

# CS-E4891 Deep Generative Models

## Lecture 6: Generative Adversarial Networks

Lauri Juvela

Department of Information and Communications Engineering (DICE)  
Aalto University

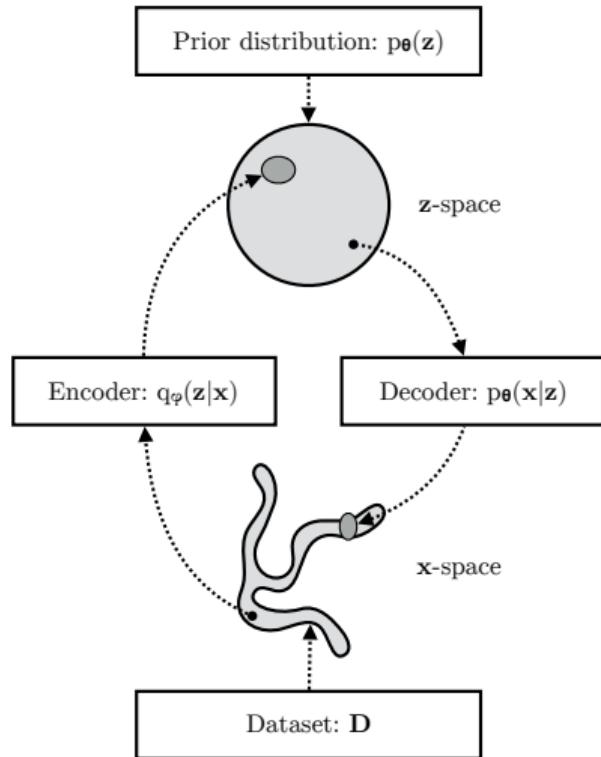
April 5, 2025

## Outline

- Generative Adversarial Networks (GANs)
- Learning by comparison
- GAN loss functions
- GAN training process
- Theoretical properties: density ratios, quality, coverage
- Conditional GAN
- GAN architectures
- Cycle GAN
- Reading: Section 22 from Murphy, 2023; Chapter 15 from Prince 2024; Chapter 17 from Bishop, 2024)

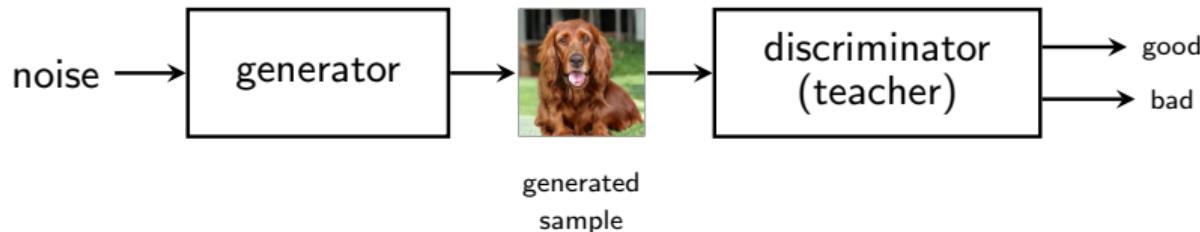
## Variational autoencoders

- VAE reconstruction loss is defined in the pixel space, this is not a good metric perceptually
- Classifier neural networks are powerful for learning representations
- Could we use another neural net as a score function?



# Generative Adversarial Networks (GANs) – Introduction

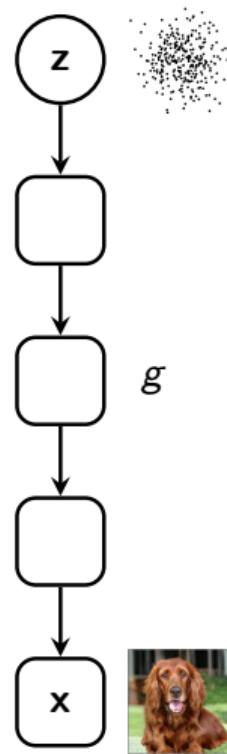
- GANs consist of two neural networks
- **Generator** transforms samples from a simple prior distribution to the data space
- **Discriminator** compares the generated samples to real data examples and scores them



- The Generator can be any parametric differentiable model (typically a deep neural network)
- Discriminator is a (deep neural network) classifier that classifies between real and generated samples

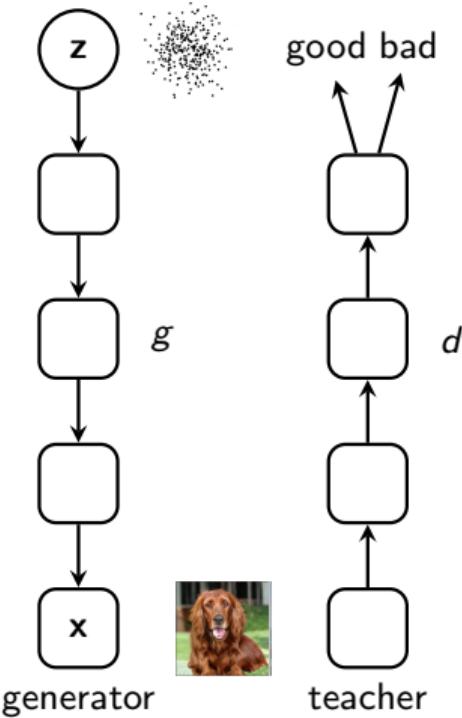
## Generator

- Sample latent variable from isotropic Gaussian  
 $z \sim \mathcal{N}(0, I)$
- Transform  $z$  to the data space using a deep neural network  
 $x = G(z; \theta)$
- Generator is an implicit generative model – model outputs generated samples but the explicit probability distribution is unknown



## Generator and Discriminator

- Generator is guided by a teacher network (discriminator or critic) that evaluates the quality of generated samples
- How to train the teacher network  $d(x; \phi)$ ?
- Classification task: classify training data examples as “good” and generator output samples as “bad”
- Generator attempts to improve its output by making its output more likely to classified as “good”



## Learning by comparison

- GANs are implicit generative models
- Learning by comparison: sample from a generative models and compare generated samples to real data examples
- For more see, Murphy section 26.2

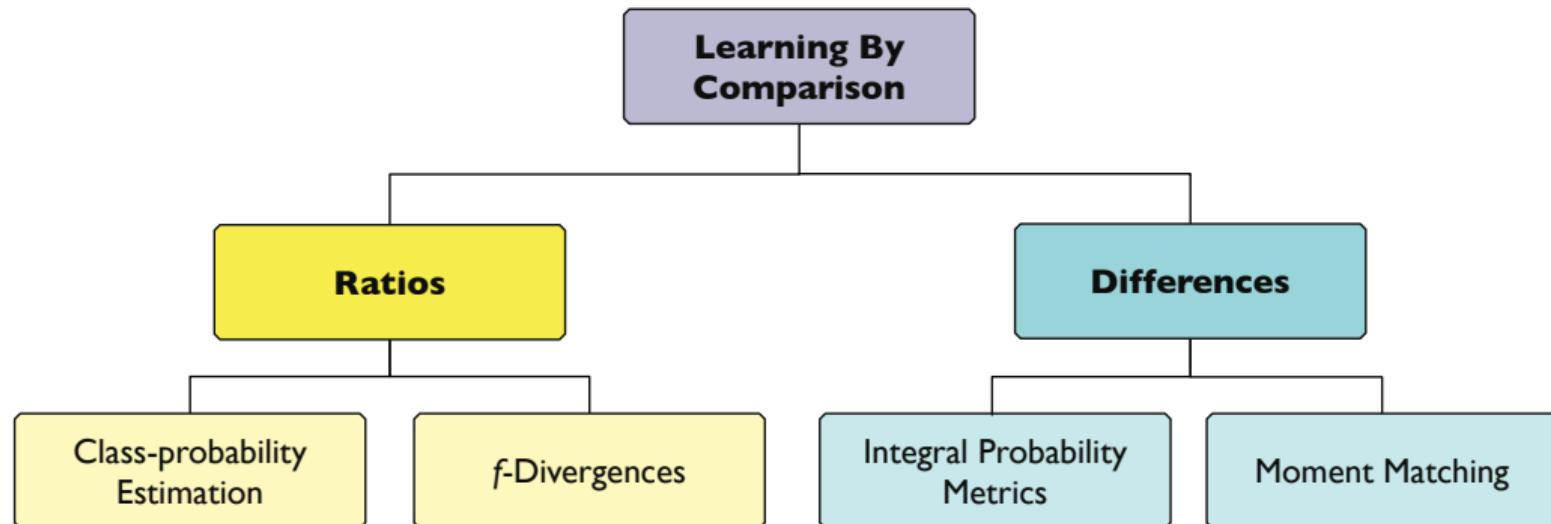


Figure: Overview of approaches for learning in implicit generative models (Murphy, Fig. 26.2)

## GAN loss function

- Discriminator  $D_\phi(\mathbf{x}) \in [0, 1]$  is a neural network parameterised by  $\phi$
- Select the Bernoulli log-loss, also known as binary cross entropy loss, to obtain the “vanilla” GAN training objective for Discriminator

$$V(q_\theta) = \arg \max_{\phi} \mathbb{E}_{p(\mathbf{x}|y)p(y)} [y \log D_\phi(\mathbf{x}) + (1 - y) \log(1 - D_\phi(\mathbf{x}))]$$

- Split the objective into two separate terms for real and fake classes

$$V(q_\theta) = \arg \max_{\phi} [\mathbb{E}_{p(\mathbf{x}|y=1)p(y=1)} \log D_\phi(\mathbf{x}) + \mathbb{E}_{p(\mathbf{x}|y=0)p(y=0)} \log(1 - D_\phi(\mathbf{x}))]$$

- Draw samples from real data distribution  $p_{\text{data}}(\mathbf{x}) = p^*(\mathbf{x})$  and generated data distribution  $p_G(\mathbf{x}) = q_\theta(\mathbf{x})$

$$V(q_\theta) = \arg \max_{\phi} [\mathbb{E}_{p^*(\mathbf{x})} \log D_\phi(\mathbf{x}) + \mathbb{E}_{q_\theta(\mathbf{x})} \log(1 - D_\phi(\mathbf{x}))]$$

- More in Murphy, section 26.2

## Non-saturating GAN loss

- Change the generator objective: instead of minimizing the probability of generator output classified as fake  $\log(1 - D(G(z)))$ , maximize the probability of generator output classified as real –  $\log D(G(z))$
- The non-saturating loss provides stronger gradients when the discriminator is easily detecting that generated samples are fake.
- Generator loss as a function of discriminator score. (left)
- The gradients of the generator loss with respect to the discriminator score. (right)

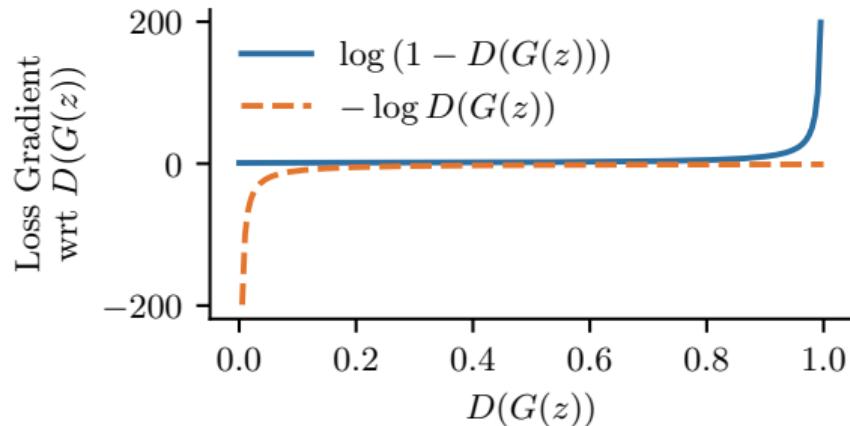
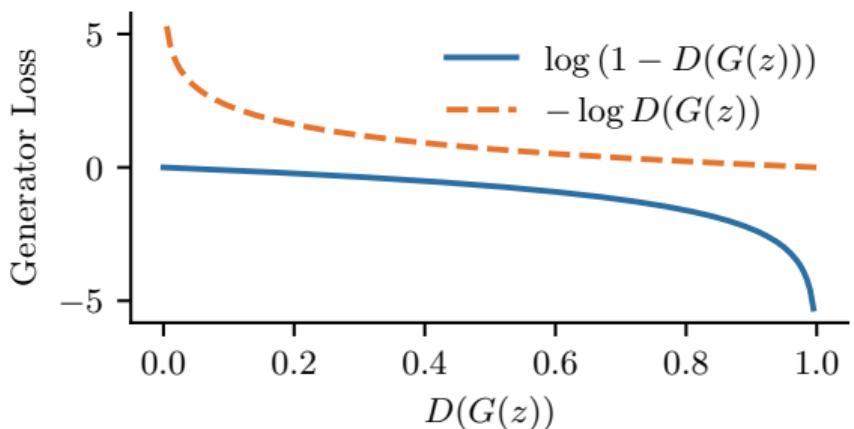


Figure: Murphy Fig. 26.5

## GAN training procedure

### 1. Update the Discriminator

- Sample N examples  $\mathbf{x}_i$  from the training set
- Generate N samples  $G(\mathbf{z}_i)$
- Compute the binary cross-entropy loss

$$\mathcal{L}_D = -\frac{1}{N} \sum_{i=1}^N \log D(\mathbf{x}_i) - \frac{1}{N} \sum_{i=1}^N \log (1 - D(G(\mathbf{z}_i)))$$

- Update  $\phi$  by SGD (or Adam):  $\phi \leftarrow \phi - \nabla_\phi \mathcal{L}_D$

### 2. Update the Generator

- Generate N samples  $G(\mathbf{z}_i)$  using the Generator
- Compute the (non-saturating) loss function

$$\mathcal{L}_G = -\frac{1}{N} \sum_{i=1}^N \log D(G(\mathbf{z}_i))$$

- Update  $\theta$  by SGD (or Adam):  $\theta \leftarrow \theta - \nabla_\theta \mathcal{L}_G$  (gradients flow through the discriminator)

## Optimal discriminator (home exercise)

Write the expected values in the objective function in the integral form

$$V(G, D) = \int_x p^*(x) \log D(x) dx + \int_z p(z) \log(1 - D(G(z))) dz$$

Treat  $G$  as constant and apply “the law of unconscious statistician” to combine the expectations to a single integrand

$$= \int_x p^*(\mathbf{x}) \log D(\mathbf{x}) + q_{\theta}(\mathbf{x}) \log(1 - D(\mathbf{x})) dx$$

Upper bound: the integral is less than or equal to an integral over the maximum of an integrand

$$V(G, D) \leq \int_x \max [p^*(\mathbf{x}) \log D(\mathbf{x}) + q_{\theta}(\mathbf{x}) \log(1 - D(\mathbf{x}))] dx$$

Check the extreme values of the integrand, equality is fulfilled at

$$D^*(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_{\theta}(\mathbf{x})}$$

This is the optimal discriminator

## GANs implicitly minimise the Jensen-Shannon divergence

- For the optimal discriminator the GAN objective becomes

$$\begin{aligned} V^*(q_{\theta}, p^*) &= \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} \left[ \log \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_{\theta}(\mathbf{x})} \right] + \frac{1}{2} \mathbb{E}_{q_{\theta}(\mathbf{x})} \left[ \log \left( 1 - \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_{\theta}(\mathbf{x})} \right) \right] \\ &= \frac{1}{2} \mathbb{E}_{p^*(\mathbf{x})} \left[ \log \frac{p^*(\mathbf{x})}{\frac{p^*(\mathbf{x}) + q_{\theta}(\mathbf{x})}{2}} \right] + \frac{1}{2} \mathbb{E}_{q_{\theta}(\mathbf{x})} \left[ \log \frac{q_{\theta}(\mathbf{x})}{\frac{p^*(\mathbf{x}) + q_{\theta}(\mathbf{x})}{2}} \right] - \log 2 \\ &= \frac{1}{2} D_{\text{KL}} \left( p^* \middle\| \frac{p^* + q_{\theta}}{2} \right) + \frac{1}{2} D_{\text{KL}} \left( q_{\theta} \middle\| \frac{p^* + q_{\theta}}{2} \right) - \log 2 \\ &= \text{JSD}(p^*, q_{\theta}) - \log 2 \end{aligned}$$

- In practice, we can only approximate the optimal discriminator

## Quality and coverage

- **Coverage:** divergence is small where true density is high and the mixture  $p^*(\mathbf{x}) + q_\theta(\mathbf{x})$  has high probability.
- Penalize regions with real examples but no generated samples

$$D_{\text{KL}} \left( p^* \left\| \frac{p^* + q_\theta}{2} \right. \right) = \frac{1}{2} \int p^*(\mathbf{x}) \log \frac{2p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})} d\mathbf{x}$$

- **Quality:** divergence is small where generated sample density  $q_\theta(\mathbf{x})$  is high and the mixture  $p^*(\mathbf{x}) + q_\theta(\mathbf{x})$  has high probability;
- Penalize regions with generated samples but no real samples

$$D_{\text{KL}} \left( q_\theta(\mathbf{x}) \left\| \frac{p^* + q_\theta}{2} \right. \right) = \frac{1}{2} \int q_\theta(\mathbf{x}) \log \frac{2q_\theta(\mathbf{x})}{p^*(\mathbf{x}) + q_\theta(\mathbf{x})} d\mathbf{x}$$

- Balancing quality and diversity (coverage) are fundamental problems in generative modeling
- Mode collapse is a typical problem with GANs

## Problems with density ratios

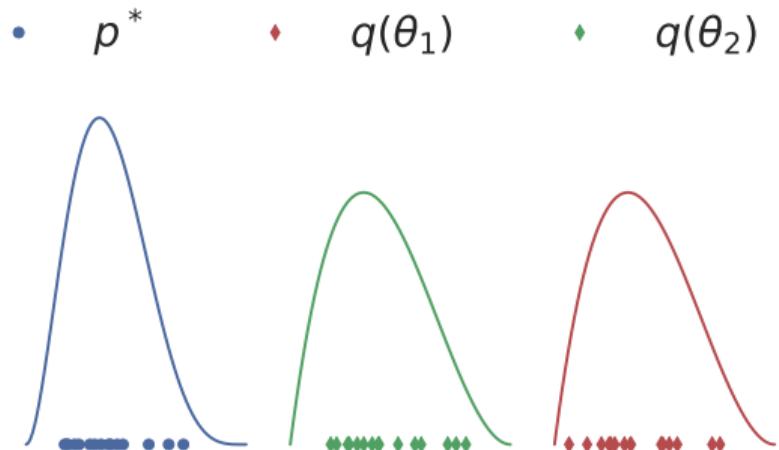


Figure: Murphy, Fig. 26.4a

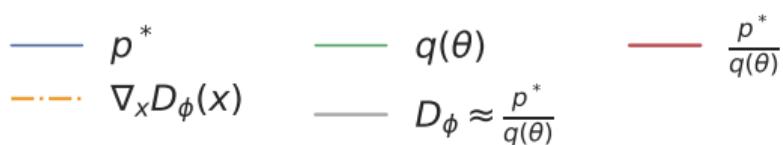


Figure: Murphy, Fig. 26.4b

- Failure of the KL divergence to distinguish between distributions with non-overlapping support
- $D_{\text{KL}}(p^* \| q_{\theta_1}) = D_{\text{KL}}(p^* \| q_{\theta_2}) = \inf$
- $q_{\theta_1}$  is clearly closer to  $p^*$  than  $q_{\theta_2}$

- The density ratio  $\frac{p^*}{q_\theta}$  used by the KL divergence and a smooth estimate given by an MLP, together with the gradient it provides with respect to the input variable.

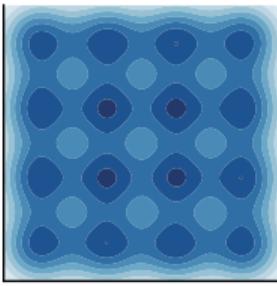
## GAN min-max game

- Global minimum of  $\mathcal{L}_G$  is achieved if and only if  $q_{\theta} = p^*$
- This is the Nash-equilibrium of GAN min-max game – no player can gain anything by changing just their own strategy
- Optimal discriminator is then

$$D^*(\mathbf{x}) = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + q_{\theta}(\mathbf{x})} = \frac{p^*(\mathbf{x})}{p^*(\mathbf{x}) + p^*(\mathbf{x})} = \frac{1}{2}$$

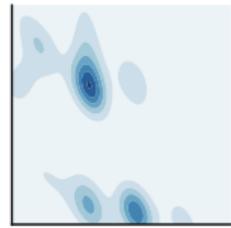
- This means that Discriminator can not distinguish generated samples from real data examples

## GAN mode collapse

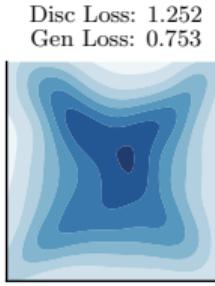


(a) Data

Disc Loss: 0.892  
Gen Loss: 1.852

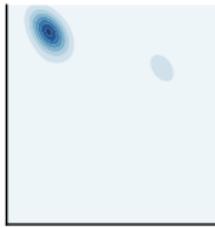


(d) 8000 iterations

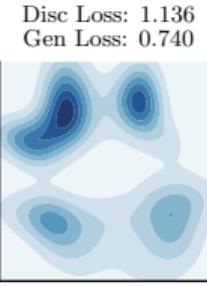


(b) 2000 iterations

Disc Loss: 1.252  
Gen Loss: 0.753

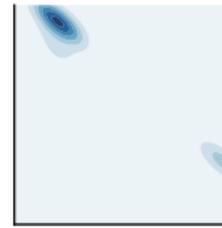


(e) 12000 iterations



(c) 4000 iterations

Disc Loss: 1.136  
Gen Loss: 0.740



(f) 14000 iterations

Disc Loss: 0.305  
Gen Loss: 3.016

Figure: Mode collapse with GAN trying to fit a dataset of 16 Gaussians in 2D (Murphy, Fig. 26.6)

## Least squares GAN and non-zero-sum GANs

- General GAN loss functions can have independent terms for Discriminator and Generator

$$\min_{\phi} \mathcal{L}_D(\phi, \theta); \min_{\theta} \mathcal{L}_G(\phi, \theta)$$

- Zero-sum formulations  $-L_D(\phi, \theta) = L_G(\phi, \theta)$  are not strictly needed for GANs to work
- See Murphy section 26.3.2 for a wide range of various GAN loss functions and how they connect to f-divergences

### Least-squares GAN

- Discriminator tries to output 1 (real data label) for the real examples and 0 (fake data label) for generated samples

$$\mathcal{L}_D = \mathbb{E}_{p^*(\mathbf{x})}(D(\mathbf{x}) - 1)^2 + \mathbb{E}_{q_\theta(\mathbf{x})} D(\mathbf{x})^2$$

- Generator tries to generate samples that get a discriminator output close to 1 (real data label)

$$\mathcal{L}_G = \mathbb{E}_{q_\theta(\mathbf{x})}(D(\mathbf{x}) - 1)^2$$

## Discriminator Feature Matching and Moment Matching

- Often the Discriminator becomes too powerful and it's helpful to give the generator a sneak-peek into Discriminator activations
- Compute discriminator hidden activations  $D(\mathbf{x})$  for real and generated samples and try to match those
- Moment matching tries to match the statistics of discriminator activations using a moment generating function  $s(\cdot)$ , for example, the mean and variance

$$\mathcal{L}_{\text{MM}} = \min_{\theta} \|\mathbb{E}_{p^*(\mathbf{x})} s(D(\mathbf{x})) - \mathbb{E}_{q_{\theta}(\mathbf{x})} s(D(\mathbf{x}))\|_2^2$$

- Feature matching refers to trying to directly match the hidden activations without first computing statistics
  - this assumes paired data

$$\mathcal{L}_{\text{FM}} = \min_{\theta} \mathbb{E}_{p^*(\mathbf{x}), q_{\theta}(\mathbf{x})} \|D(\mathbf{x}) - D(G(\mathbf{z}))\|_2^2$$

- Feature matching is heavily used in speech synthesis, especially HiFi-GAN and related models
- Murphy 26.2.5

# Deep Convolutional GAN (DCGAN)

- In the generator, a 100D latent variable  $z$  is mapped by a linear transformation to a  $4 \times 4$  representation with 1024 channels.
- This is then passed through convolutional layers that gradually upsample the representation and decrease the number of channels.
- At the end a tanh function maps the representation to a fixed range so that it can represent an image.
- The discriminator consists of a standard convolutional net that classifies the input as either a real example or a generated sample.

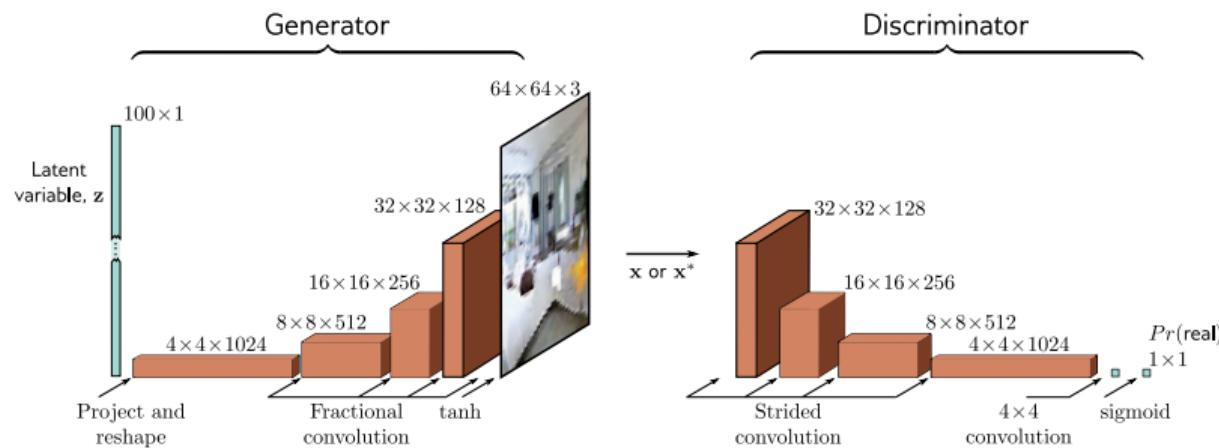
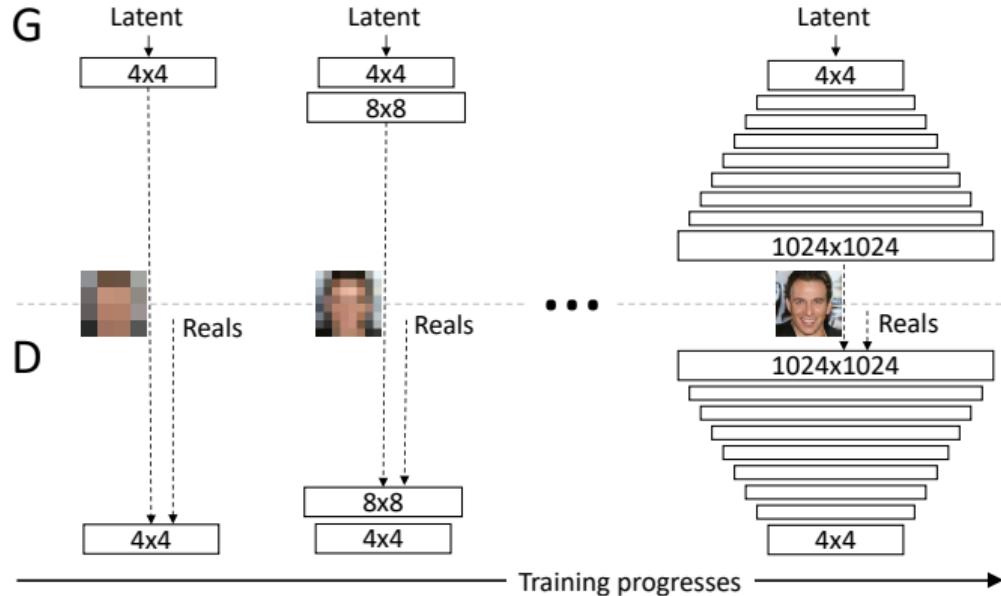


Figure: DCGAN architecture. Fig. 15.3 from Prince

# Progressive growing of GAN

- Start by generating and discriminating images at low resolutions, progressively add more upsampling and downsampling layers
- Training is more stable for small images
- Training time is reduced compared to iterating at full resolution from the beginning



## Images generated by progressive growing of GANs



Figure: From Karras et al., 2018

## Conditional GAN

- The generator of the conditional GAN also receives an attribute vector  $c$  describing some aspect of the image.
- As usual, the discriminator receives either a real example or a generated sample, but now it also receives the attribute vector; this encourages the samples both to be realistic and compatible with the attribute.

a) Conditional GAN

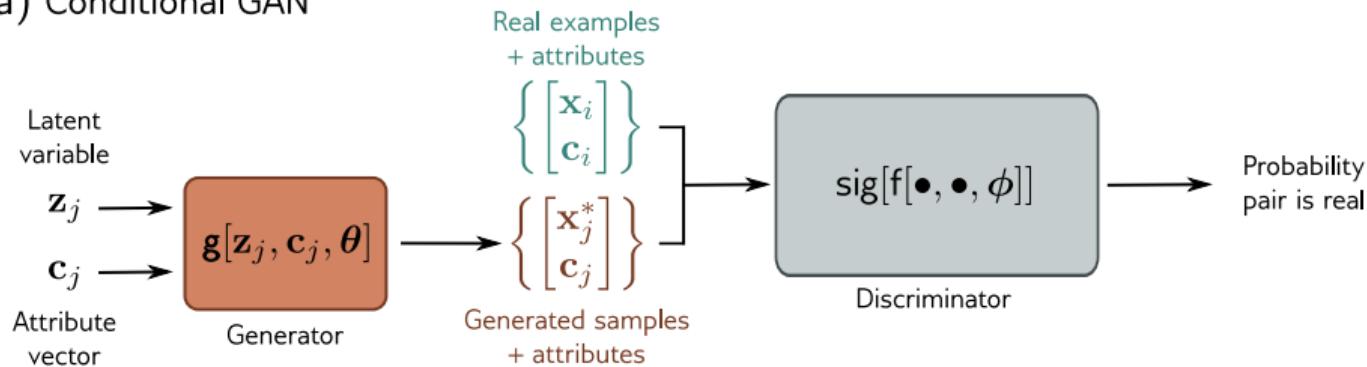


Figure: Fig. 15.13 from Prince, 2024

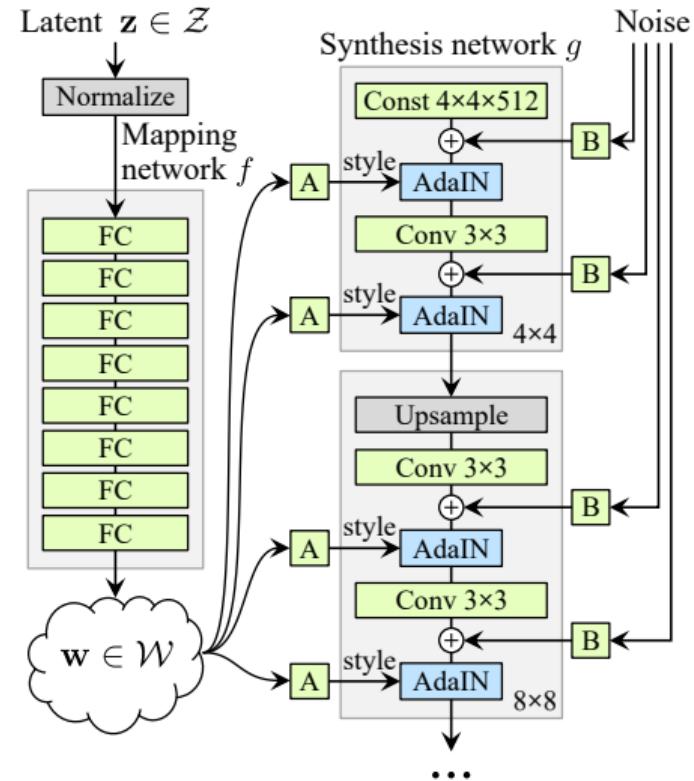
# StyleGAN

- Mapping network learns different conditioning features (styles) for different spatial resolutions in the network
- Additional noise inputs are provided for texture
- Conditioning combines FiLM and batch norm into Adaptive Instance Normalisation (AdaIN)

$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \odot \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i}$$

where  $\mathbf{x}_i$  is one feature map,  $\mu(\mathbf{x}_i)$  and  $\sigma(\mathbf{x}_i)$  are minibatch mean and standard deviation, respectively

- Style vectors  $(\mathbf{y}_s, \mathbf{y}_b)$  are multiplicative and additive conditioning (similar to FiLM) derived from the latent  $\mathbf{z}$



## Cycle GAN – unpaired image-to-image translation

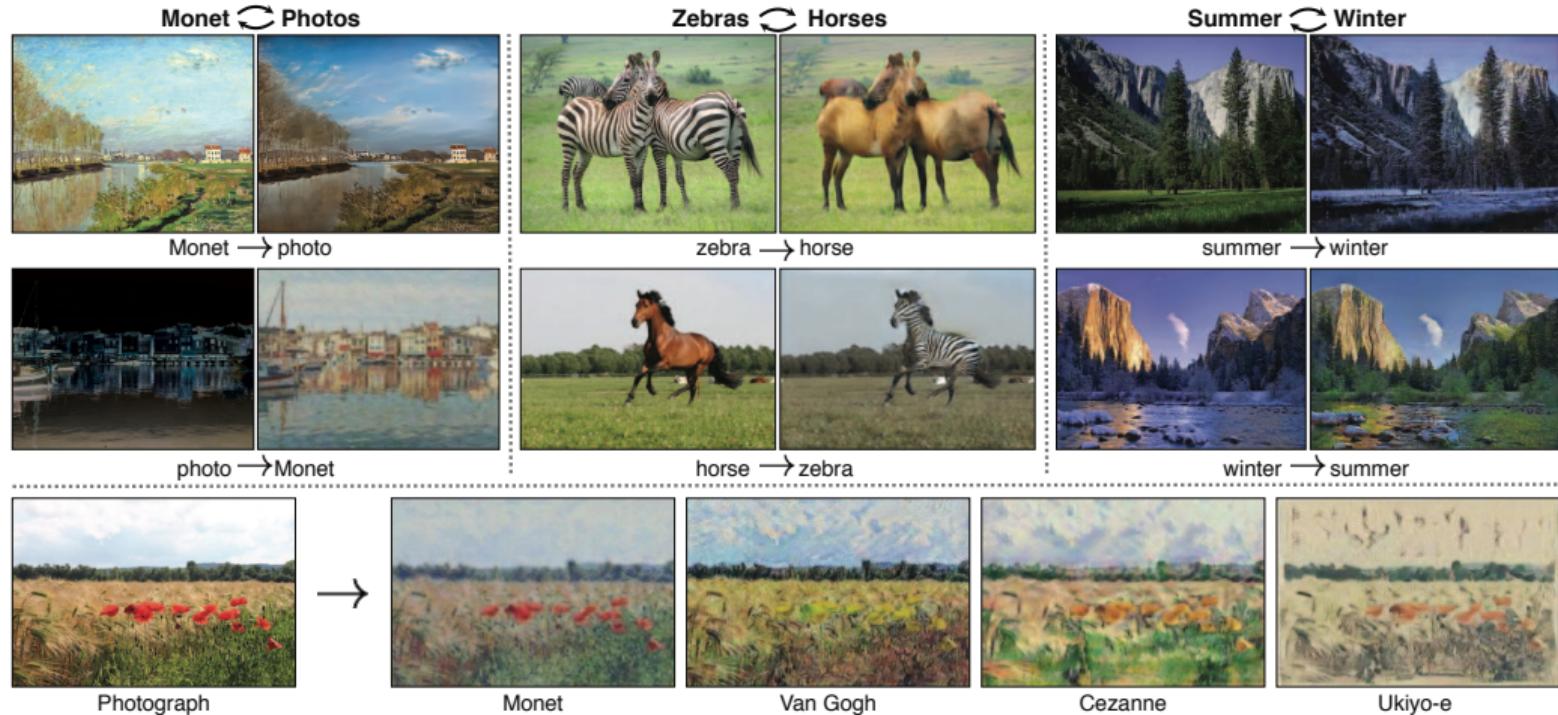


Figure: Examples of unpaired image-to-image translation. (Murphy Fig. 26.18, orig. Zhu et al., 2017)

# Cycle GAN

- Cycle consistency loss: round-trip converted samples should be similar to the original sample  
$$\mathcal{L}_{\text{cycle}} = \mathbb{E}_{p(\mathbf{x})} \|F(G(\mathbf{x})) - \mathbf{x}\|_1 + \mathbb{E}_{p(\mathbf{y})} \|G(F(\mathbf{y})) - \mathbf{y}\|_1$$
- Identity loss: converting images back to its own domain should not change the image  
$$\mathcal{L}_{\text{identity}} = \mathbb{E}_{p(\mathbf{x})} \|\mathbf{x} - F(\mathbf{x})\|_1 + \mathbb{E}_{p(\mathbf{y})} \|\mathbf{y} - G(\mathbf{y})\|_1$$
- Total loss

$$\mathcal{L} = \mathcal{L}_{GAN} + \lambda_1 \mathcal{L}_{\text{cycle}} + \lambda_2 \mathcal{L}_{\text{identity}}$$

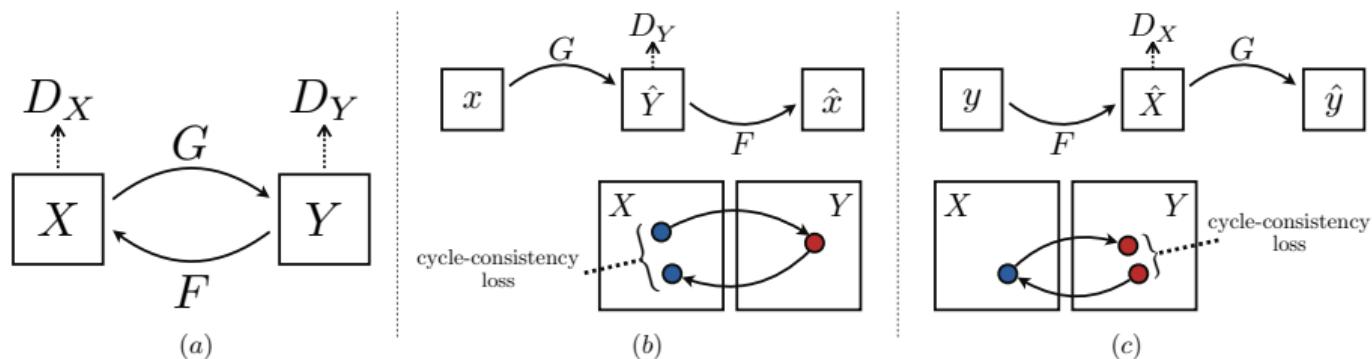


Figure: CycleGAN training scheme. (Murphy Fig. 26.18, orig. Zhu et al. 2017)

## Evaluation of generated samples

- (a) Generated samples are of high quality, but not diverse
- (b) Generated samples are both high quality and diverse
  - How to capture diversity in an evaluation metric?

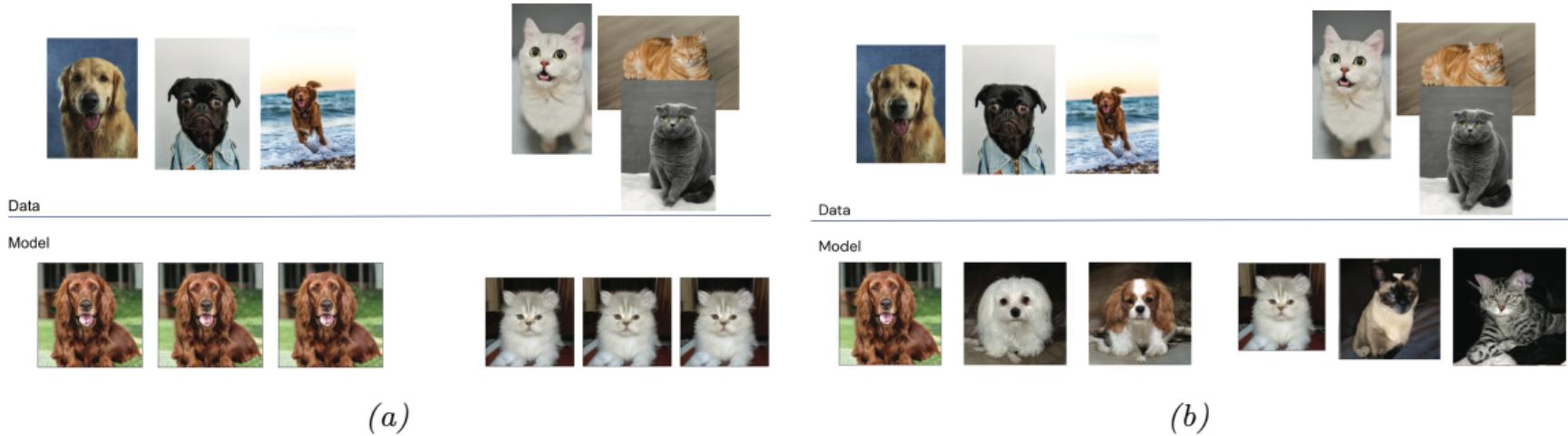


Figure: Murphy Fig. 20.11

## Fréchet distances

- Fréchet Inception Distance (FID) is calculated from statistics of hidden activations in a pre-trained image classified (Inception Net in this case)
- Calculate mean  $\mu_m$  and covariance  $\Sigma_m$  from generated model output samples
- Calculate mean  $\mu_d$  and covariance  $\Sigma_d$  from real data examples

$$\text{FID} = \|\mu_m - \mu_d\|_2^2 - \text{tr} \left( \Sigma_d + \Sigma_m - 2(\Sigma_d \Sigma_m)^{1/2} \right)$$

- The idea of Fréchet distances is applicable to audio (FAD) and other domains – choose the suitable pretrained classifier for your application
- Similar to moment matching and discriminator feature matching, but typically only used for evaluation
- More in Murphy, Section 20.4.2

## Summary

- Generative Adversarial Networks (GANs)
- Learning by comparison
- GAN loss functions
- GAN training process
- Theoretical properties: density ratios, quality, coverage
- Conditional GAN
- GAN architectures
- Cycle GAN
- Reading: Section 22 from Murphy, 2023; Chapter 15 from Prince 2024; Chapter 17 from Bishop, 2024)

## References

- Murphy K, Probabilistic Machine Learning: Advanced Topics, 2023.
- Prince SJD, Understanding Deep Learning, The MIT Press, 2024.
- Bishop, Christopher M., and Hugh Bishop. Deep learning: Foundations and concepts. Springer Nature, 2024.