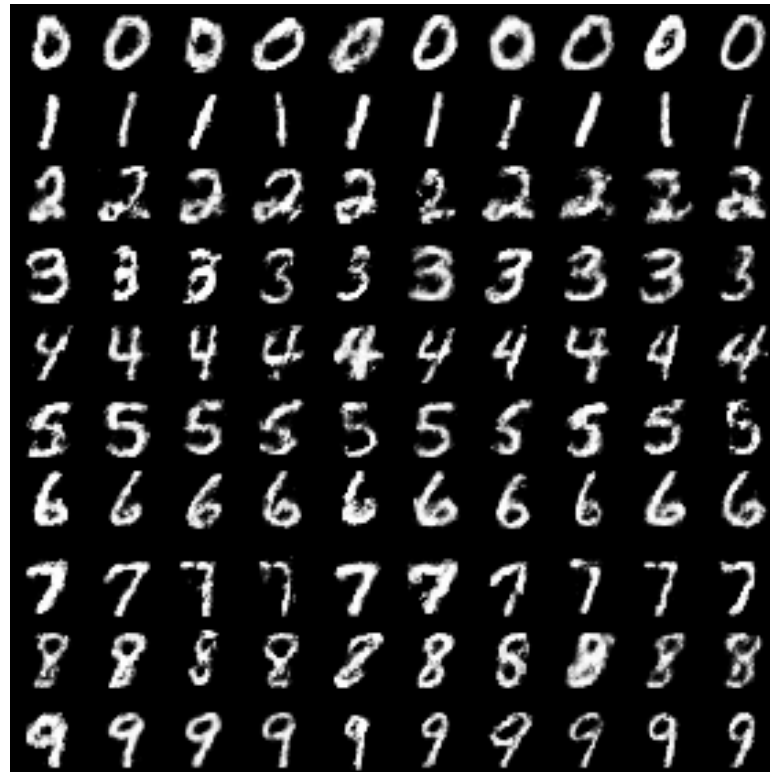
- *Conditional Generative Adversarial Nets* <https://arxiv.org/abs/1411.1784>
- unconditional vs. conditional GAN, y – *condition*

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] &+ \mathbb{E}_{\mathbf{x} \sim p_z(\mathbf{x})} [\log(1 - D(G(\mathbf{z})))] \\ \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] &+ \mathbb{E}_{\mathbf{x} \sim p_z(\mathbf{x})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))] \end{aligned}$$

- conditioning by extending latent variable of generator

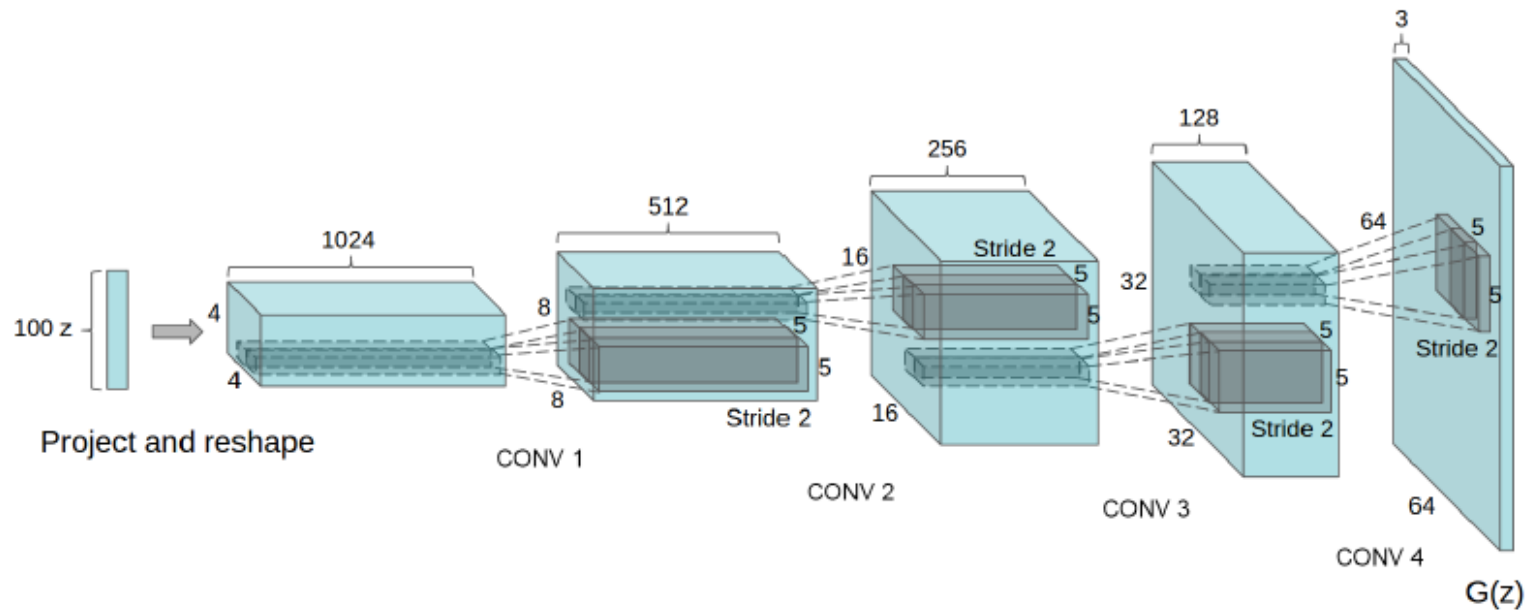


MNIST dataset



DCGAN - 2015

- *Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks* <https://arxiv.org/abs/1511.06434>
- architecture - uses convolutional layers

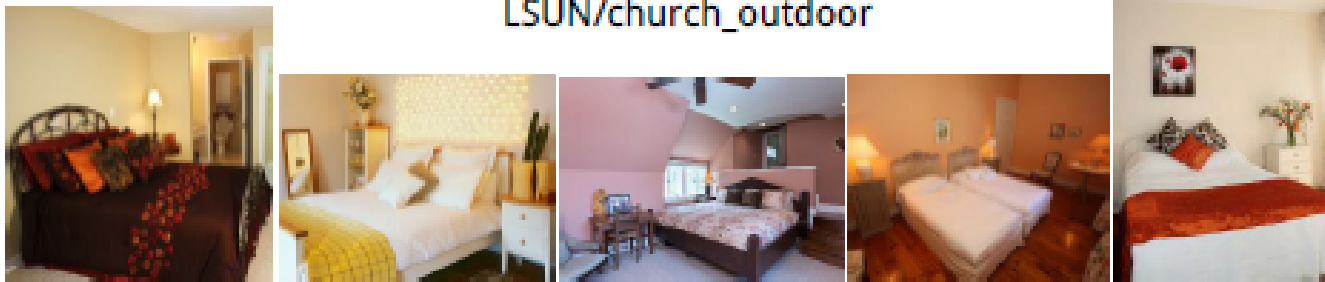


LSUN dataset

- 10 - categories, (church_outdoor, bedroom, bridge ...)
https://www.yf.io/p/l_sun



LSUN/church_outdoor



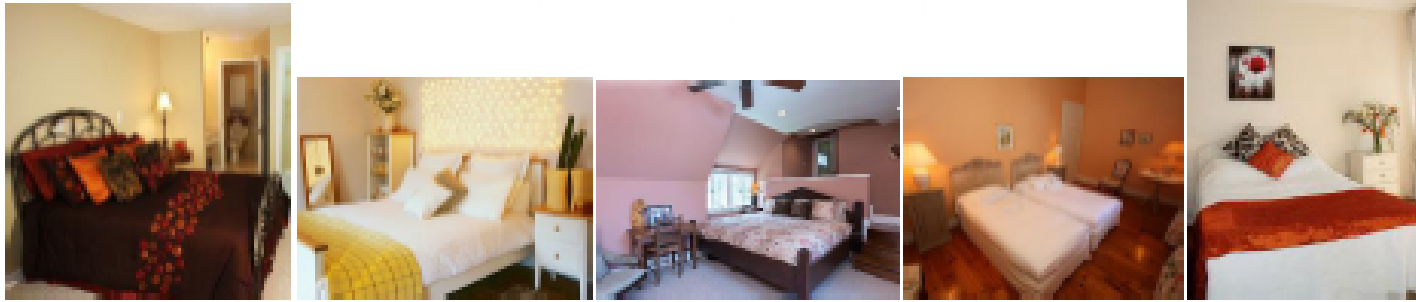
LSUN/bedroom

DCGAN - 2015



Figure 3: Generated bedrooms after five epochs of training. There appears to be evidence of visual under-fitting via repeated noise textures across multiple samples such as the base boards of some of the beds.

DCGAN - 2015



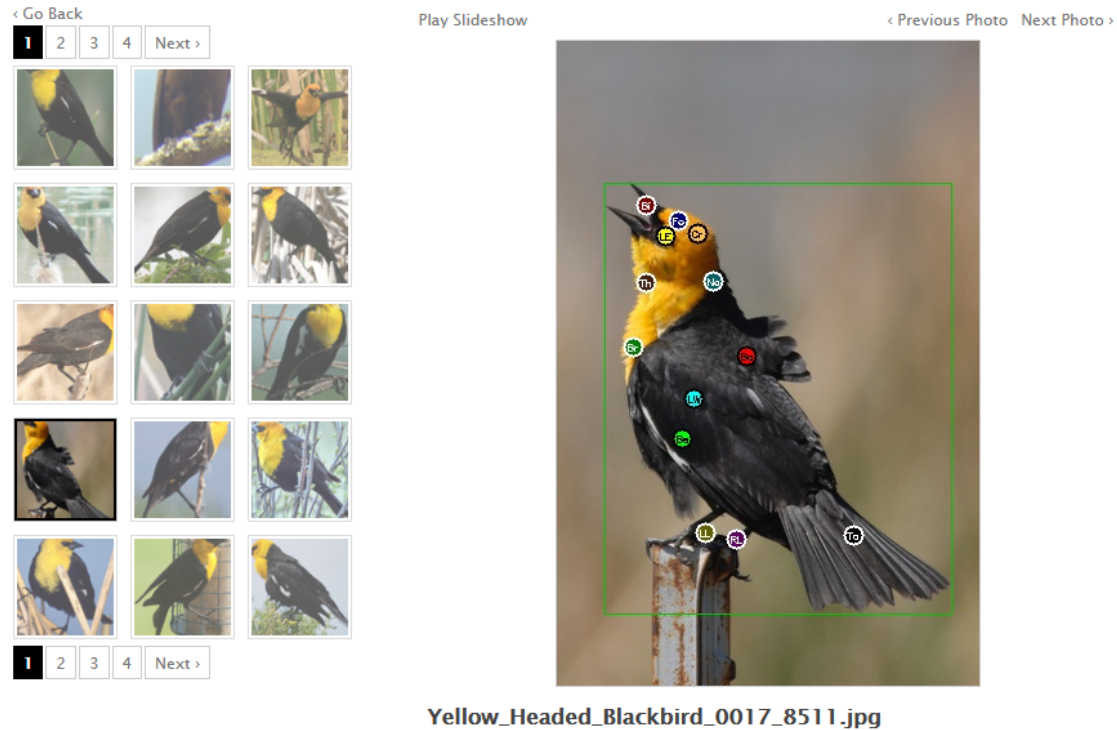
LSUN/bedroom



StackGAN - 2016

- *StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks* <https://arxiv.org/abs/1612.03242>
- **Caltech-UCSD Birds 200 Dataset**
<http://www.vision.caltech.edu/visipedia/CUB-200-2011.html>
- **102 Category Flower Dataset**
<https://www.robots.ox.ac.uk/vgg/data/flowers/102/>

StackGAN - 2016



- a bird has a bright golden crown and throat, it's breast is yellow, and back is black
- upper body yellow and lower black with black color around beak
- this bird has a bright yellow crown, a long straight bill, and white wingbars
- this is a black bird with a yellow head and breast ...

StackGAN - 2016

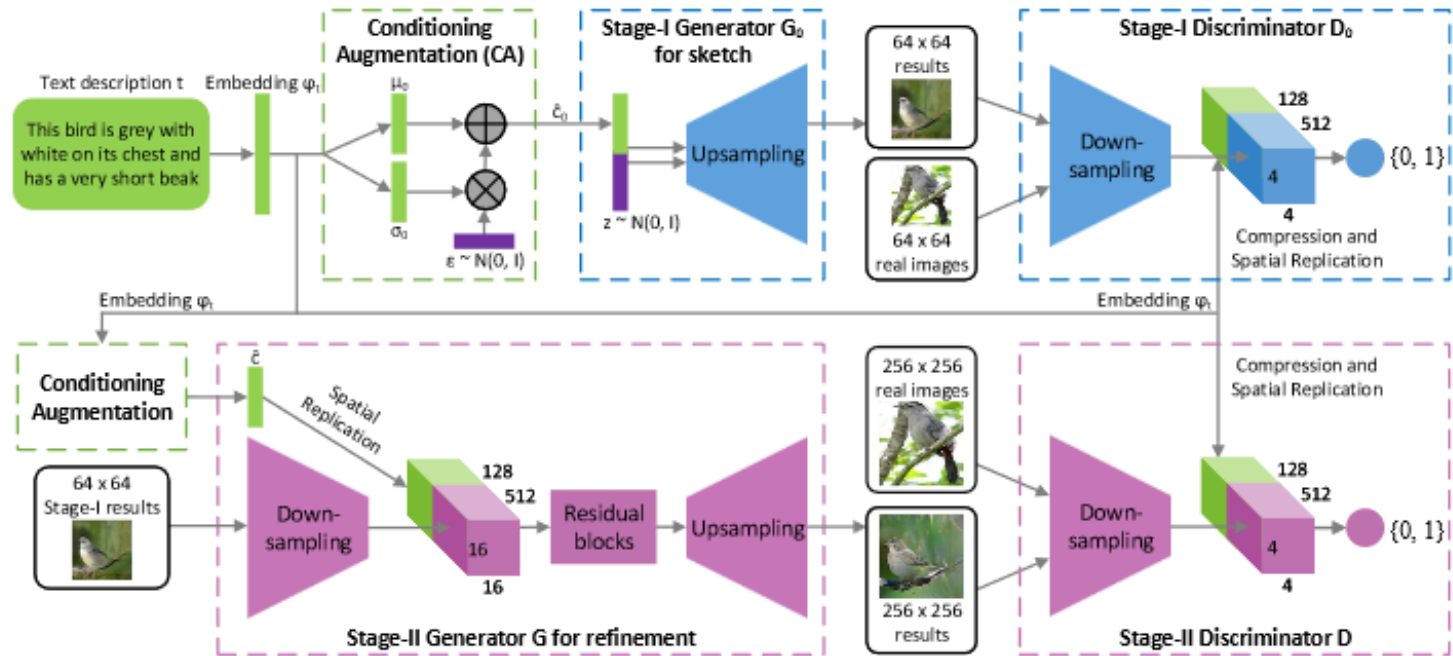


Figure 2. The architecture of the proposed StackGAN. The Stage-I generator draws a low-resolution image by sketching rough shape and basic colors of the object from the given text and painting the background from a random noise vector. Conditioned on Stage-I results, the Stage-II generator corrects defects and adds compelling details into Stage-I results, yielding a more realistic high-resolution image.

StackGAN - 2016

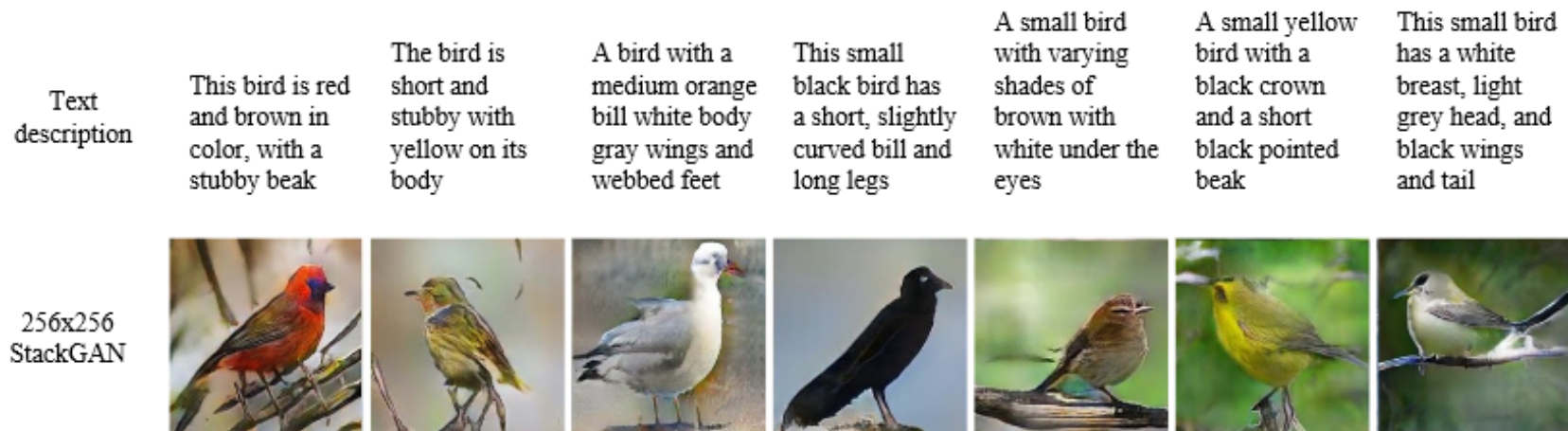


Figure 3. Example results by our StackGAN conditioned on text descriptions from CUB test set.

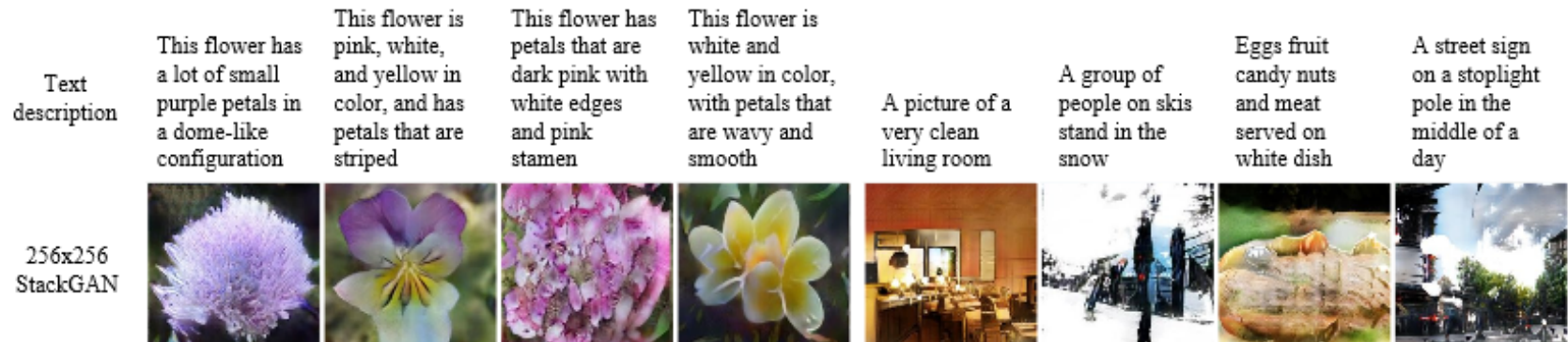


Figure 4. Example results by our StackGAN conditioned on text descriptions from Oxford-102 test set and COCO validation set

StackGAN - 2016






Text description	This bird is blue with white and has a very short beak	This bird has wings that are brown and has a yellow belly	A white bird with a black crown and yellow beak	This bird is white, black, and brown in color, with a brown beak	The bird has small beak, with reddish brown crown and gray belly	This is a small, black bird with a white breast and white on the wingbars.	This bird is white black and yellow in color, with a short black beak
Stage-I images							
Stage-II images							

Figure 5. Samples generated by our StackGAN from unseen texts in CUB test set. Each column lists the text description, images generated from the text by Stage-I and Stage-II of StackGAN.

- <https://github.com/hanzhanggit/StackGAN>

BEGAN - 2017

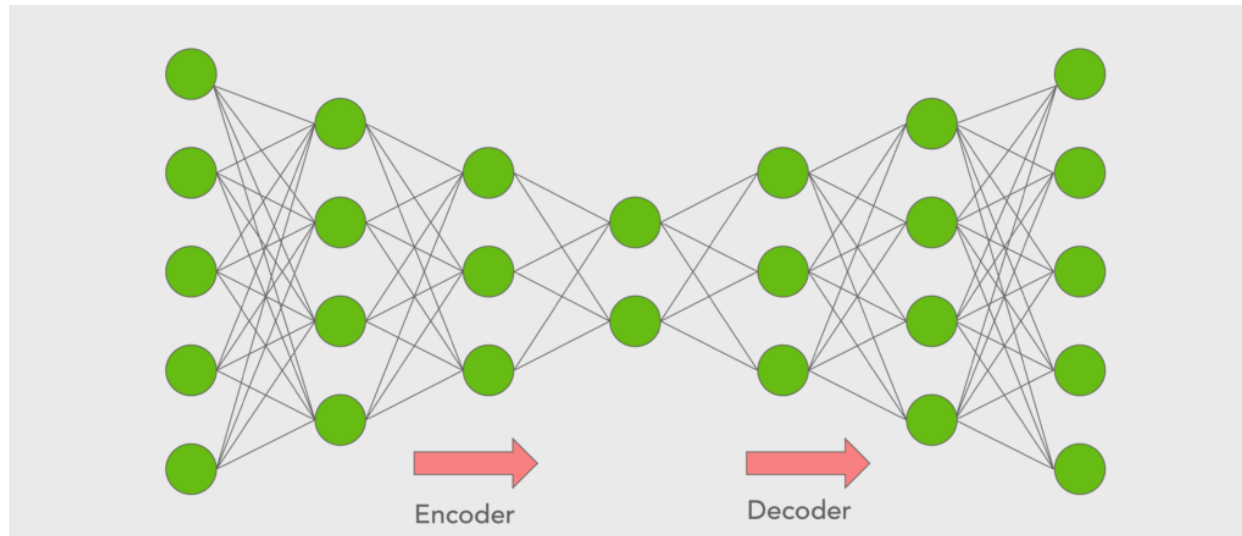
- *BEGAN: Boundary Equilibrium Generative Adversarial Networks*
<https://arxiv.org/abs/1703.10717>
- **energy based GAN**, discriminator assigns **low energy values to real data** and **high to fake ones** - generalized view of loss functions, training - loss minimization

$$V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [D_{\theta_d}(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_z(\mathbf{x})} [(m - D_{\theta_d}(G_{\theta_g}(z)))_+]$$

where m is a positive margin, $(\cdot)_+ = \max(0, \cdot)$ and $0 \leq D_{\theta_d}$

BEGAN - 2017

- discriminator as autonecoder



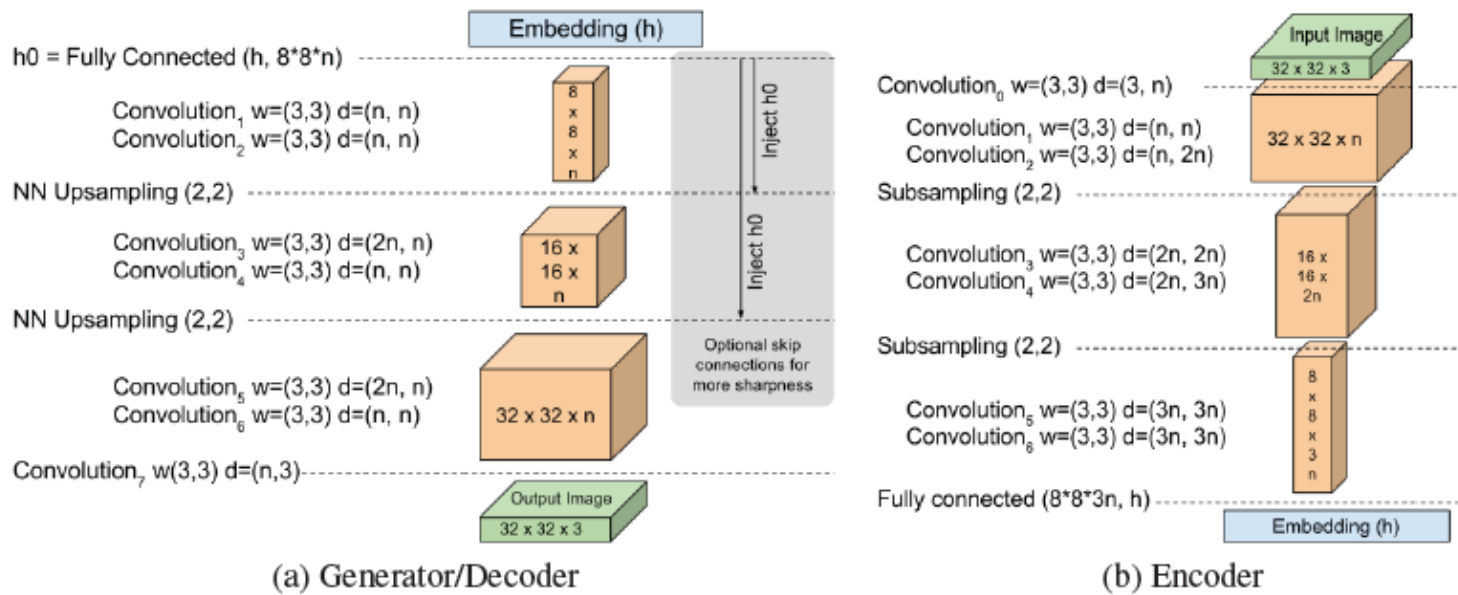
source: <https://www.mygreatlearning.com/blog/autoencoder/>

- **loss - reconstruction errors** for real and fake images

$$D_{\theta}(\mathbf{x}_{real}) = \|\text{Dec}(\text{Enc}(\mathbf{x}_{real})) - \mathbf{x}_{real}\| \rightarrow 0$$
$$D_{\theta}(\mathbf{x}_{fake}) = \|\text{Dec}(\text{Enc}(\mathbf{x}_{fake})) - \mathbf{x}_{fake}\| \rightarrow \infty$$

BEGAN - 2017

- architecture of generator/decoder and encoder



CelebA dataset

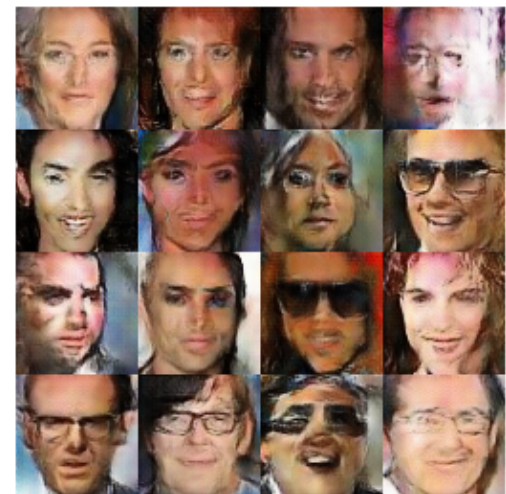
- CelebA dataset - 202599 annotated (40 attributes) celebrity portraits



- <http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

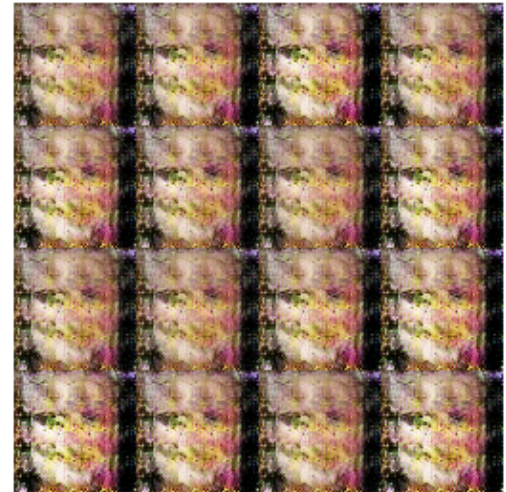
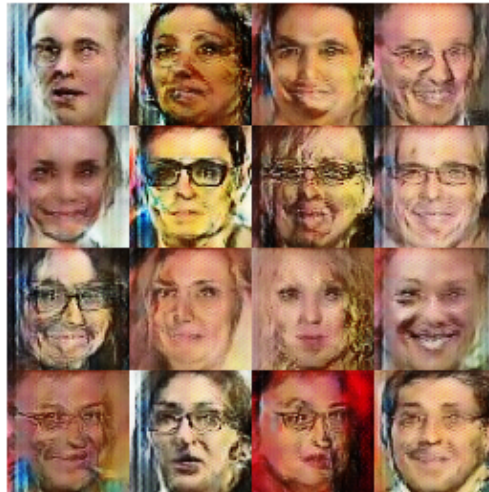
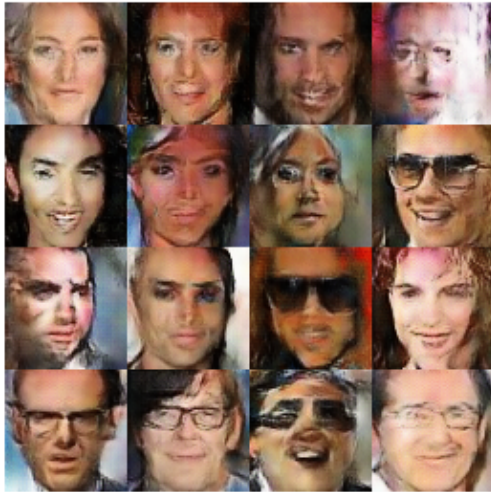
CelebA - DCGAN

- learning progress



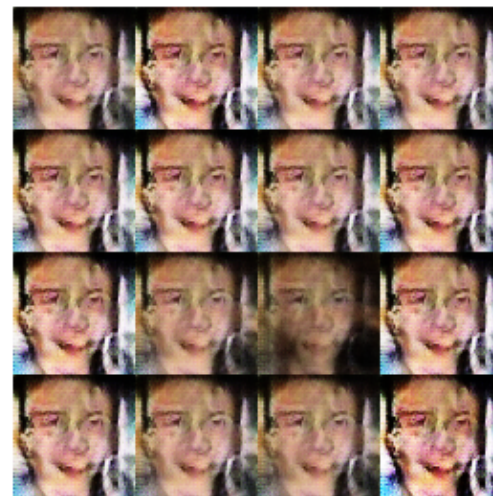
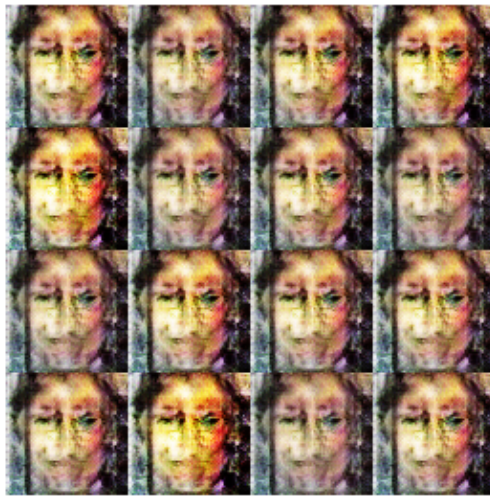
CelebA - DCGAN

- deterioration starts



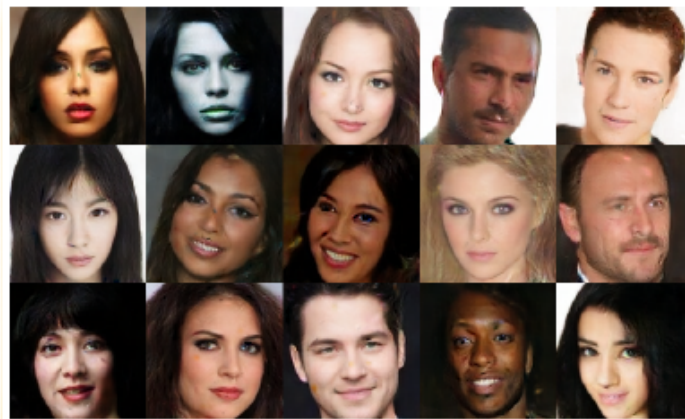
DCGAN

- mode collapse



BEGAN - 2017

- generated fake images



(b) Our results (128x128)



Figure 3: Random 64x64 samples at varying $\gamma \in \{0.3, 0.5, 0.7\}$

PGGAN - 2017

- *Progressive Growing of GANs for Improved Quality, Stability, and Variation* <https://arxiv.org/abs/1710.10196> (NVIDIA)
- **CelabA HQ dataset** - 30000 imgs at 1024x1024 resolution

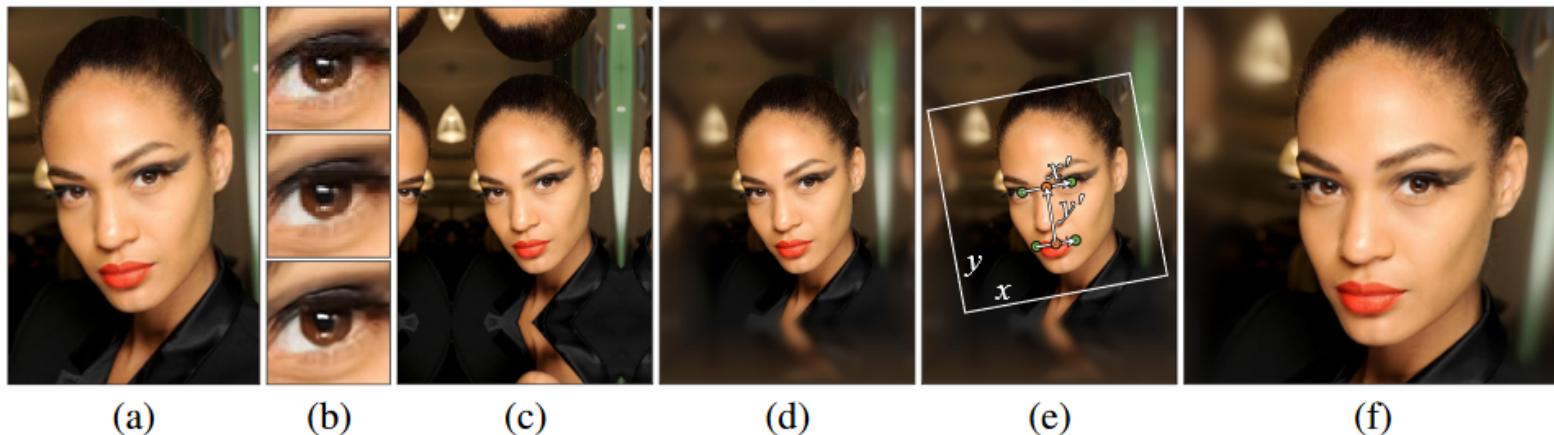


Figure 8: Creating the CELEBA-HQ dataset. We start with a JPEG image (a) from the CelebA in-the-wild dataset. We improve the visual quality (b,top) through JPEG artifact removal (b,middle) and 4x super-resolution (b,bottom). We then extend the image through mirror padding (c) and Gaussian filtering (d) to produce a visually pleasing depth-of-field effect. Finally, we use the facial landmark locations to select an appropriate crop region (e) and perform high-quality resampling to obtain the final image at 1024×1024 resolution (f).

PGGAN - 2017

- *Progressive Growing of GANs for Improved Quality, Stability, and Variation* <https://arxiv.org/abs/1710.10196>
- **architecture - progressive growing** of convolutional layers

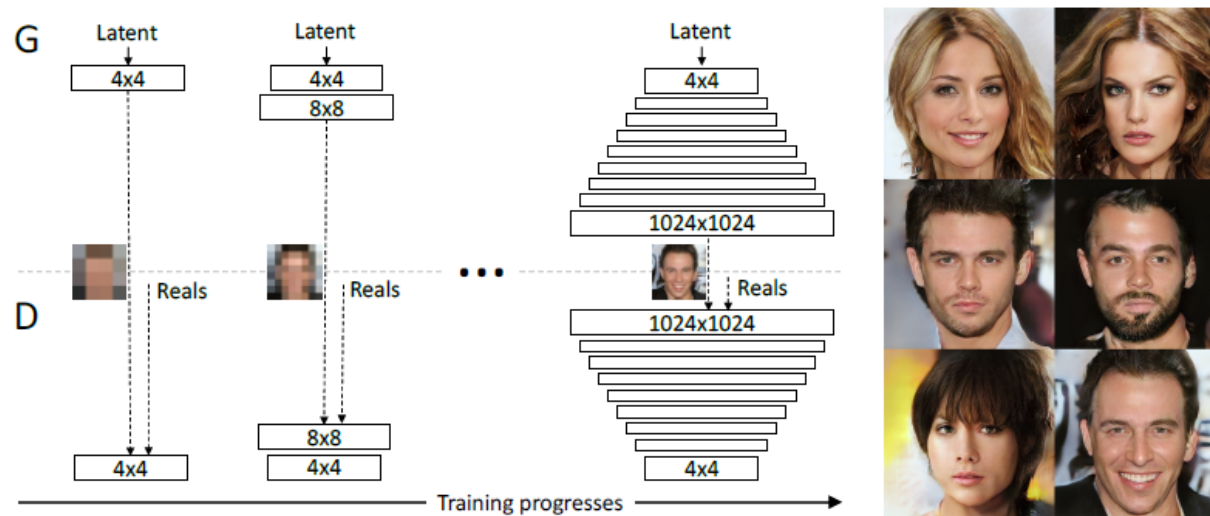


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout 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. On the right we show six example images generated using progressive growing at 1024×1024 .

PGGAN - 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.

- https://github.com/tkarras/progressive_growing_of_gans

StyleGAN - 2018

- A Style-Based Generator Architecture for Generative Adversarial Networks <https://arxiv.org/abs/1812.04948>

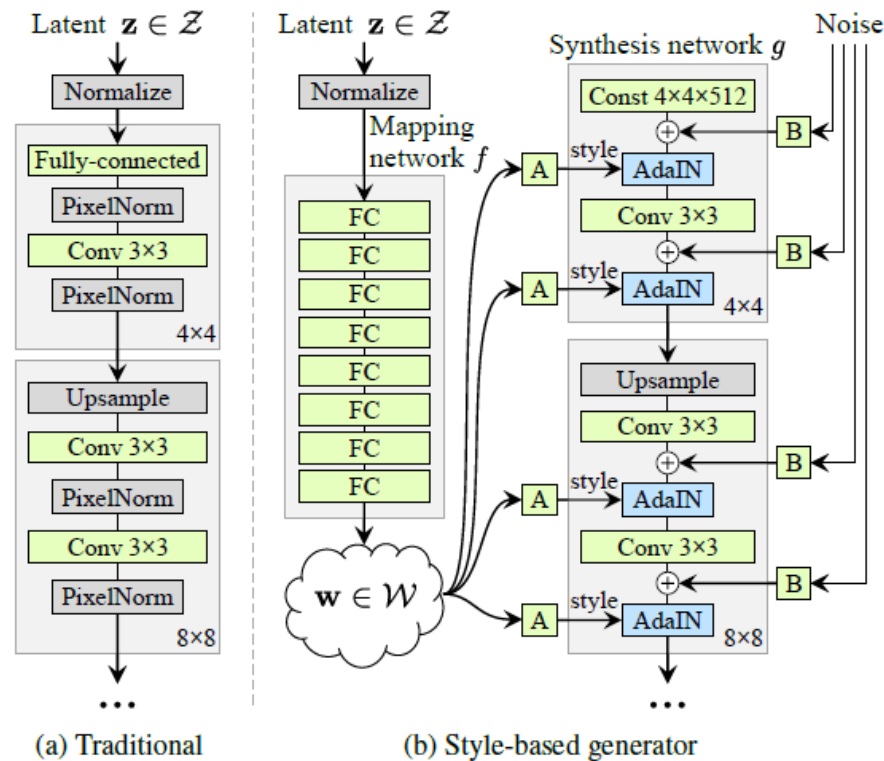


Figure 1. While a traditional generator [30] feeds the latent code through the input layer only, we first map the input to an intermediate latent space \mathcal{W} , which then controls the generator through adaptive instance normalization (AdaIN) at each convolution layer. Gaussian noise is added after each convolution

StyleGAN - 2018

- style disentanglement

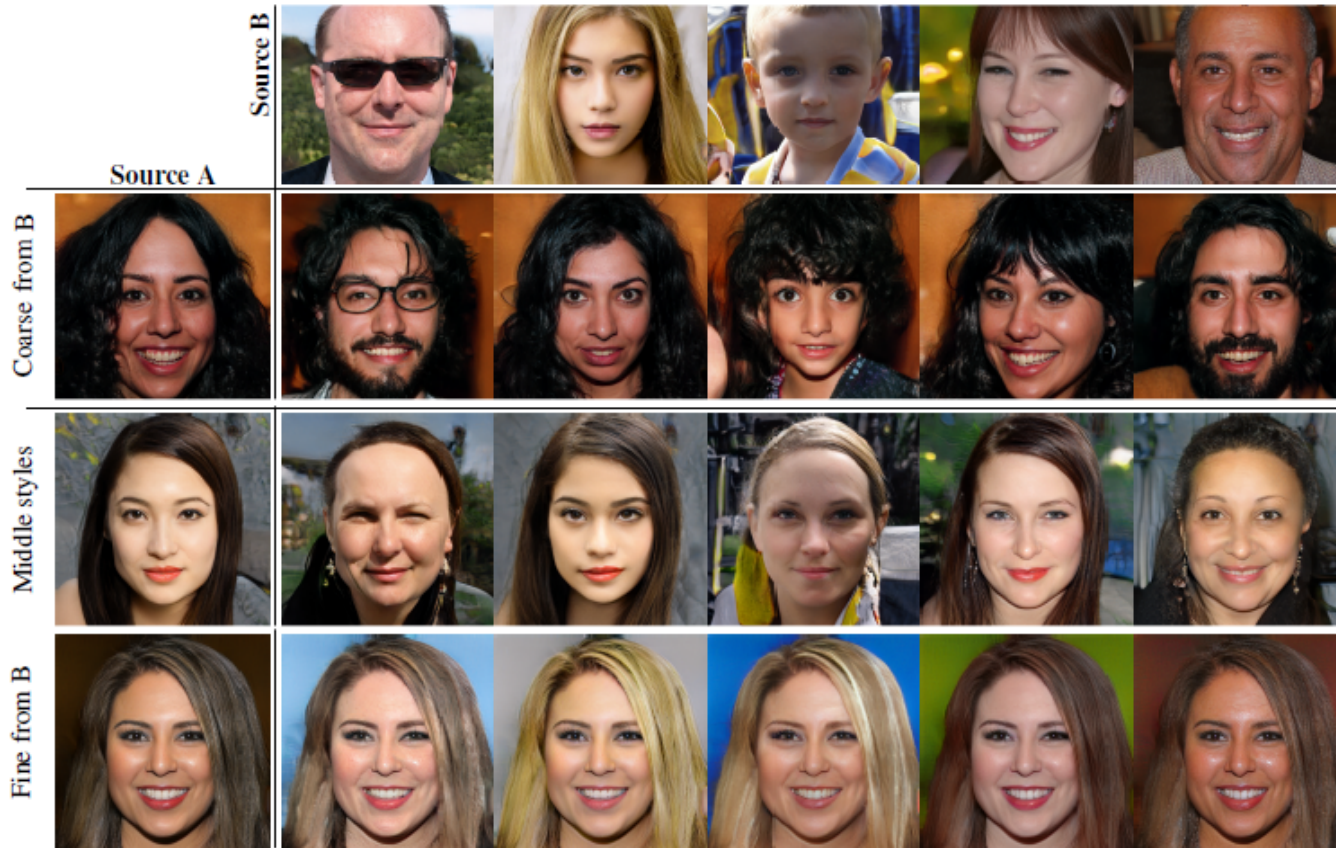


Figure 3. Two sets of images were generated from their respective latent codes (sources A and B); the rest of the images were generated by copying a specified subset of styles from source B and taking the rest from source A. Copying the styles corresponding to coarse spatial resolutions ($4^2 - 8^2$) brings high-level aspects such as pose, general hair style, face shape, and eyeglasses from source B, while all colors

StyleGAN2 - 2019

- *Analyzing and Improving the Image Quality of StyleGAN*
<https://arxiv.org/abs/1912.04958>



Figure 1. Instance normalization causes water droplet -like artifacts in StyleGAN images.

StyleGAN2 - 2019

- Analyzing and Improving the Image Quality of StyleGAN

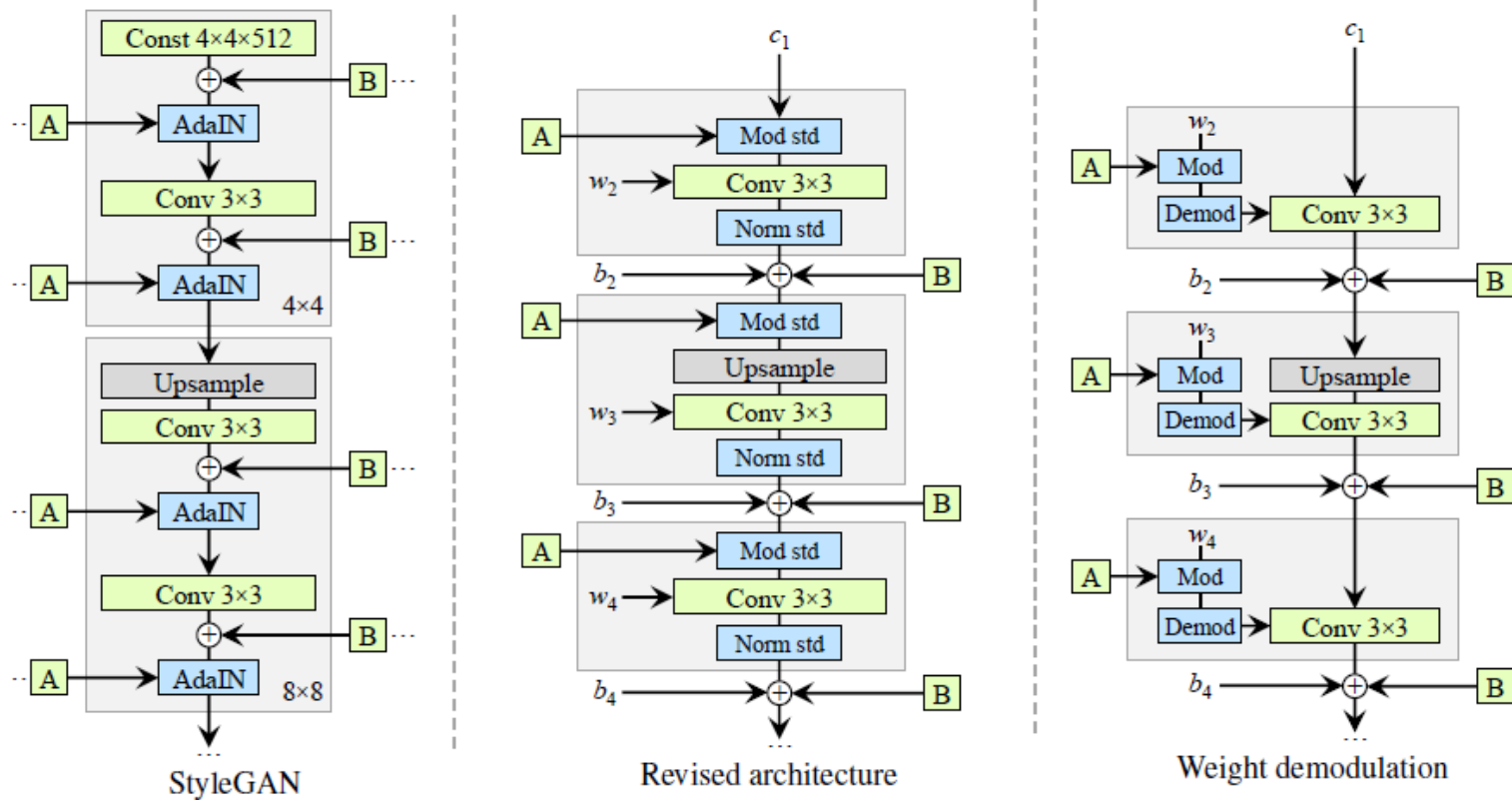


Figure 2. We redesign the architecture of the StyleGAN synthesis network.

StyleGAN2 - 2019

- *abandonment of progressive growing*



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.

- <https://github.com/NVlabs/stylegan2>

StyleGAN-ADA - 2020

- *Training Generative Adversarial Networks with Limited Data Using Adaptive Discriminator Adaptation - (ADA)*

<https://arxiv.org/abs/2006.06676>

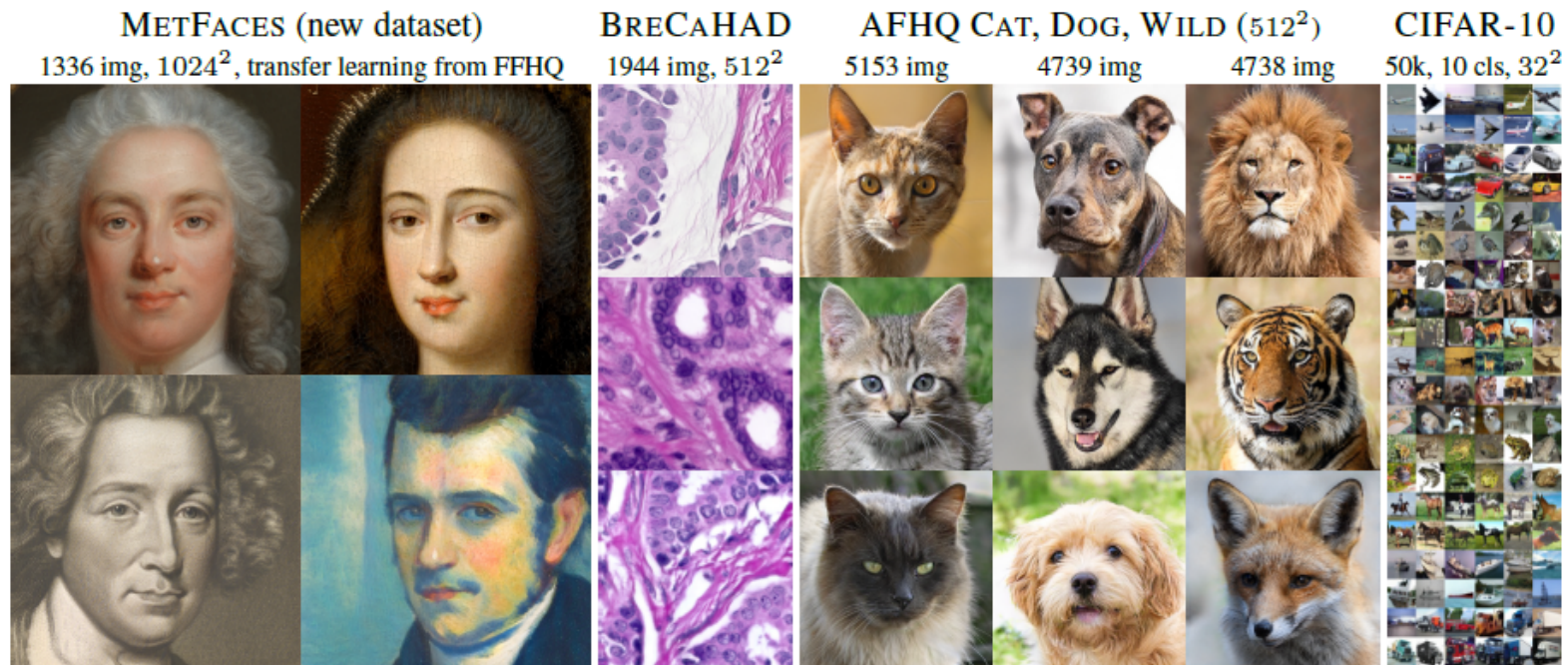
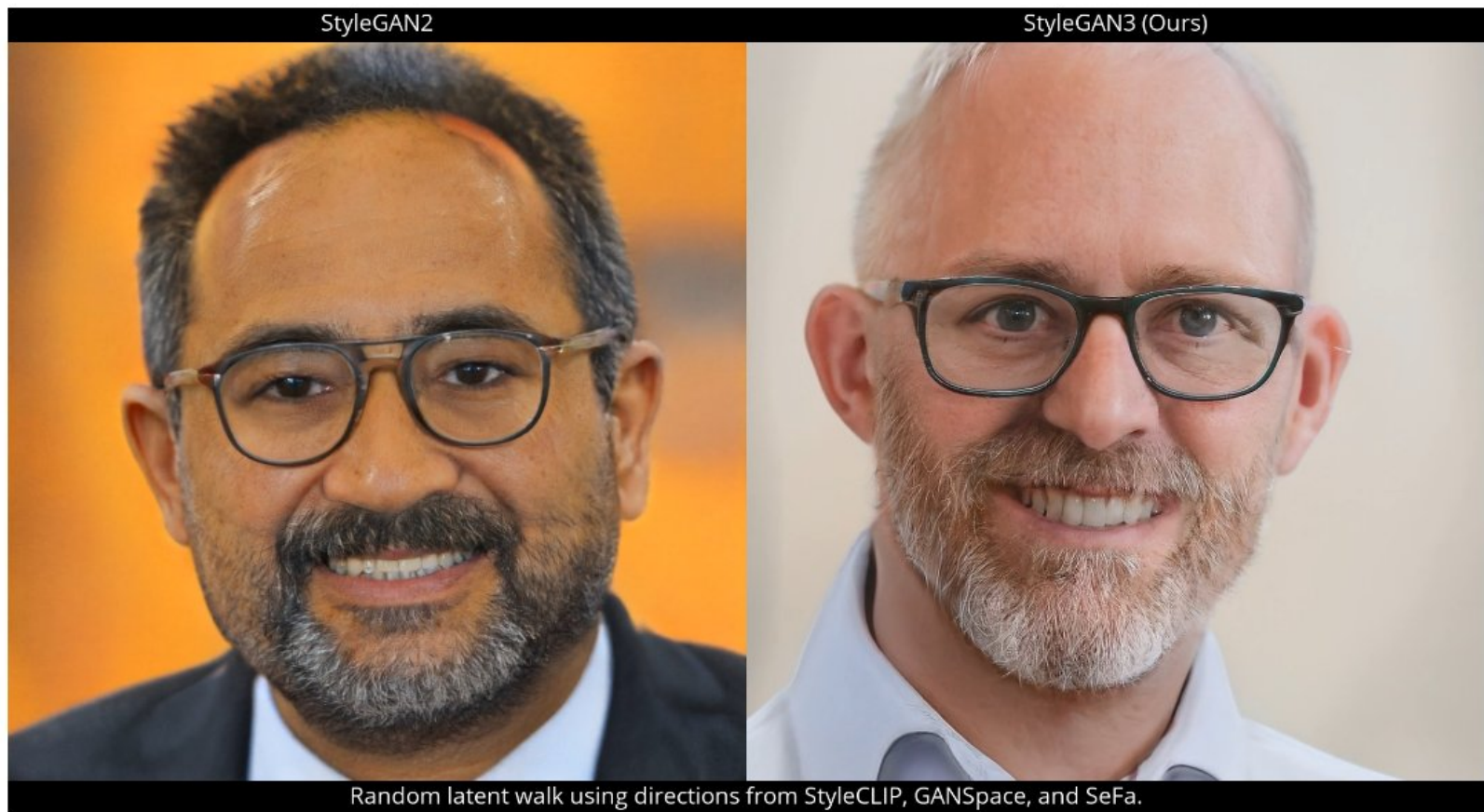


Figure 10: Example generated images for several datasets with limited amount of training data, trained using ADA.

- <https://github.com/NVlabs/stylegan>

StyleGAN3 - 2021

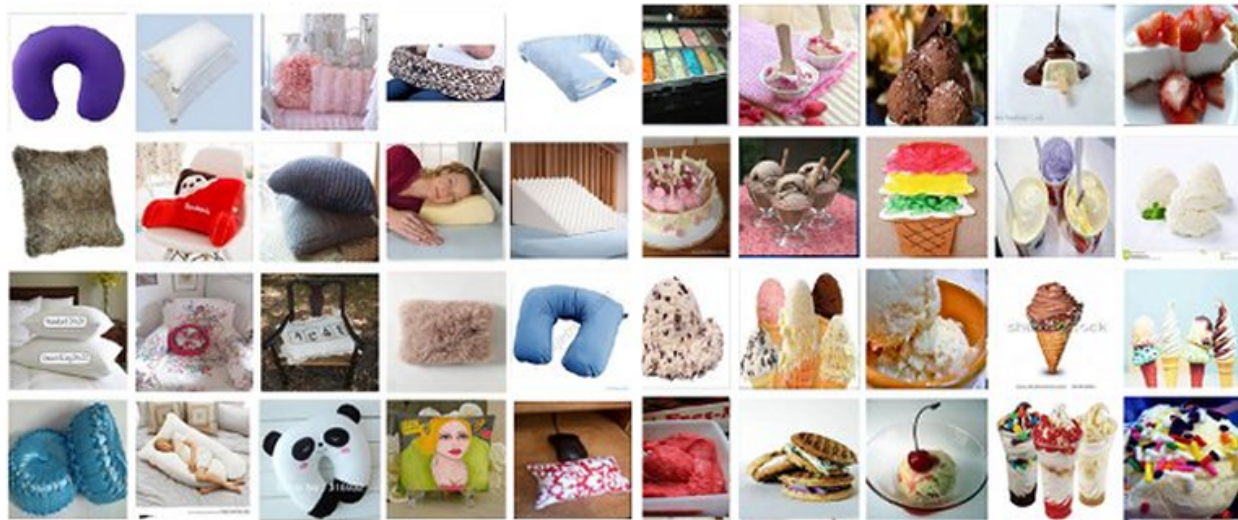
- *Alias-Free Generative Adversarial Networks (StyleGAN3)*
<https://arxiv.org/abs/2106.12423>



- <https://nvlabs.github.io/stylegan3>

ImageNet

- over 14 mil. of images from 20 thousand categories based on the WordNet database (a dictionary)



Pillow

Icecream

BigGAN - 2019

- *Large Scale GAN Training for High Fidelity Natural Image Synthesis*
<https://arxiv.org/abs/1809.11096>
- we show that GANs benefit dramatically from **scaling**, and train models with **two to four times as many parameters** and **eight times the batch size** compared to prior art
- training on 128 to 512 cores of a **Google TPUv3 Pod**

Batch	Ch.	Param (M)	Shared	Skip-z	Ortho.	Itr $\times 10^3$	FID	IS
256	64	81.5	SA-GAN Baseline			1000	18.65	52.52
512	64	81.5	✗	✗	✗	1000	15.30	58.77(± 1.18)
1024	64	81.5	✗	✗	✗	1000	14.88	63.03(± 1.42)
2048	64	81.5	✗	✗	✗	732	12.39	76.85(± 3.83)
2048	96	173.5	✗	✗	✗	295(± 18)	9.54(± 0.62)	92.98(± 4.27)
2048	96	160.6	✓	✗	✗	185(± 11)	9.18(± 0.13)	94.94(± 1.32)
2048	96	158.3	✓	✓	✗	152(± 7)	8.73(± 0.45)	98.76(± 2.84)
2048	96	158.3	✓	✓	✓	165(± 13)	8.51(± 0.32)	99.31(± 2.10)
2048	64	71.3	✓	✓	✓	371(± 7)	10.48(± 0.10)	86.90(± 0.61)

Table 1: Fréchet Inception Distance (FID, lower is better) and Inception Score (IS, higher is better) for ablations of our proposed modifications. *Batch* is batch size, *Param* is total number of parameters, *Ch.* is the channel multiplier representing the number of units in each layer, *Shared* is using shared embeddings, *Skip-z* is using skip connections from the latent to multiple layers, *Ortho.* is Orthogonal Regularization, and *Itr* indicates if the setting is stable to 10^6 iterations, or it collapses at the given iteration. Other than rows 1-4, results are computed across 8 random initializations.

BigGAN - 2019

- architecture - convolutional layers, no pg

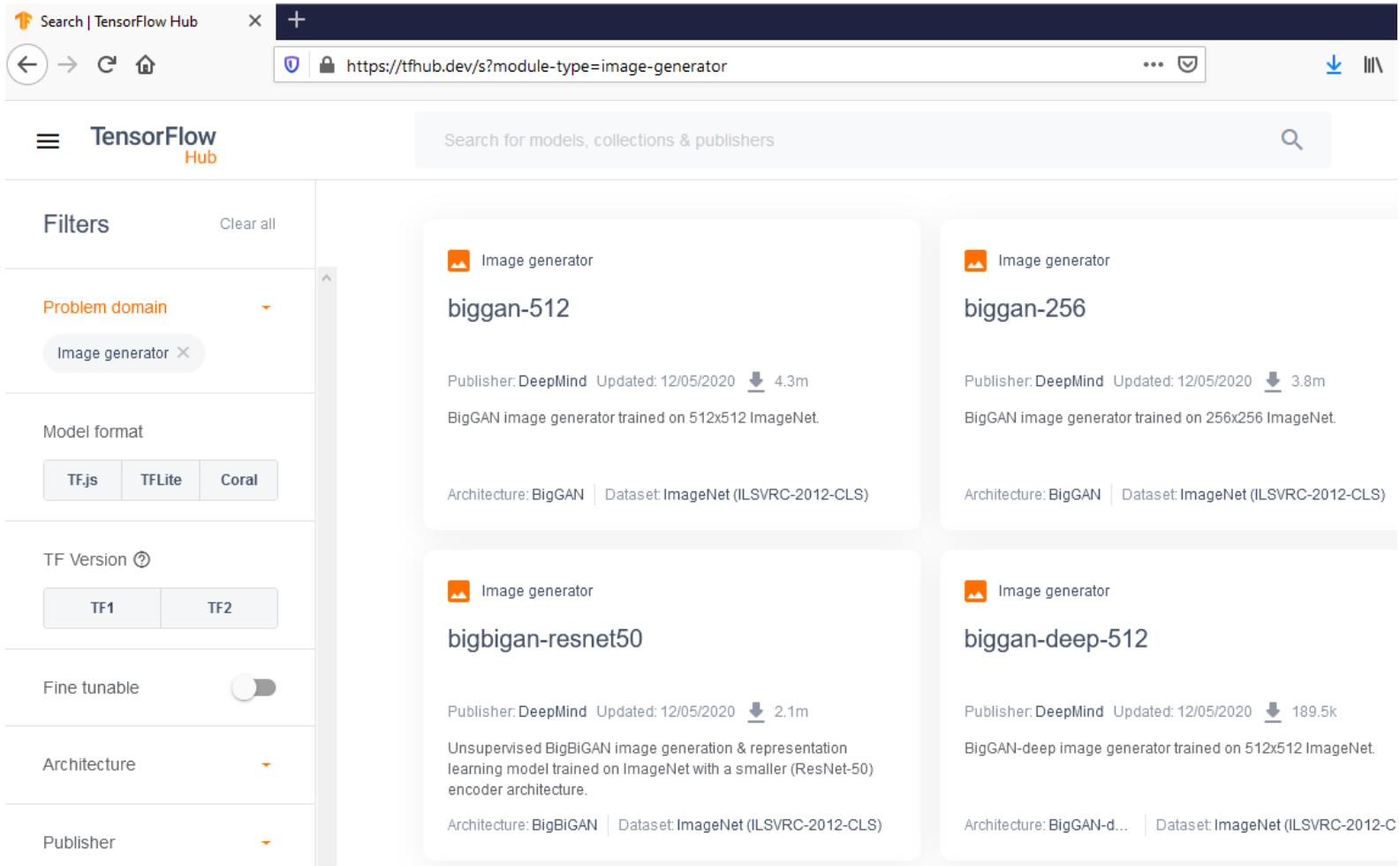


Figure 1: Class-conditional samples generated by our model.



BigGAN - 2019

- TensorFlow Hub - pretrained weights



The screenshot shows the TensorFlow Hub interface with search results for image generators. The browser address bar shows the URL <https://tfhub.dev/s?module-type=image-generator>. The search bar contains the text "Search for models, collections & publishers". The left sidebar shows filters for "Problem domain" (Image generator), "Model format" (TF.js, TFLite, Coral), "TF Version" (TF1, TF2), "Fine tunable" (toggle), "Architecture", and "Publisher". The main content area displays four model cards:

- biggan-512**: Image generator, Publisher: DeepMind, Updated: 12/05/2020, 4.3m downloads. Description: BigGAN image generator trained on 512x512 ImageNet. Architecture: BigGAN | Dataset: ImageNet (ILSVRC-2012-CLS).
- biggan-256**: Image generator, Publisher: DeepMind, Updated: 12/05/2020, 3.8m downloads. Description: BigGAN image generator trained on 256x256 ImageNet. Architecture: BigGAN | Dataset: ImageNet (ILSVRC-2012-CLS).
- bigbigan-resnet50**: Image generator, Publisher: DeepMind, Updated: 12/05/2020, 2.1m downloads. Description: Unsupervised BigBiGAN image generation & representation learning model trained on ImageNet with a smaller (ResNet-50) encoder architecture. Architecture: BigBiGAN | Dataset: ImageNet (ILSVRC-2012-CLS).
- biggan-deep-512**: Image generator, Publisher: DeepMind, Updated: 12/05/2020, 189.5k downloads. Description: BigGAN-deep image generator trained on 512x512 ImageNet. Architecture: BigGAN-d... | Dataset: ImageNet (ILSVRC-2012-C).

OpenAI DALL-E - 2021

- DALL-E is a 12-billion parameter version of GPT-3 trained to generate images from text descriptions, using a dataset of text-image pairs.

TEXT PROMPT

an armchair in the shape of an avocado. . . .

AI-GENERATED
IMAGES



Edit prompt or view more images ↓

- *The supercomputer developed for OpenAI is a single system with more than 285,000 CPU cores, 10,000 GPUs and 400 gigabits per second of network connectivity for each GPU server.*
- <https://openai.com/blog/dall-e>

Open questions

- What **sorts of distributions** can GANs model?
- How can we scale GANs **beyond image synthesis?**
(text, audio, **computer-aided drug design** - <https://insilico.com>)
- What can we say about the **global convergence** of the training dynamics?
- How does GAN training **scale with batch size?**

source: <https://distill.pub/2019/gan-open-problems>