

GENERATIVE AI

Anders Eklund

anders.eklund@liu.se

Department of Biomedical Engineering (IMT)
Department of Computer and Information Science (IDA)
Center for Medical Image Science and Visualization (CMIV)
Linköping University, Sweden

February 7, 2023

OUTLINE

- ▶ Data augmentation, increase amount of training data
(can only modify existing images)
- ▶ Image synthesis, create completely new images
 - ▶ Noise-to-image GANs
 - ▶ Image-to-image GANs (image translation)
- ▶ GAN problems (mode collapse)
- ▶ How to evaluate GANs
- ▶ Training CNNs with synthetic images
- ▶ Diffusion models
- ▶ Synthetic data raise new ethical questions

IMAGE SYNTHESIS / DATA AUGMENTATION

- ▶ Add more realistic images to improve CNN training
- ▶ Rotations / Mirroring / Flipping
- ▶ Changing colours or image intensity
- ▶ Scaling (change size)
- ▶ Elastic (non-linear) deformations
- ▶ For “classical” image processing,
features were designed to be invariant
to scale and rotation (e.g. log polar transform)
- ▶ In deep learning, this is solved with data augmentation
(CNNs are lazy, will only learn what they need to learn)

ROTATION



What happens if you train a network on 1 million cat images,
and then show an upside down cat?

Add randomly rotated versions of all images as training data,
add horizontal and vertical flips of training images

SCALE



What happens if you train a network on 1 million cat images,
and then show a much smaller cat?

Add randomly scaled versions of all images as training data

QUALITY



What happens if you train a network on 1 million cat images,
and then show a blurry cat?

Add blurred versions of all images as training data

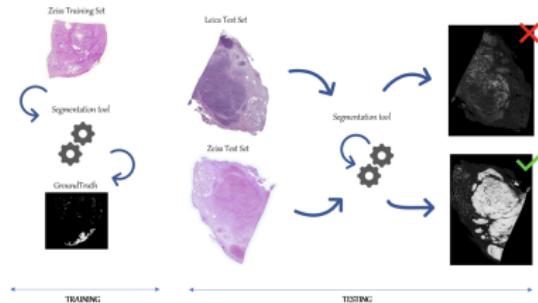
Add noise to all training images

DIFFERENT CAMERAS / SCANNERS

- ▶ Deep networks are often sensitive to the input data (overfitting)
- ▶ A CNN trained on images from camera / scanner A will normally not perform well on images from camera / scanner B
- ▶ Different cameras / scanners have different noise properties
- ▶ Different cameras / scanners have different color histograms
- ▶ For improved robustness, train on data from different cameras / scanners or use 'style transfer' techniques to 'translate' images

DIFFERENT CAMERAS / SCANNERS

Problem:

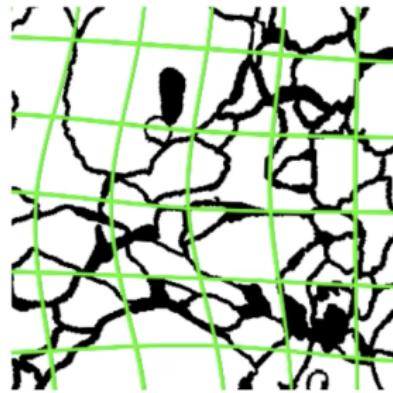
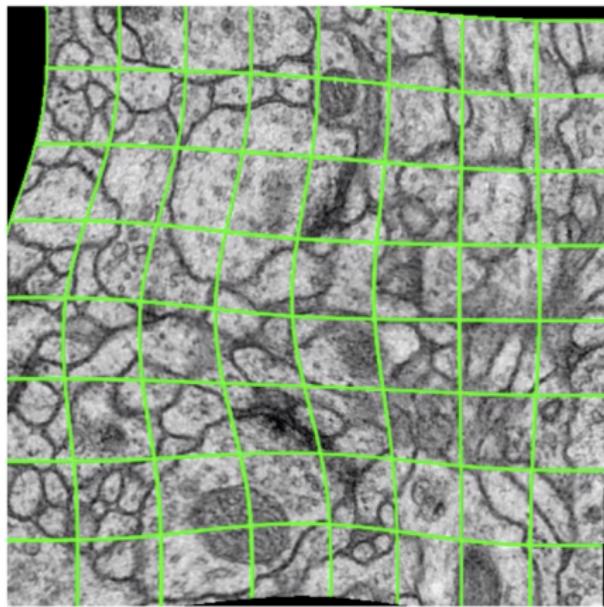


Solution:



Alessia de Biase, Generative Adversarial Networks to enhance decision support in digital pathology, LIU-IDA/STAT-A-19/007-SE

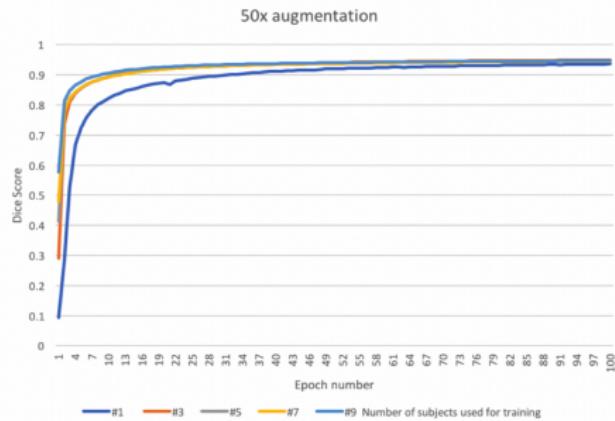
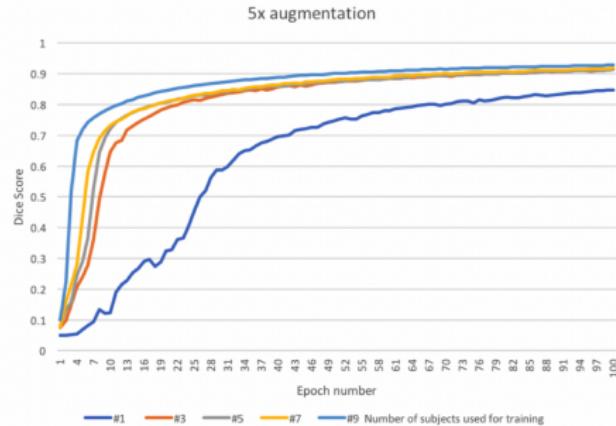
ELASTIC (NON-LINEAR) DEFORMATIONS



correspondingly deformed
manual labels

<https://lmb.informatik.uni-freiburg.de/people/ronneber/u-net/>

DEGREE OF AUGMENTATION



Gaonkar, Bilwaj, et al. Extreme Augmentation: Can deep learning based medical image segmentation be trained using a single manually delineated scan?, arXiv:1810.01621, 2018

GENERATIVE ADVERSARIAL NETWORKS

- ▶ Standard data augmentation cannot create new images, can only modify existing images through rotations, scaling, ...
- ▶ Generative adversarial networks (GANs) can be trained to synthesize new images, given a large set of training images
- ▶ A GAN consists of a generator G (e.g. generates new images), and a discriminator D (e.g. classifies an image as real or fake)
- ▶ For images, both the generator and discriminator are CNNs
- ▶ Goodfellow et al. (2014). Generative adversarial nets. In Advances in neural information processing systems (pp. 2672-2680)

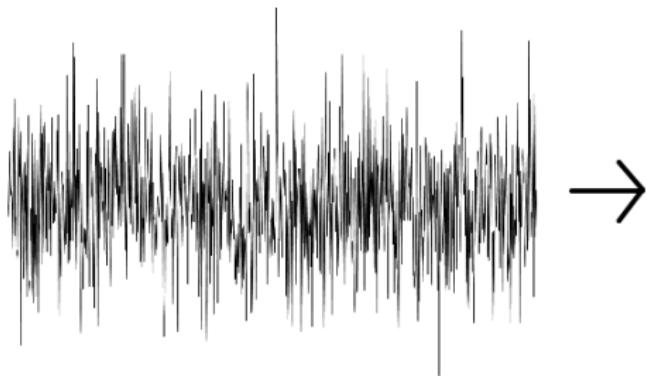
GANS - ANY DATA

- ▶ This presentation will focus on synthesis of images, but GANs can be used to synthesize any kind of data
- ▶ <https://chat.openai.com>, ChatGPT
- ▶ Can be seen as text-to-text GAN / generative AI (haven't looked into details)

GANS - BASICS

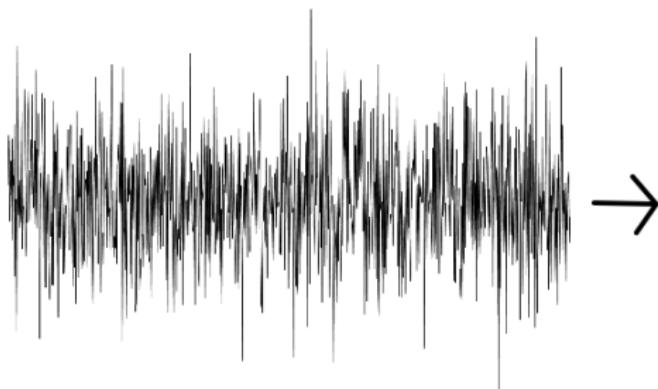
- ▶ Adversarial training, generator and discriminator compete
- ▶ The generator tries to generate better and better images
- ▶ The discriminator tries to be better and better at discriminating real & fake images
- ▶ $G(z, \theta_g)$, generator takes noise vector z as input, uses parameters θ_g to generate an image from the data distribution
- ▶ $D(x, \theta_d)$, discriminator takes sample x , uses parameters θ_d to output a scalar; the probability that the sample x is real and not from the generator
- ▶ Latent (noise) space z , e.g. a noise vector of 128 dimensions, generator maps this noise vector to a manifold of realistic images

NOISE-TO-IMAGE GAN



Karras, T., Aila, T., Laine, S., & Lehtinen, J. Progressive growing of GANs for improved quality, stability, and variation, ICLR, 2018

NOISE-TO-IMAGE GAN



Karras, T., Aila, T., Laine, S., & Lehtinen, J. Progressive growing of GANs for improved quality, stability, and variation, ICLR, 2018

NOISE-TO-IMAGE GAN

- ▶ <https://thispersondoesnotexist.com> (faces)
- ▶ <https://thisxdoesnotexist.com> (everything)

GANS - BASICS

- ▶ Train D to maximize the probability of assigning the correct label to both real (x) and fake (z) samples
- ▶ Train G to minimize $\log(1 - D(G(z)))$, i.e. to fool the discriminator as often as possible, as $D(G(z))$ will be 1 if the discriminator thinks that the fake sample is real
- ▶ Adversarial loss function
 - ▶
$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))].$$

GANS - BASICS

- ▶ In almost all cases it is easier to train the discriminator
- ▶ When starting the training, the fake images have very low quality, very easy for the discriminator to classify images as real or fake, will not provide much information for the generator to learn from
- ▶ Instead of minimizing $\log(1 - D(G(z)))$, maximize $\log(D(G(z)))$, will provide stronger gradients in the beginning of the training

GANS - BASICS

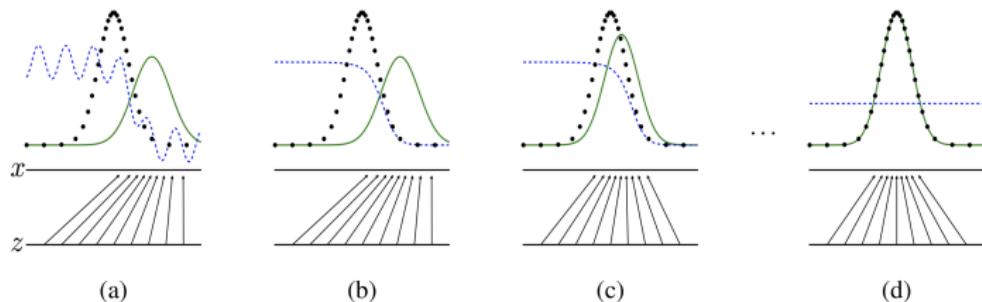


Figure 1: Generative adversarial nets are trained by simultaneously updating the discriminative distribution (D , blue, dashed line) so that it discriminates between samples from the data generating distribution (black, dotted line) p_{data} from those of the generative distribution p_g (G) (green, solid line). The lower horizontal line is the domain from which \mathbf{z} is sampled, in this case uniformly. The horizontal line above is part of the domain of \mathbf{x} . The upward arrows show how the mapping $\mathbf{x} = G(\mathbf{z})$ imposes the non-uniform distribution p_g on transformed samples. G contracts in regions of high density and expands in regions of low density of p_g . (a) Consider an adversarial pair near convergence: p_g is similar to p_{data} and D is a partially accurate classifier. (b) In the inner loop of the algorithm D is trained to discriminate samples from data, converging to $D^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$. (c) After an update to G , gradient of D has guided $G(\mathbf{z})$ to flow to regions that are more likely to be classified as data. (d) After several steps of training, if G and D have enough capacity, they will reach a point at which both cannot improve because $p_g = p_{\text{data}}$. The discriminator is unable to differentiate between the two distributions, i.e. $D(\mathbf{x}) = \frac{1}{2}$.

Goodfellow et al. (2014). Generative adversarial nets. In Advances in neural information processing systems (pp. 2672-2680)

GANS - ADVERSARIAL TRAINING

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 $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ 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

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

Goodfellow et al. (2014). Generative adversarial nets. In Advances in neural information processing systems (pp. 2672-2680)

GANS - FIRST RESULTS

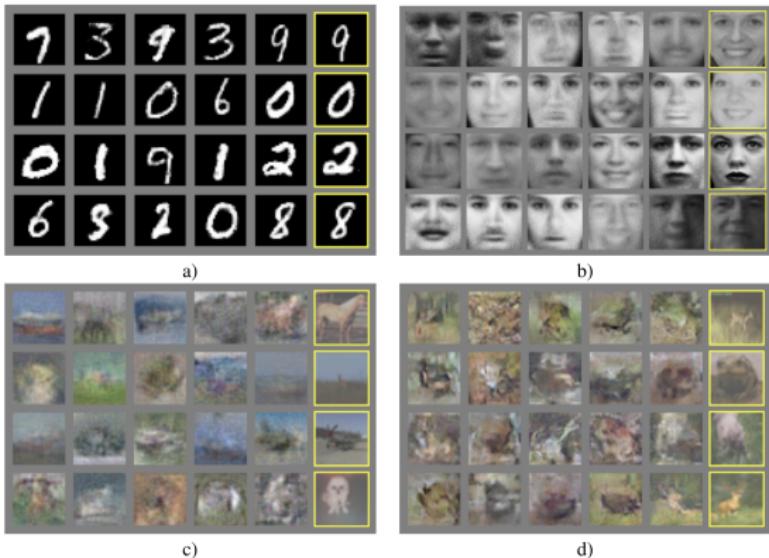


Figure 2: Visualization of samples from the model. Rightmost column shows the nearest training example of the neighboring sample, in order to demonstrate that the model has not memorized the training set. Samples are fair random draws, not cherry-picked. Unlike most other visualizations of deep generative models, these images show actual samples from the model distributions, not conditional means given samples of hidden units. Moreover, these samples are uncorrelated because the sampling process does not depend on Markov chain mixing. a) MNIST b) TFD c) CIFAR-10 (fully connected model) d) CIFAR-10 (convolutional discriminator and “deconvolutional” generator)

Goodfellow et al. (2014). Generative adversarial nets. In Advances in neural information processing systems (pp. 2672-2680)

DEEP CONVOLUTIONAL GAN - DCGAN

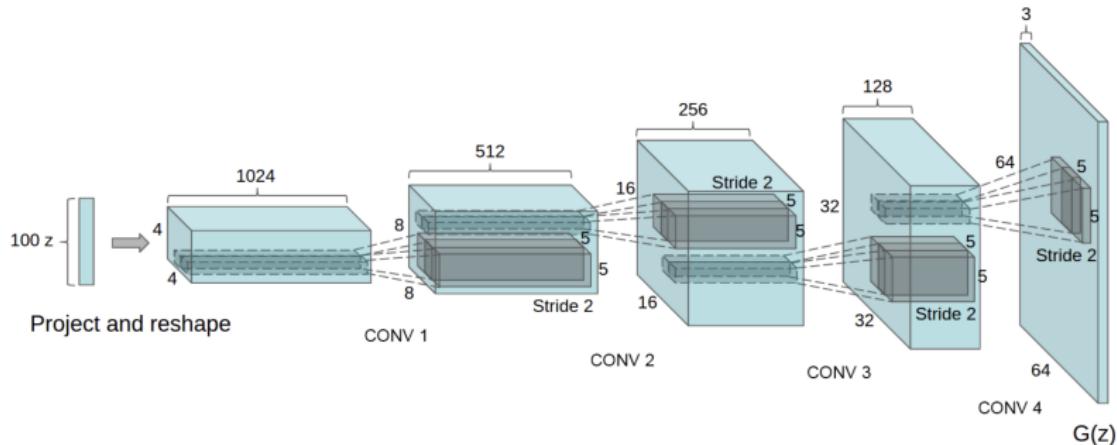


Figure 1: DCGAN generator used for LSUN scene modeling. A 100 dimensional uniform distribution Z is projected to a small spatial extent convolutional representation with many feature maps. A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions) then convert this high level representation into a 64×64 pixel image. Notably, no fully connected or pooling layers are used.

Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv:1511.06434.

DCGAN - FULL ARCHITECTURE

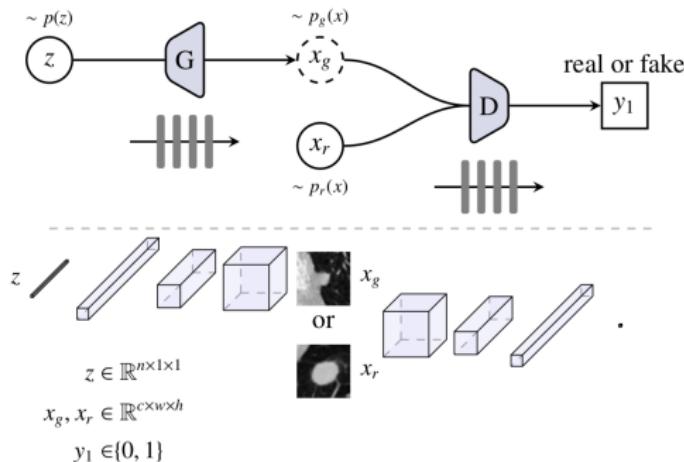


Figure 2: Schematic view of the vanilla GAN for synthesis of lung nodule on CT images. Top of the figure shows the network configuration. The part below shows the input, output and the internal feature representations of the generator G and discriminator D . G transforms a sample z from $p(z)$ into a generated nodule x_g . D is a binary classifier that differentiates the generated and real images of lung nodule formed by x_g and x_r respectively.

Yi, X., Walia, E., & Babyn, P. (2019). Generative adversarial network in medical imaging: A review, Medical Image Analysis

DEEP CONVOLUTIONAL GAN - DCGAN

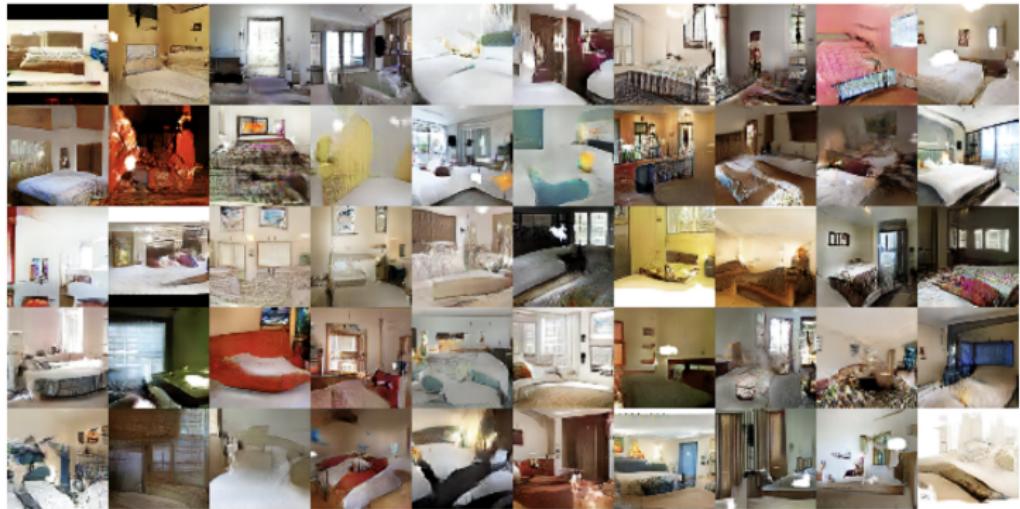


Figure 3: Generated bedrooms after five epochs of training. There appears to be evidence of visual under-fitting via repeated noise textures across multiple samples such as the base boards of some of the beds.

Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv:1511.06434.

DEEP CONVOLUTIONAL GAN - DCGAN

- ▶ DCGAN can rather easily be used to create low resolution images (64 x 64 pixels)
- ▶ Producing images of higher resolution is (much) harder
- ▶ Mode collapse problem: generator will generate a single image
- ▶ Difficult to balance generator and discriminator,
if one is much better than the other the training will stop
- ▶ Wasserstein GAN, another algorithm for more stable training
- ▶ Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein GAN, arXiv:1701.07875.

NOISE-TO-IMAGE GAN - PROGRESSIVE

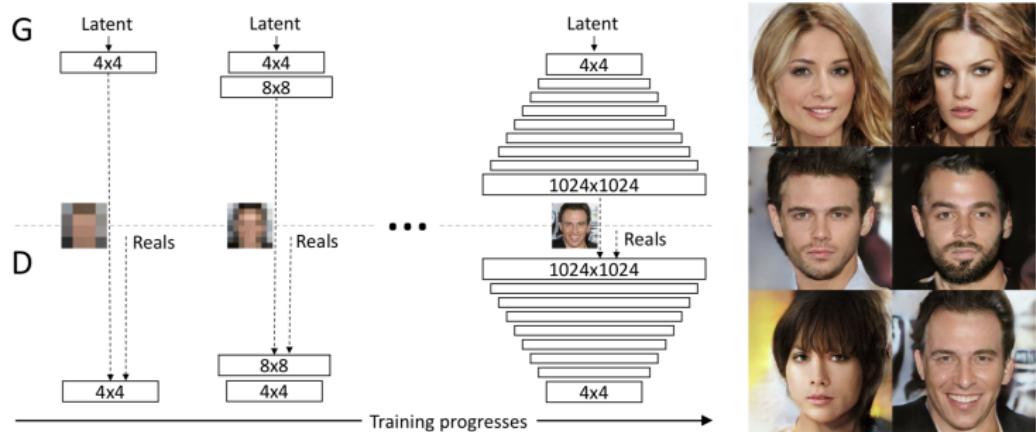
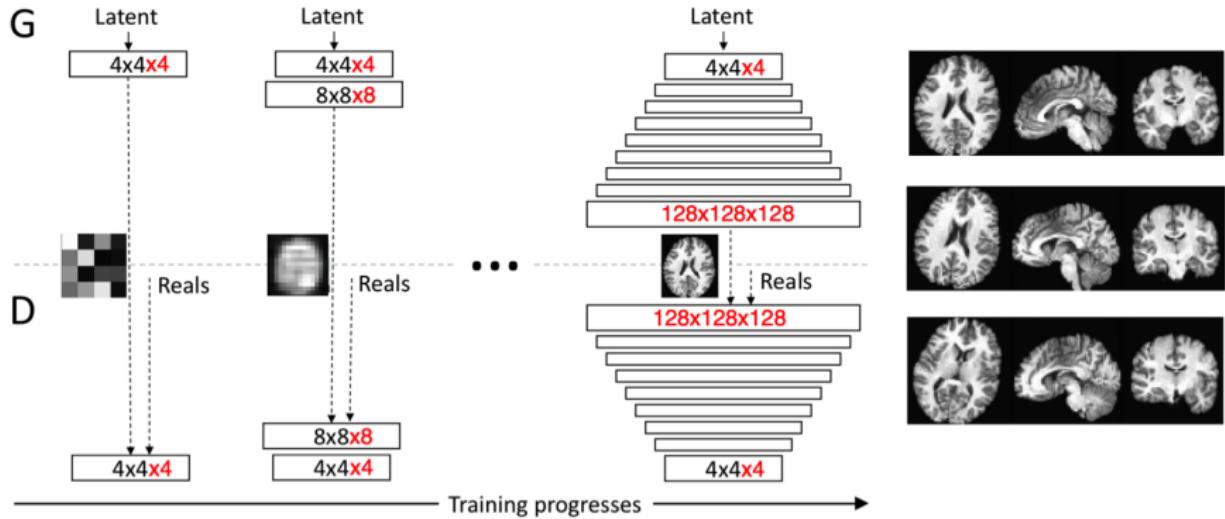


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout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

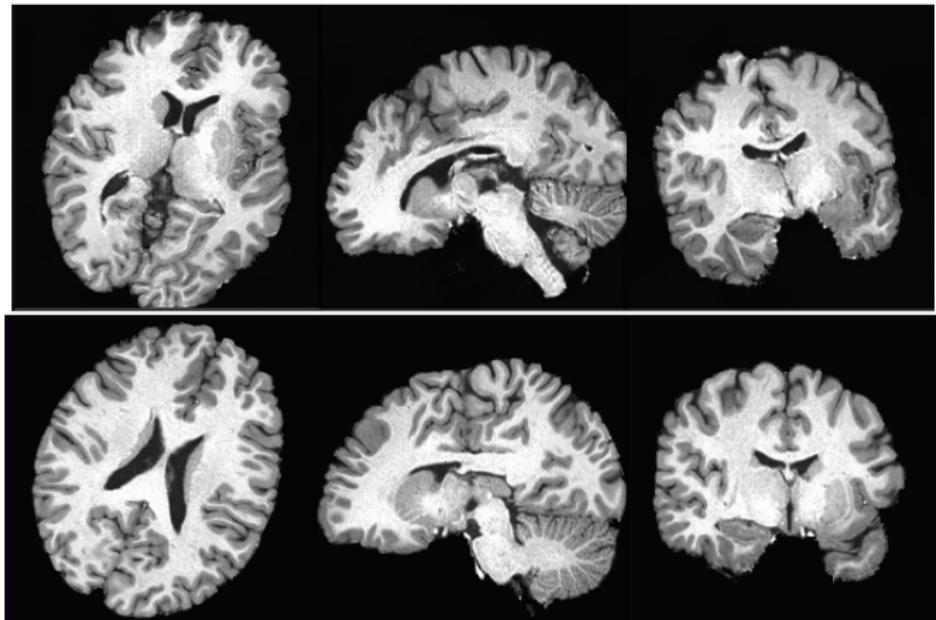
Karras, T., Aila, T., Laine, S., & Lehtinen, J. Progressive growing of GANs for improved quality, stability, and variation, ICLR, 2018

3D PROGRESSIVE NOISE-TO-IMAGE GAN



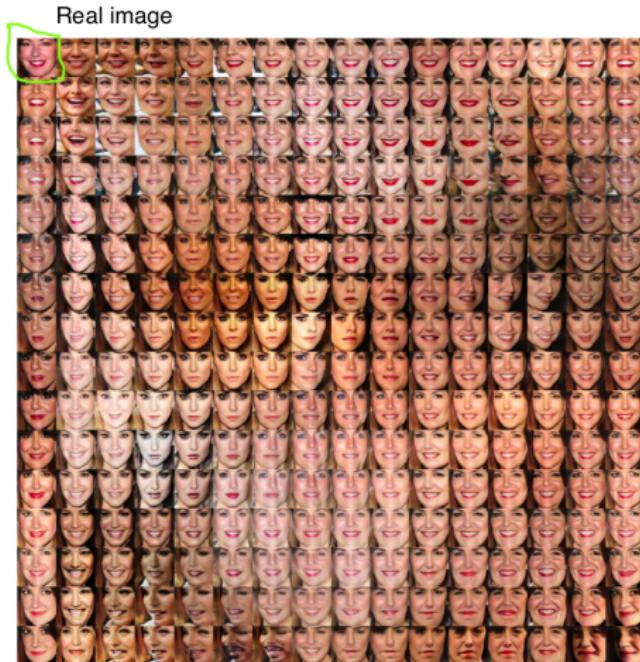
Eklund, A. (2019), Feeding the zombies: Synthesizing brain volumes using a 3D progressive growing GAN. arXiv:1912.05357.

3D PROGRESSIVE NOISE-TO-IMAGE GAN



256 × 256 × 256, Which brain is fake? Code by Wazeer Zulfikar © MIT

GANS FOR DATA AUGMENTATION



Antoniou, A., Storkey, A., Edwards, H., Data Augmentation Generative Adversarial Networks, arXiv:1711.04340

GANS FOR DATA AUGMENTATION

Omniglot DAGAN Augmented Classification		
Experiment ID	Samples Per Class	Test Accuracy
Omni_5_Standard	5	0.689904
Omni_5_DAGAN_Augmented	5	0.821314
Omni_10_Standard	10	0.794071
Omni_10_DAGAN_Augmented	10	0.862179
Omni_15_Standard	15	0.819712
Omni_15_DAGAN_Augmented	15	0.874199

EMNIST DAGAN Augmented Classification		
Experiment ID	Samples Per Class	Test Accuracy
EMNIST_Standard	15	0.739353
EMNIST_DAGAN_Augmented	15	0.760701
EMNIST_Standard	25	0.783539
EMNIST_DAGAN_Augmented	25	0.802598
EMNIST_Standard	50	0.815055
EMNIST_DAGAN_Augmented	50	0.827832
EMNIST_Standard	100	0.837787
EMNIST_DAGAN_Augmented	100	0.848009

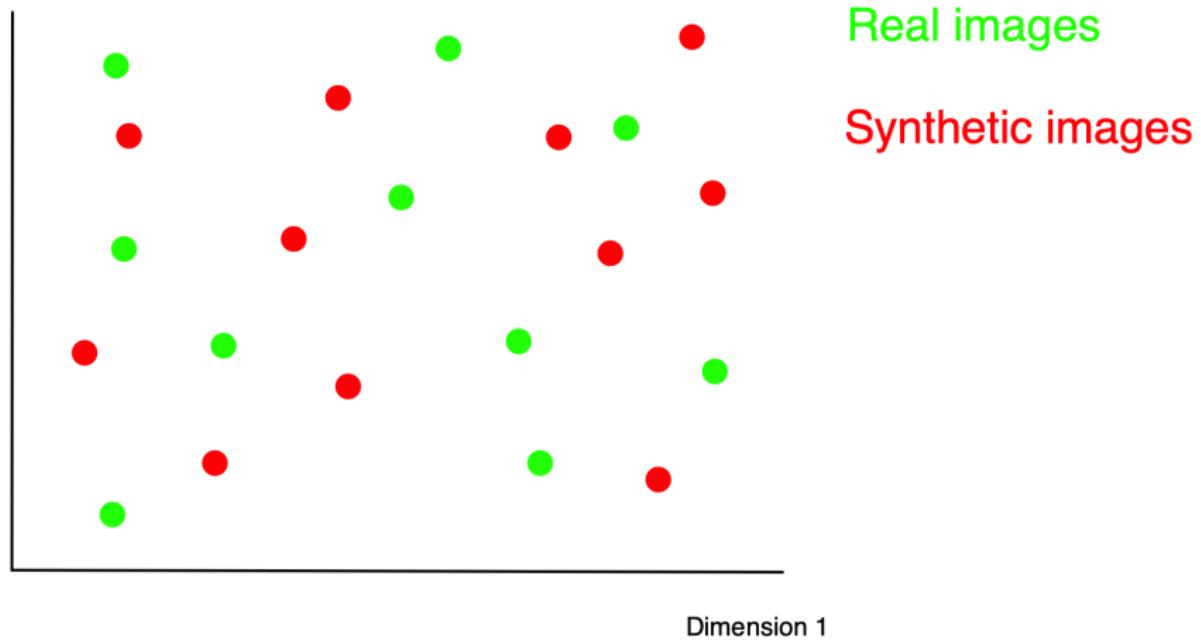
Face DAGAN Augmented Classification		
Experiment ID	Samples Per Class	Test Accuracy
VGG-Face_Standard	5	0.0446948
VGG-Face_DAGAN_Augmented	5	0.125969
VGG-Face_Standard	15	0.39329
VGG-Face_DAGAN_Augmented	15	0.429385
VGG-Face_Standard	25	0.579942
VGG-Face_DAGAN_Augmented	25	0.584666

Table 1: Vanilla Classification Results: All results are averages over 5 independent runs. The DAGAN augmentation improves the classifier performance in all cases. Test accuracy is the result on the test cases in the test domain

Antoniou, A., Storkey, A., Edwards, H., Data Augmentation Generative Adversarial Networks, arXiv:1711.04340

DATA AUGMENTATION - LOW DIM PLOT

Dimension 2



HOW TO EVALUATE GANS

- ▶ No objective loss function is used when training GANs
- ▶ Not so easy to say how good a GAN is,
hard to compare different GAN architectures / configurations
- ▶ Want to compare GANs in terms of image quality, diversity
- ▶ Manual / visual evaluation;
synthesize images and let humans evaluate them
 - ▶ Subjective, can lead to biased results
 - ▶ What is a realistic image in a certain domain?
 - ▶ Time consuming
 - ▶ Low reproducibility

QUALITATIVE GAN GENERATOR EVALUATION

- ▶ Qualitative ways to evaluate a GAN generator
 - ▶ Look at nearest neighbours in training set
(has the GAN memorized the training data?)
 - ▶ Rating and preference judgment
 - ▶ Evaluating mode drop and mode collapse
 - ▶ Investigating and visualizing the internals of network
- ▶ Borji, A. (2019). Pros and cons of GAN evaluation measures.
Computer Vision and Image Understanding, 179, 41-65.

EVALUATING MODE DROP AND MODE COLLAPSE

- ▶ GANs often fail to model the entire data distribution
- ▶ Mode collapse; many input noise vectors are mapped to the same synthetic image (i.e. low diversity of the synthetic images)
- ▶ Mode drop; some modes of the data distribution are ignored
- ▶ May be hard to detect mode collapse for GANs trained on a large number of real images, easier for synthetic datasets
- ▶ A mode is considered lost if there is no sample in the generated test data within a certain standard deviations from the center of that mode

EVALUATING MODE DROP AND MODE COLLAPSE

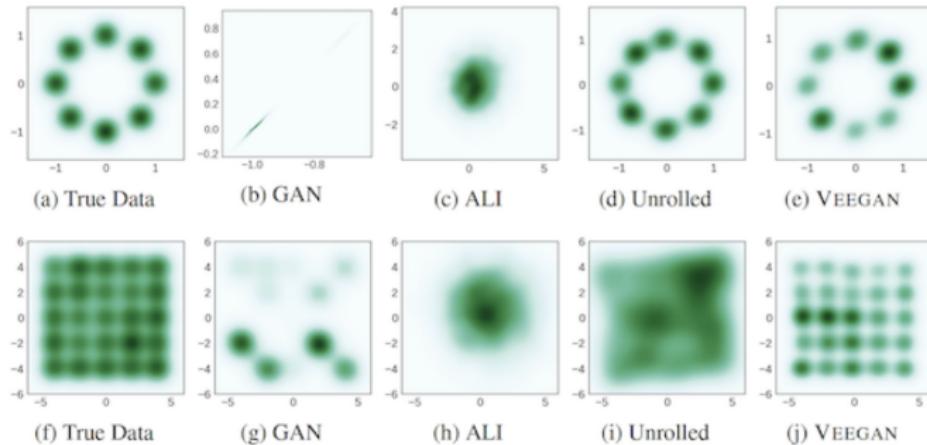


Figure 19: Density plots of the true data and generator distributions from different GAN methods trained on mixtures of Gaussians arranged in a ring (top) or a grid (bottom). Figure from [58].

Srivastava et al., Reducing mode collapse in GANs using implicit variational learning, in: Advances in Neural Information Processing Systems, 2017, pp. 3310–3320

EVALUATING MODE DROP AND MODE COLLAPSE

- ▶ For real datasets, train GAN using a well balanced dataset (equal number of samples from each class)
- ▶ Train a multi-class classifier using the same dataset
- ▶ Generate images from GAN, use classifier to obtain labels
- ▶ Is the distribution of synthetic images also uniform?

EVALUATING MODE DROP AND MODE COLLAPSE

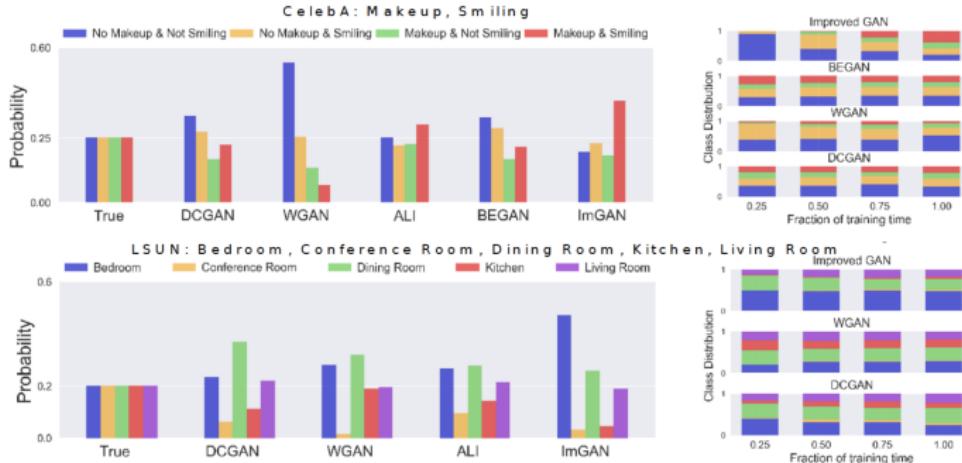


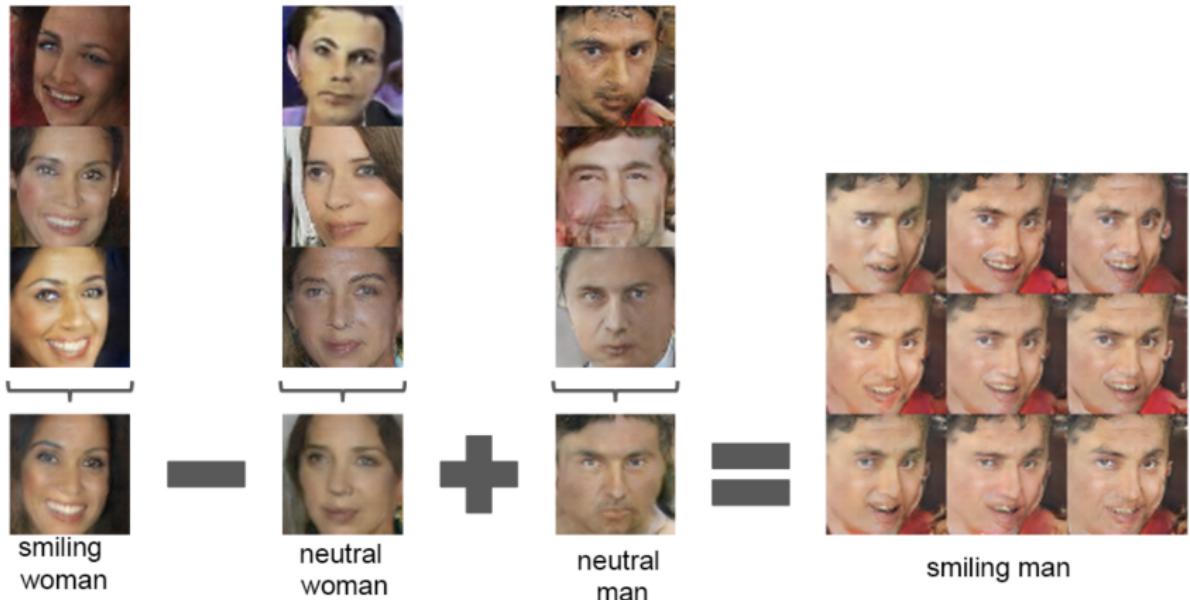
Figure 18: Illustration of mode collapse in GANs trained on select subsets of CelebA and LSUN datasets using the technique in [42]. Left panel shows the relative distribution of modes in samples drawn from the GANs, and compares it to the true data distribution (leftmost plots). Right panel shows the evolution of class distributions in different GANs over the course of training. It can be seen that these GANs introduce covariate shift through mode collapse. Figure compiled from [42].

S. Santurkar, L. Schmidt, A. Madry, A classification-based study of covariate shift in GAN distributions, in: International Conference on Machine Learning, 2018, pp. 4487–4496

INVESTIGATING AND VISUALIZING THE INTERNALS OF NETWORKS

- ▶ Understand the latent space; what does each dimension of the noise vector mean?
- ▶ Disentangled representations. “Disentanglement” regards the alignment of “semantic” visual concepts to axes in the latent space. Some tests can check the existence of semantically meaningful directions in the latent space, meaning that varying the seed along those directions leads to predictable changes (e.g. changes in facial hair, or pose).
- ▶ Perform arithmetic on noise vectors
- ▶ Smiling woman - neutral woman + neutral man = smiling man

INVESTIGATING AND VISUALIZING THE INTERNALS OF NETWORKS



Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv:1511.06434.

INVESTIGATING AND VISUALIZING THE INTERNALS OF NETWORKS

- ▶ Space continuity, given two noise vectors z_1, z_2 interpolate between z_1 and z_2 , check resulting images
- ▶ If the “interpolated” images are reasonable, the model can produce new images, has not only memorized training images



Figure 20: Top: Interpolations on z_r between real images at 128×128 resolution (from BEGAN [124]). These images were not part of the training data. The first and last columns contain the real images to be represented and interpolated. The images immediately next to them are their corresponding approximations while the images in between are the results of linear interpolation in z_r .

- ▶ A. Odena, C. Olah, J. Shlens, Conditional image synthesis with auxiliary classifier GANs, arXiv:1610.09585

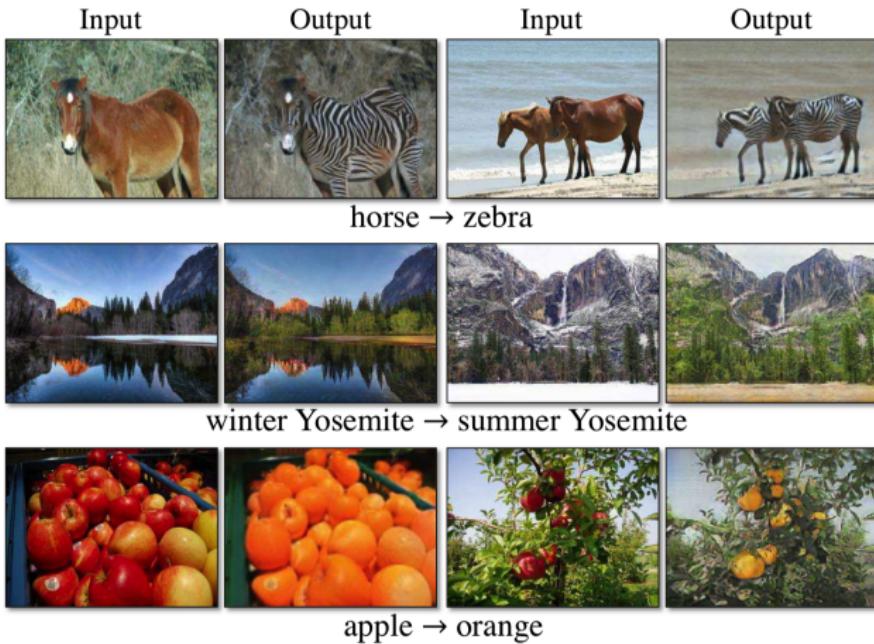
QUANTITATIVE GAN EVALUATION

- ▶ Average log-likelihood, does not work well for high-dimensional signals like images, requires very large number of samples
- ▶ Inception score, uses a pre-trained CNN (Inception Net trained on ImageNet), are images classifiable and diverse with respect to class labels?
- ▶ Fréchet inception distance, embeds a set of generated samples into a feature space given by a specific layer of Inception Net (or any CNN). Viewing the embedding layer as a continuous multivariate Gaussian, the mean and covariance are estimated for both the generated data and the real data. The Fréchet distance between these two Gaussians (a.k.a Wasserstein-2 distance) is then used to quantify the quality of generated samples
- ▶ $FID(r, g) = ||\mu_r - \mu_g||^2 + Tr (\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2})$
- ▶ Borji, A. (2019). Pros and cons of GAN evaluation measures. Computer Vision and Image Understanding, 179, 41-65.

IMAGE-TO-IMAGE GAN (CONDITIONAL GANS)

- ▶ Noise-to-image GAN: synthesize image from noise (latent vector)
- ▶ Image-to-image GAN: synthesize type A image from type B image
- ▶ Image-to-image GANs are normally much easier to train,
since you start from an image and not from noise

IMAGE-TO-IMAGE GAN (CONDITIONAL GANS)



Zhu et al., Unpaired Image-to-Image Translation using
Cycle-Consistent Adversarial Networks, ICCV, 2017

APPLICATIONS - DAY TO NIGHT

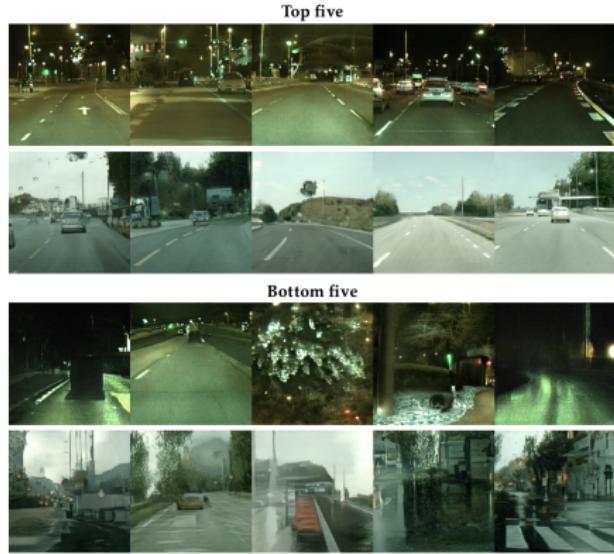


Figure 5.15: Synthetic results from the CycleGAN baseline model on 256x256 pixel images. The model is trained using the large day and night dataset. Top and bottom results are from city and open road street view environments.

S. Karlsson, P. Welander, Generative Adversarial Networks for Image-to-Image Translation on Street View and MR Images,
LIU-IMT-TFK-A—18/554—SE, 2018

APPLICATIONS - MR TO CT

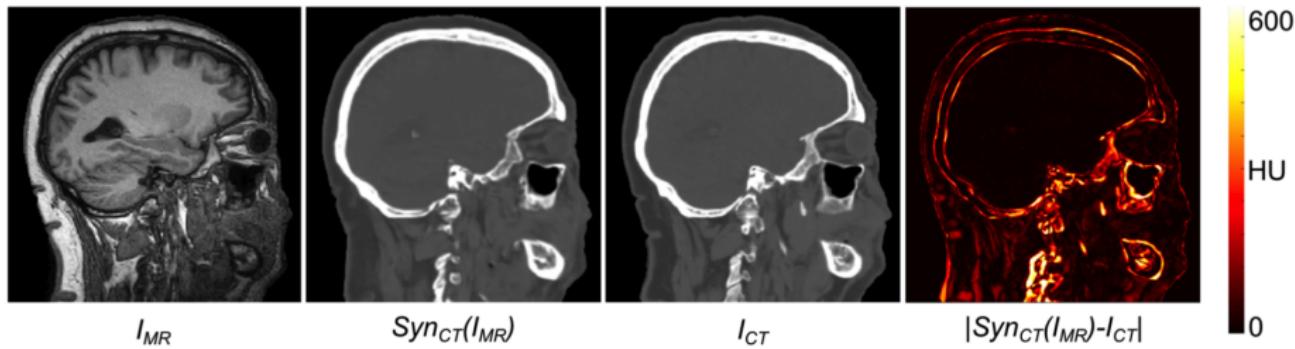


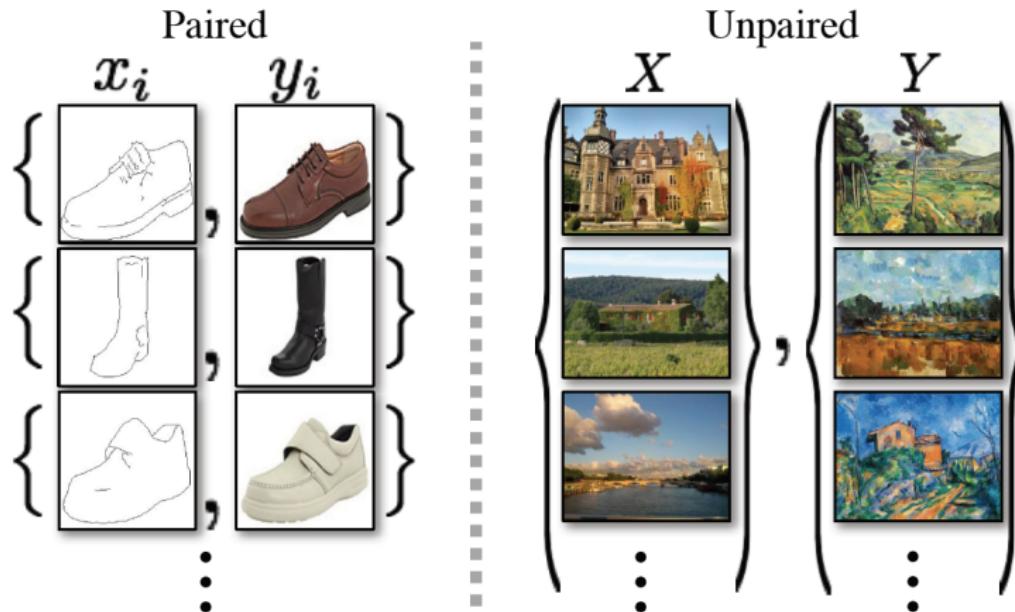
Fig. 4: *From left to right* Input MR image, synthesized CT image, reference real CT image, and absolute error between real and synthesized CT image.

Wolterink et al., Deep MR to CT synthesis using unpaired data.
International Workshop on Simulation and Synthesis in Medical Imaging,
2017

SUPERVISED VS UNSUPERVISED GANS

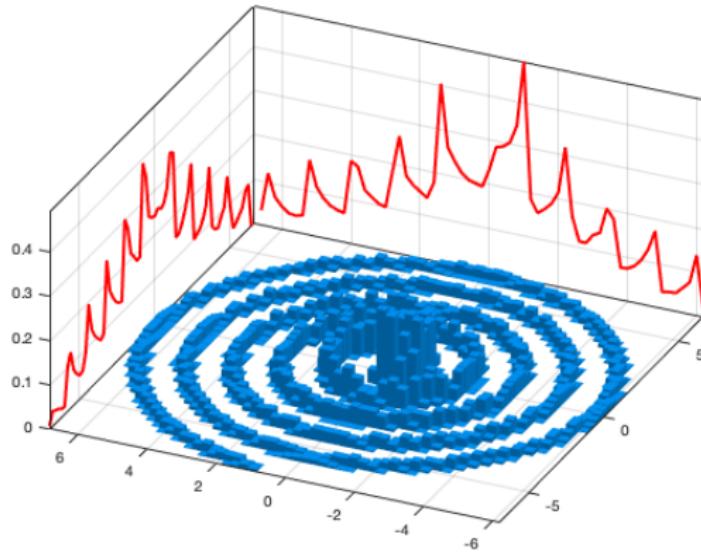
- ▶ We have images from two domains, X & Y
- ▶ We can train a GAN to convert from domain X to domain Y
 - ▶ If we have every image represented in both domains $(x_1; y_1); (x_2; y_2)$; we can train a *supervised* GAN, which learns from the joint distribution of the data in both domains
 - ▶ If we have unmatched images in both domains $(x_1; ?); (?; y_2)$; we can train an *unsupervised* GAN, which learns from the marginal distributions of the data in both domains.

PAIRED AND UNPAIRED DATA



Zhu et al., Unpaired Image-to-Image Translation using
Cycle-Consistent Adversarial Networks, ICCV, 2017

SUPERVISED VS UNSUPERVISED GANS

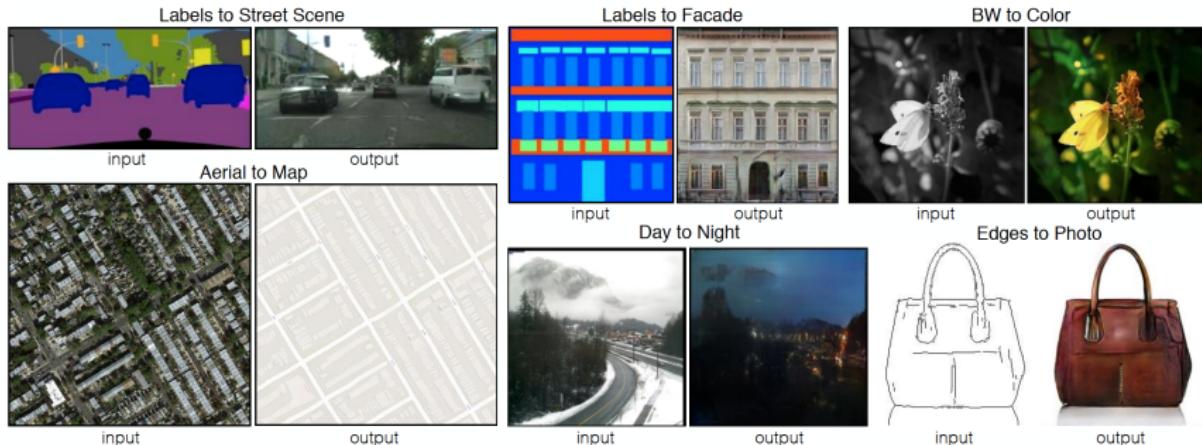


Supervised GAN: Learn the blue distribution
from paired examples from the blue distribution

Unsupervised GAN: Learn the blue distribution
from examples from the red distributions

PAIRED TRAINING - PIX2PIX

- ▶ Given images in two domains, X and Y, the GAN learns to convert an image in domain X to domain Y
- ▶ Drawback: We need paired (and registered) training images



Isola et al., Image-to-Image Translation with
Conditional Adversarial Networks, CVPR, 2017

UNPAIRED TRAINING - CYCLEGAN

- ▶ Unsupervised image-to-image translation is more difficult
- ▶ Learning a joint distribution from marginal distributions is a vastly underdetermined problem, need regularization.
- ▶ One solution: CycleGAN. Learn transformations in both directions $X \leftrightarrow Y$ and impose a consistency constraint in forward-backward conversion
- ▶

$$\begin{cases} x_i \rightarrow \hat{y}_i \rightarrow \tilde{x}_i \approx x_i; \\ y_i \rightarrow \hat{x}_i \rightarrow \tilde{y}_i \approx y_i \end{cases}$$

LOSS FUNCTIONS FOR TRAINING

- ▶ Mappings $G : X \rightarrow Y$ and $F : Y \rightarrow X$,
discriminators D_X and D_Y

- ▶ Adversarial loss function (GAN), L_{GAN}

$$L_{GAN}(G, D_Y, X, Y) = E[\log D_Y(y)] + E[\log(1 - D_Y(G(x)))]$$

- ▶ Cycle consistency loss function (CycleGAN), L_{cyc}

$$L_{cyc}(G, F) = E[||F(G(x)) - x||_1] + E[||G(F(y)) - y||_1]$$

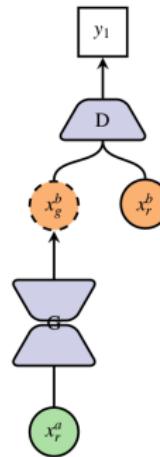
- ▶ Full loss function for CycleGAN

$$L(G, F, D_X, D_Y) =$$

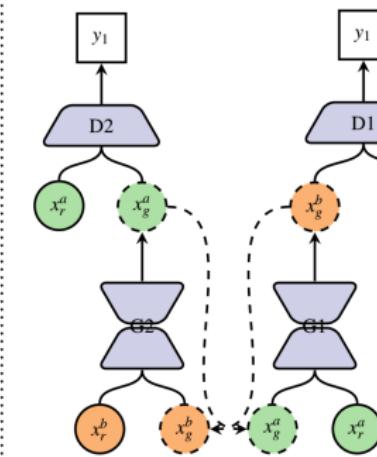
$$L_{GAN}(G, D_Y, X, Y) + L_{GAN}(F, D_X, Y, X) + \lambda L_{cyc}(G, F)$$

PIX2PIX VS CYCLEGAN VS UNIT

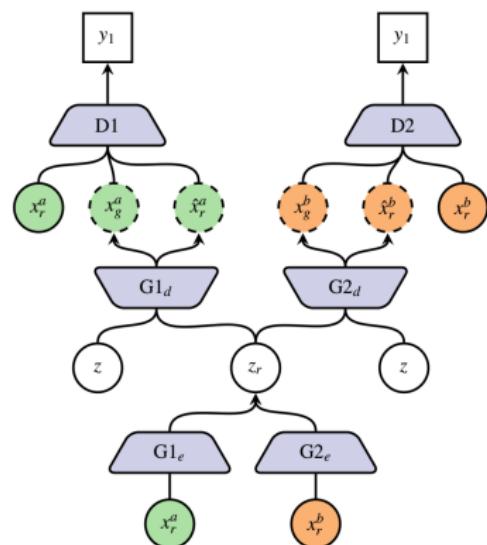
(a) pix2pix



(b) CycleGAN



(c) UNIT



● domain A

● domain B

○ real image

○ generated fake image

Yi, X., Walia, E., & Babyn, P. (2018). Generative adversarial network in medical imaging: A review. Medical Image Analysis

CYCLEGAN - GENERATOR

- ▶ What does the generator look like in Keras?

Generator architectures We adopt our architectures from Johnson et al. [23]. We use 6 residual blocks for 128×128 training images, and 9 residual blocks for 256×256 or higher-resolution training images. Below, we follow the naming convention used in the Johnson et al.'s Github repository.

Let $c7s1-k$ denote a 7×7 Convolution-InstanceNorm-ReLU layer with k filters and stride 1. dk denotes a 3×3 Convolution-InstanceNorm-ReLU layer with k filters and stride 2. Reflection padding was used to reduce artifacts. Rk denotes a residual block that contains two 3×3 convolutional layers with the same number of filters on both layers. uk denotes a 3×3 fractional-strided-Convolution-InstanceNorm-ReLU layer with k filters and stride $\frac{1}{2}$.

The network with 6 residual blocks consists of:

$c7s1-64, d128, d256, R256, R256, R256,$
 $R256, R256, R256, u128, u64, c7s1-3$

The network with 9 residual blocks consists of:

$c7s1-64, d128, d256, R256, R256, R256,$
 $R256, R256, R256, R256, R256, R256, u128$
 $u64, c7s1-3$

```
def modelGenerator(self, name=None):
    # Specify input
    input_img = Input(shape=self.img_shape)
    # Layer 1
    x = ReflectionPadding2D((3, 3))(input_img)
    x = self.c7Ak(x, 32)
    # Layer 2
    x = self.dk(x, 64)
    # Layer 3
    x = self.dk(x, 128)

    # Layer 4-12: Residual layer
    for _ in range(4, 13):
        x = self.Rk(x)

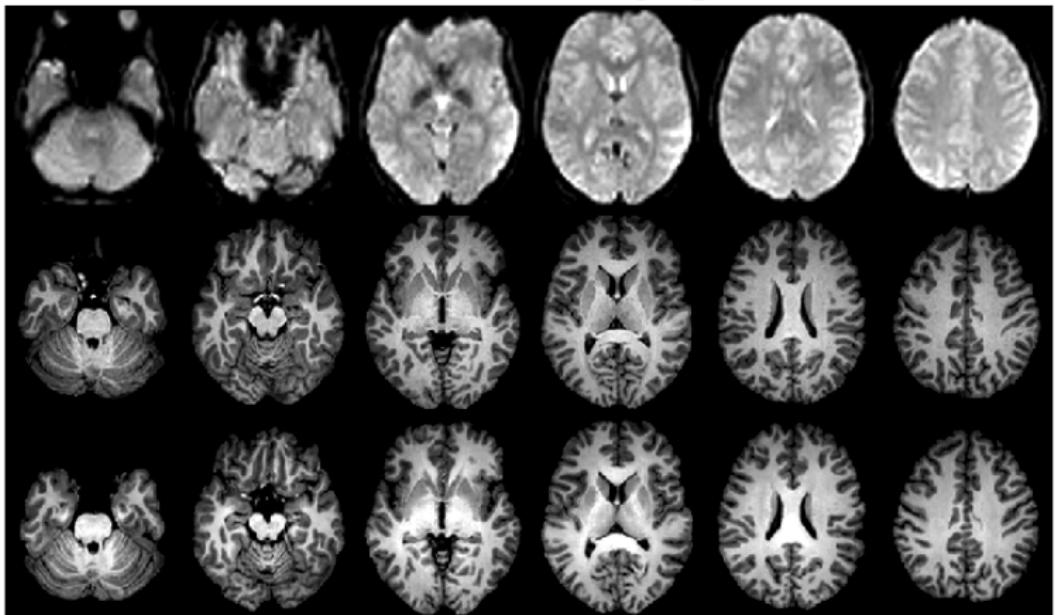
    # Layer 13
    x = self.uk(x, 64)
    # Layer 14
    x = self.uk(x, 32)
    x = ReflectionPadding2D((3, 3))(x)
    x = Conv2D(self.channels, kernel_size=7, strides=1)(x)
    x = Activation('tanh')(x) # They say they use Relu but really they do not
    return Model(inputs=input_img, outputs=x, name=name)
```

- ▶ <https://github.com/simontomaskarlsson/CycleGAN-Keras>

CYCLEGAN - 3D

fMRI to T1, Beijing

fMRI
real



Abramian, D., & Eklund, A. (2019). Generating fMRI volumes from T1-weighted volumes using 3D CycleGAN. arXiv:1907.08533.

TRAINING WITH SYNTHETIC IMAGES

- ▶ Medical images can be difficult to share (ethics, GDPR)
- ▶ Synthetic images do not belong to a specific person (~~GDPR~~) can potentially be shared (legal investigation required)
- ▶ For training segmentation networks, we need images AND the annotations
- ▶ This can be achieved with a multi-channel GAN, but how good are the synthetic images for training CNNs?
- ▶ Some experiments with progressive growing GAN

TRAINING WITH SYNTHETIC IMAGES

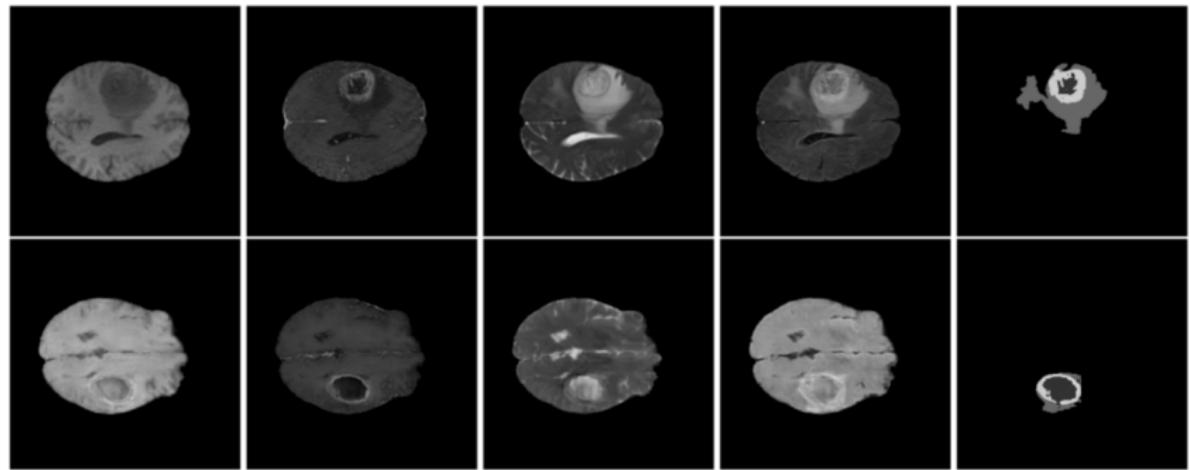


Fig. 1. Top: a real 5-channel image. Bottom: a synthetic 5-channel image. From left to right: T1-weighted, T1-weighted after gadolinium contrast, T2-weighted, FLAIR, segmentation mask.

Larsson, M., Akbar, M. U., & Eklund, A. (2022). Does an ensemble of GANs lead to better performance when training segmentation networks with synthetic images?. arXiv:2211.04086

TRAINING WITH SYNTHETIC IMAGES

Original data	# GANs	ET	ED	NCR/NET	Mean
✓	0	0.791 ± 0.009	0.785 ± 0.003	0.610 ± 0.008	0.729 ± 0.004
✓	1	0.789 ± 0.010	0.790 ± 0.008	0.609 ± 0.013	0.730 ± 0.008
✓	5	0.795 ± 0.012	0.791 ± 0.004	0.617 ± 0.008	0.734 ± 0.006
✓	10	0.791 ± 0.008	0.794 ± 0.005	0.620 ± 0.009	0.735 ± 0.003
✓	20	0.799 ± 0.008	0.791 ± 0.005	0.613 ± 0.011	0.734 ± 0.006
	1	0.665 ± 0.011	0.632 ± 0.027	0.478 ± 0.019	0.592 ± 0.016
	5	0.689 ± 0.012	0.665 ± 0.023	0.529 ± 0.015	0.628 ± 0.014
	10	0.696 ± 0.012	0.702 ± 0.013	0.545 ± 0.012	0.648 ± 0.008
	20	0.696 ± 0.017	0.693 ± 0.021	0.555 ± 0.013	0.648 ± 0.014

Larsson, M., Akbar, M. U., & Eklund, A. (2022). Does an ensemble of GANs lead to better performance when training segmentation networks with synthetic images?. arXiv:2211.04086.

TRAINING WITH SYNTHETIC IMAGES

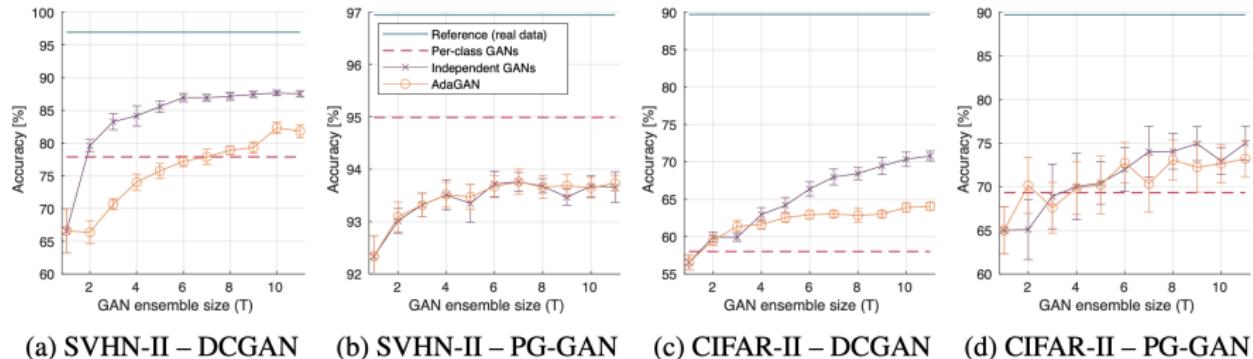


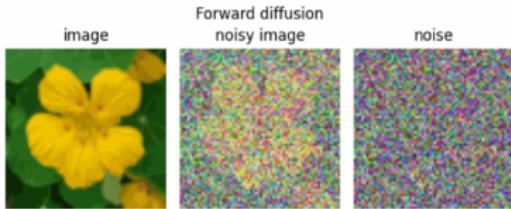
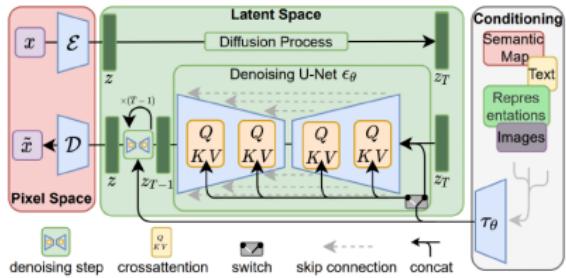
Figure 2: Classification performance on synthetic datasets, comparing different ensemble sizes and approaches. Each datapoint has been estimated from the mean of 10 separate trainings, and standard deviations are reported with error bars.

Eilertsen, G., Tsirikoglou, A., Lundström, C., & Unger, J. (2021). Ensembles of GANs for synthetic training data generation. ICLR 2021

DIFFUSION MODELS

- ▶ The GAN architecture does not work so well when combining text and image data
- ▶ Diffusion models are based on image denoising from image processing (one denoising algorithm is called anisotropic diffusion)
- ▶ Forward process, add more and more (Gaussian) noise
Backward denoising process, perform denoising
- ▶ Need many diffusion steps (iterations), 1000 - 10000, as each denoising will only give a slightly better image

DIFFUSION MODELS



One-sentence summary: diffusion models are trained to denoise noisy images, and can generate images by iteratively denoising pure noise.

<https://keras.io/examples/generative/ddim/>

Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 10684-10695).

TEXT TO IMAGE

huggingface.co/spaces/stabilityai/stable-diffusion

Hugging Face Search models, datasets, users...

Spaces: stabilityai/stable-diffusion 3.49k Running on CUSTOM ENV

App Files and versions Community 5612

Stable Diffusion Demo

Stable Diffusion is a state of the art text-to-image model that generates images from text.
For faster generation and API access you can try [DreamStudio Beta](#)

Police interrogating a jar of pickles Generate image



DIFFUSION MODELS VS GANS

- ▶ Diffusion models are easier to train, due to more standard loss functions
- ▶ GANs are much faster at generating images, diffusion models need many diffusion steps
- ▶ Diffusion models are more likely to memorize training images, the generator in a GAN never sees the training images directly

MEMORIZING TRAINING IMAGES

Extracting Training Data from Diffusion Models

*Nicholas Carlini^{*1} Jamie Hayes^{*2} Milad Nasr^{*1}
 Matthew Jagielski⁺¹ Vikash Sehwag⁺⁴ Florian Tramèr⁺³
 Borja Balle^{†2} Daphne Ippolito^{†1} Eric Wallace^{†5}*
¹Google ²DeepMind ³ETHZ ⁴Princeton ⁵UC Berkeley
^{*}Equal contribution [†]Equal contribution [‡]Equal contribution

Abstract

Image diffusion models such as DALL-E 2, Imagen, and Stable Diffusion have attracted significant attention due to their ability to generate high-quality synthetic images. In this work, we show that diffusion models memorize individual images from their training data and emit them at generation time. With a generate-and-filter pipeline, we extract over a thousand training examples from state-of-the-art models, ranging from photographs of individual people to trademarked company logos. We also train hundreds of diffusion models in various settings to analyze how different modeling and data decisions affect privacy. Overall, our results show that diffusion models are much less private than prior generative models such as GANs, and that mitigating these vulnerabilities may require new advances in privacy-preserving training.



Figure 1: Diffusion models memorize individual training examples and generate them at test time. **Left:** an image from Stable Diffusion's training set (licensed CC BY-SA 3.0, see [49]). **Right:** a Stable Diffusion generation when prompted with "Ann Graham Lotz". The reconstruction is nearly identical (ℓ_2 distance = 0.031).

Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwag, V., Tramèr, F., ... & Wallace, E. (2023). Extracting Training Data from Diffusion Models. arXiv preprint arXiv:2301.13188.

ETHICAL QUESTIONS

- ▶ Synthetic images raise new ethical questions
- ▶ Who is the owner of synthetic data?
- ▶ Can AI be an author / painter?
- ▶ If a generative AI model is trained on medical images, can the synthetic images be seen as anonymized? GDPR?
- ▶ Researchers should not fabricate data...
- ▶ Will medical doctors accept to work on synthetic images?
- ▶ There are clinically approved products, generating CT from MR (to save time)

WHO OWNS SYNTHETIC IMAGES?



Is artificial intelligence set to become art's next medium?

12 December 2018
PHOTOGRAPHS & PRINTS |
AUCTION PREVIEW

Main image:
Portrait of Edmond Belamy
(details shown) by GAN
(Generative Adversarial Network), which will be offered at Christie's on 25–26 October. Image © Christie's

Highlighted sale



Prints & Multiples

AI artwork sells for \$432,500 — nearly 45 times its high estimate — as Christie's becomes the first auction house to offer a work of art created by an algorithm

The portrait in its gilt frame depicts a portly gentleman, possibly French and — to judge by his dark frockcoat and plain white collar — a man of the church. The work appears unfinished: the facial features are somewhat indistinct and there are blank areas of canvas. Oddly, the whole composition is displaced slightly to the north-west. A label on the wall states that the sitter is a man named Edmond Belamy, but the giveaway clue as to the origins of the work is the artist's signature at the bottom right. In cursive Gallic script it reads:

$$\min_G \max_D \mathbb{E}_x[\log(D(x))] + \mathbb{E}_z[\log(1 - D(G(z)))]$$

Christie's/PA Wire/PA

<https://www.christies.com/features/A-collaboration-between-two-artists-one-human-one-a-machine-9332-1.aspx>