

Deep Learning (CS324)

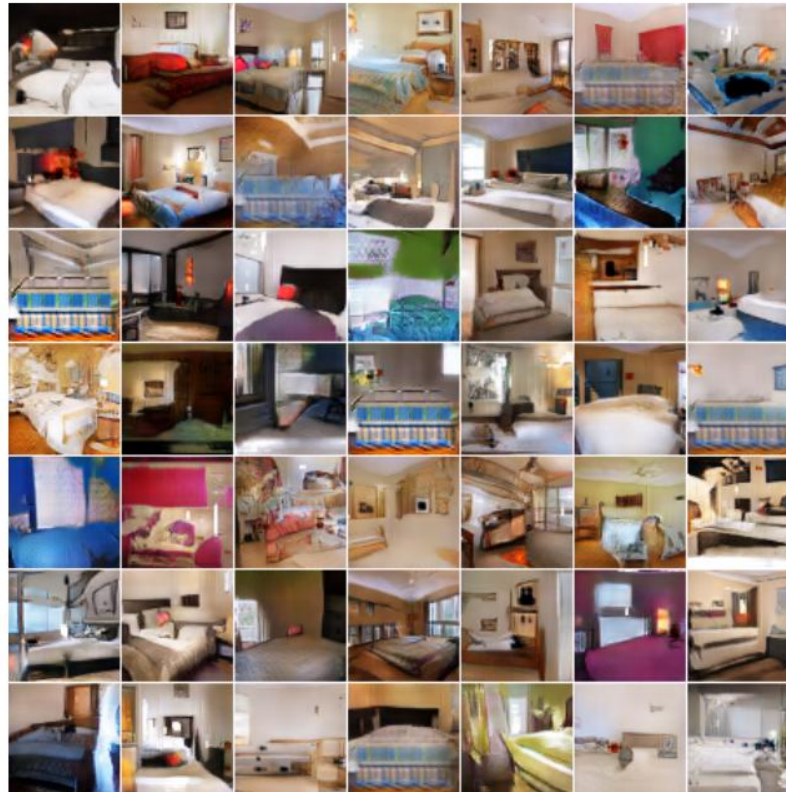
9. Generative adversarial networks*

Prof. Jianguo Zhang
SUSTech

*Based on <http://www.deeplearningbook.org> chapter 20.10.4 + various sources cited in the slides

Generative tasks

- Generation (from scratch): learn to sample from the distribution represented by the training set
- *Unsupervised learning task*



Generative tasks

- Conditional generation

goldfish



indigo
bunting



redshank



saint
bernard

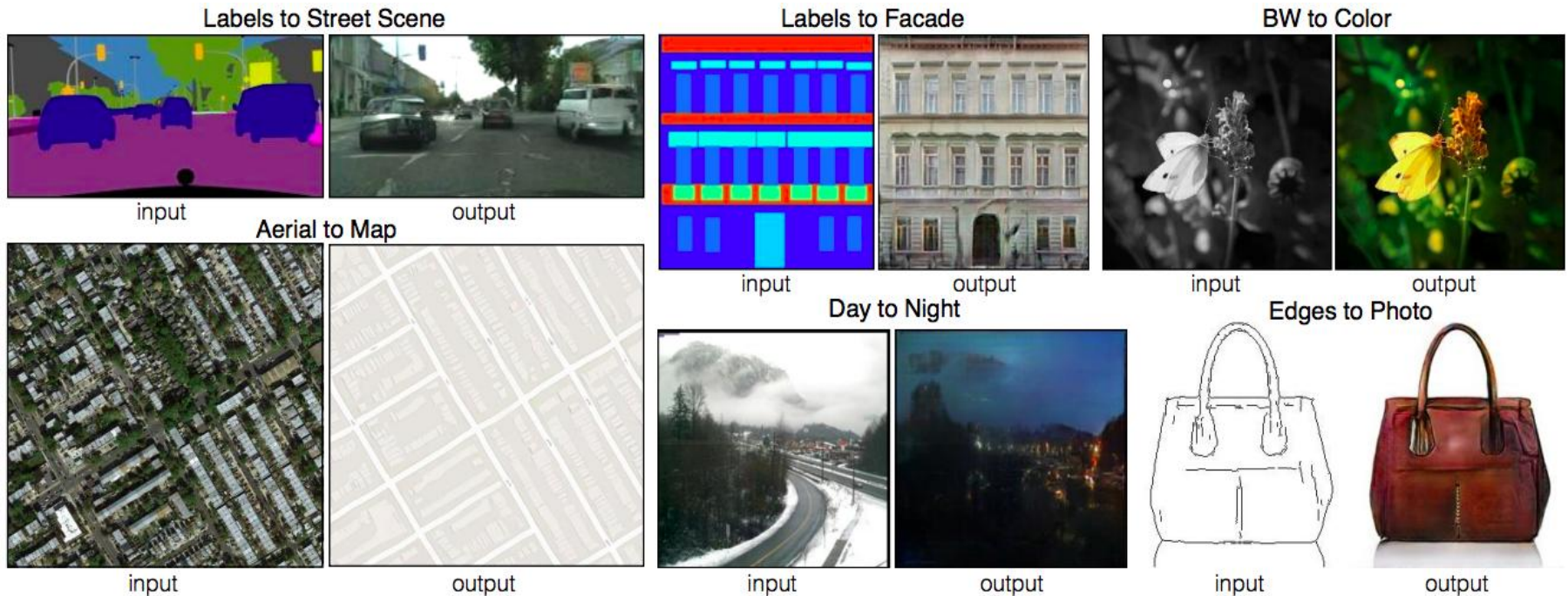


tiger
cat



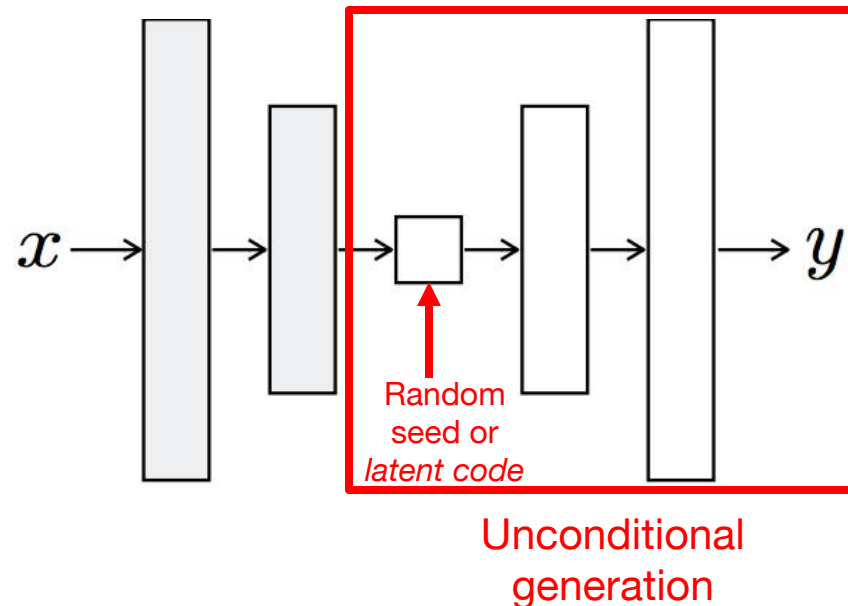
Generative tasks

- Image-to-image translation



Designing a network for generative tasks

- We saw that VAEs can learn to generate data by training both an encoder and a decoder (generator)
- Can we learn just the generator?
 - Recall the decoder of VAE



Learning to sample



Training data $x \sim p_{\text{data}}$

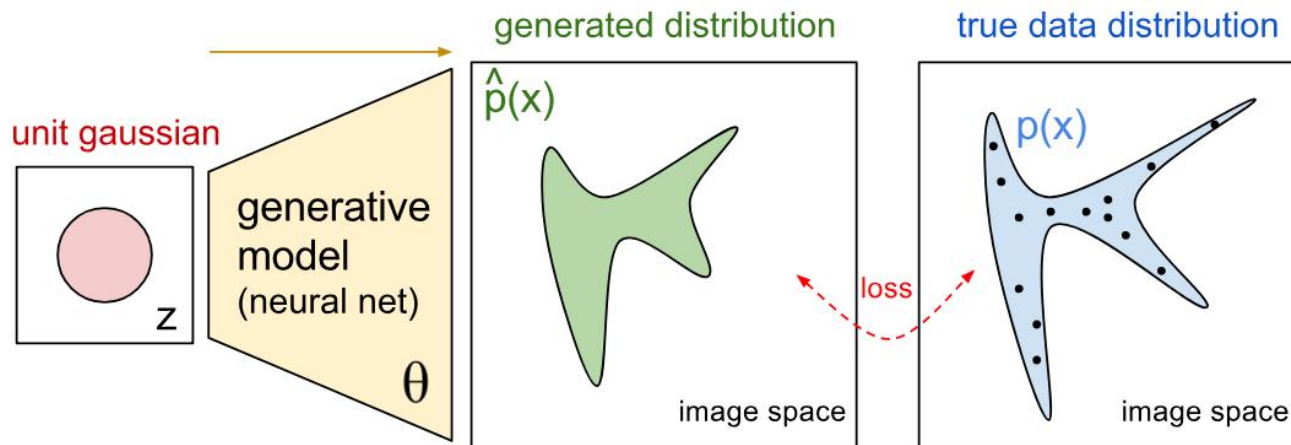


Generated samples $x \sim p_{\text{model}}$

We want to learn p_{model} that matches p_{data}

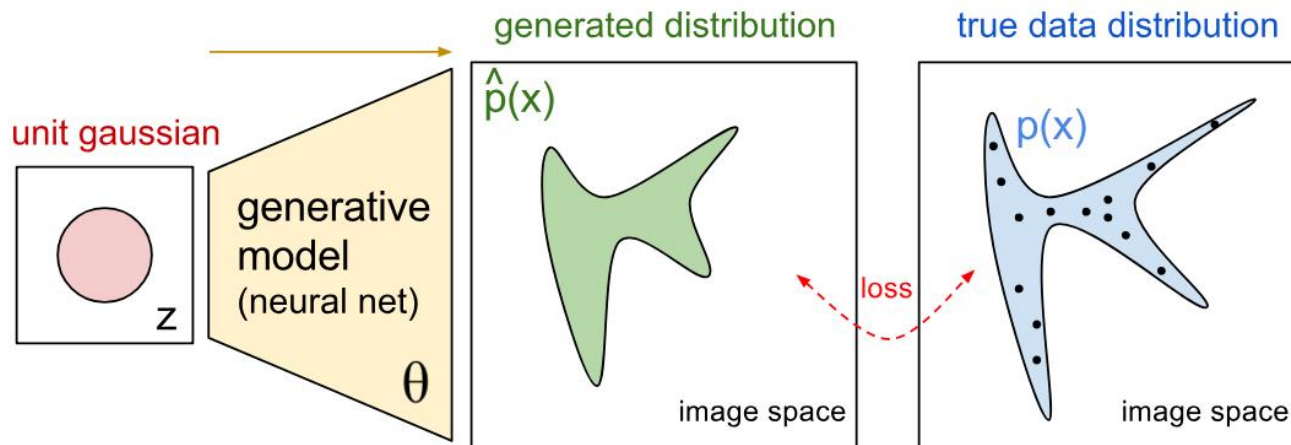
From VAEs to GANs

- We saw that VAEs can learn to generate data by training both an encoder and a decoder (generator)
- Can we learn just the generator?



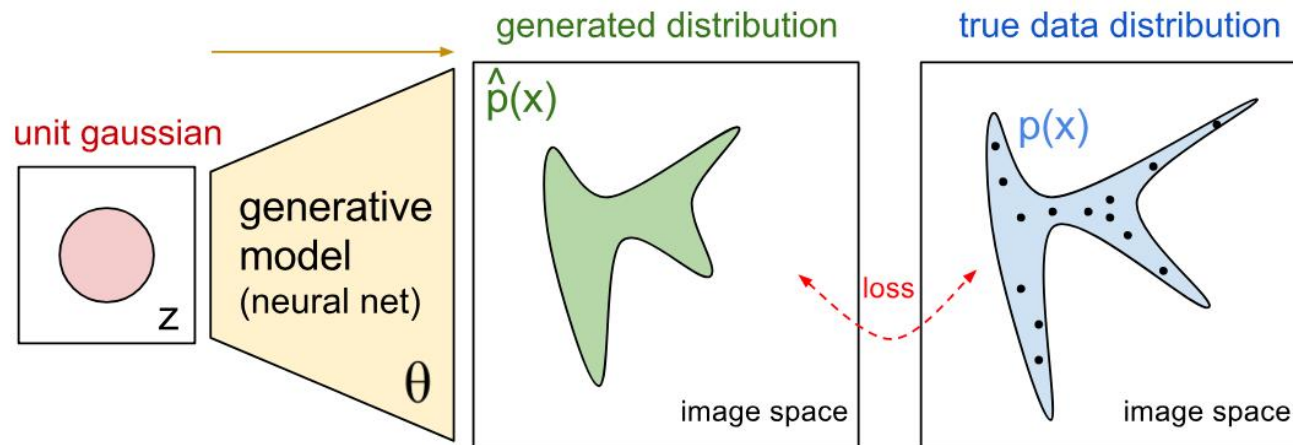
From VAEs to GANs

- But, how do we measure the quality of the generator without an encoder?
- In other words, what's the loss function?



Generative adversarial networks

- GANs propose to learn the loss function
- The training process is a game between two networks



Generative adversarial networks

- Train two networks with opposing objectives:
 - **Generator:** learns to generate samples
 - **Discriminator:** learns to distinguish between generated and real samples

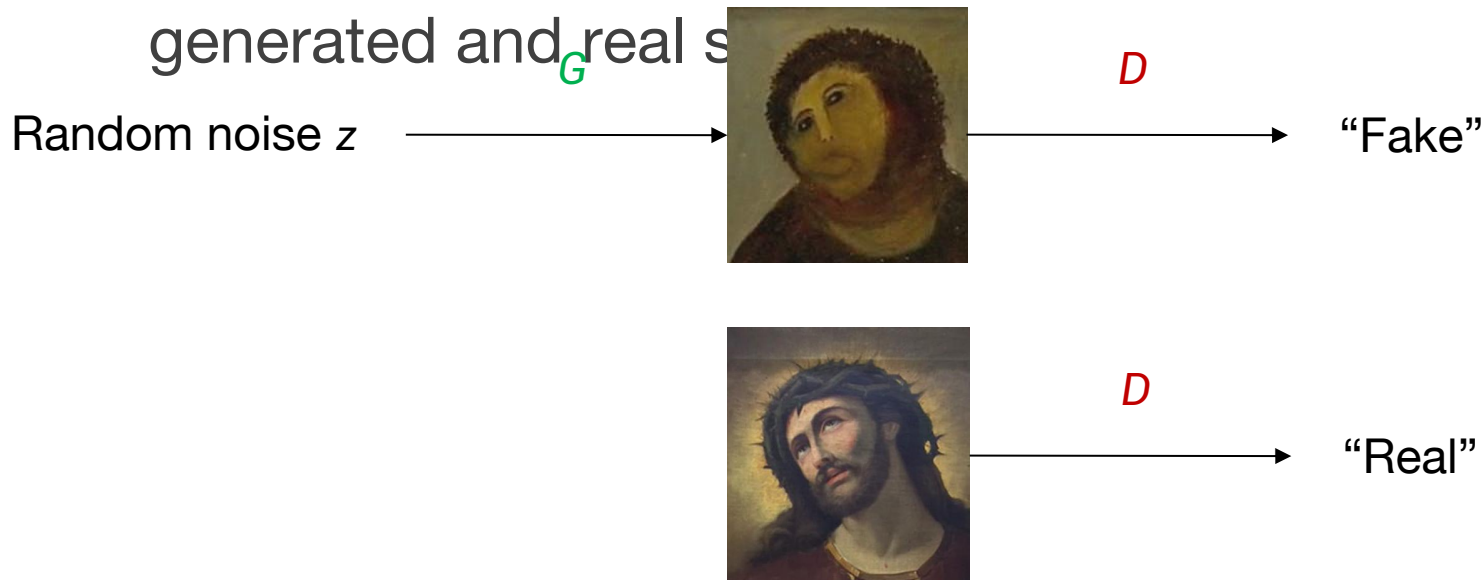
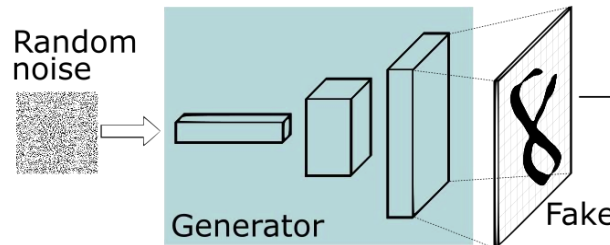


Figure adapted
from [F. Fleuret](#)

I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, [Generative adversarial nets](#), NIPS

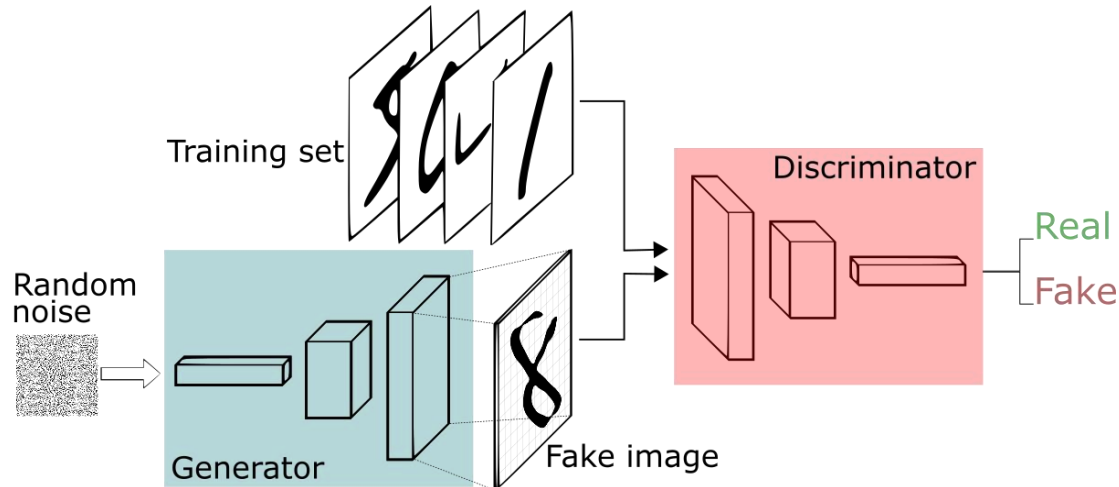
Generative adversarial networks

- **Generator:** $G(z)$ takes random noise z as input and outputs a (fake) image



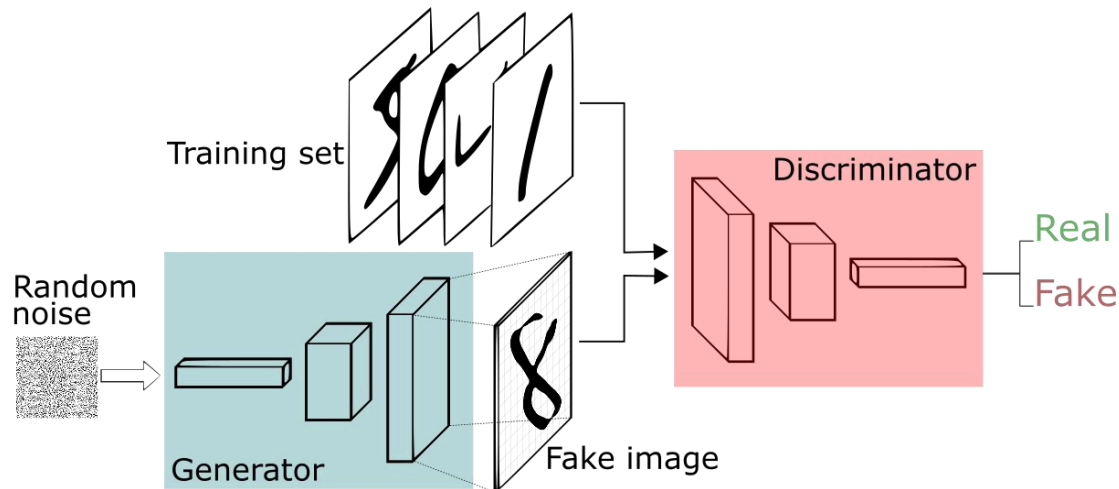
Generative adversarial networks

- **Generator:** $G(z)$ takes random noise z as input and outputs a (fake) image
- **Discriminator:** $D(x)$ receives an image x in input, real or fake, and estimates its probability to be real



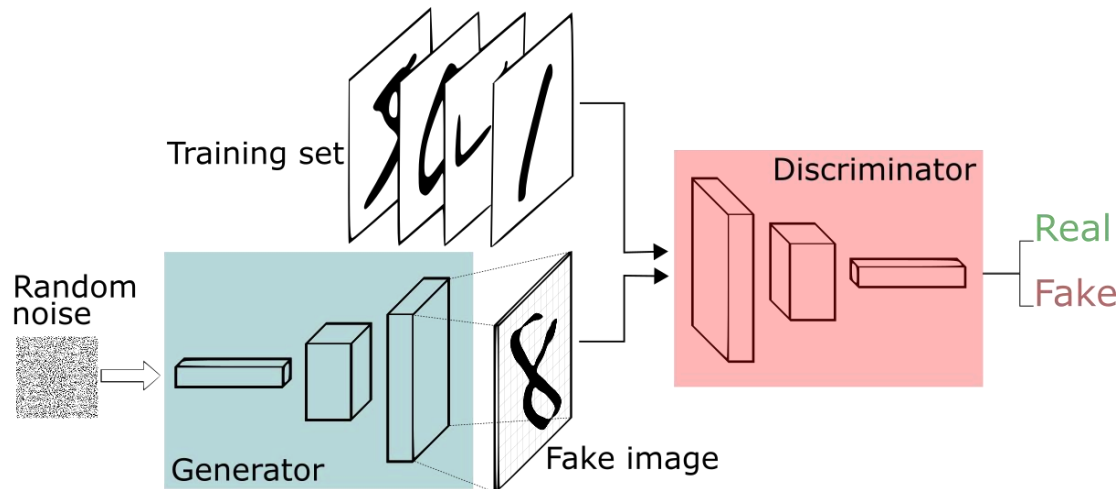
Generative adversarial networks

- **Adversarial training:** the generator tries to fool the discriminator while the discriminator tries to get better at distinguishing fake vs real images



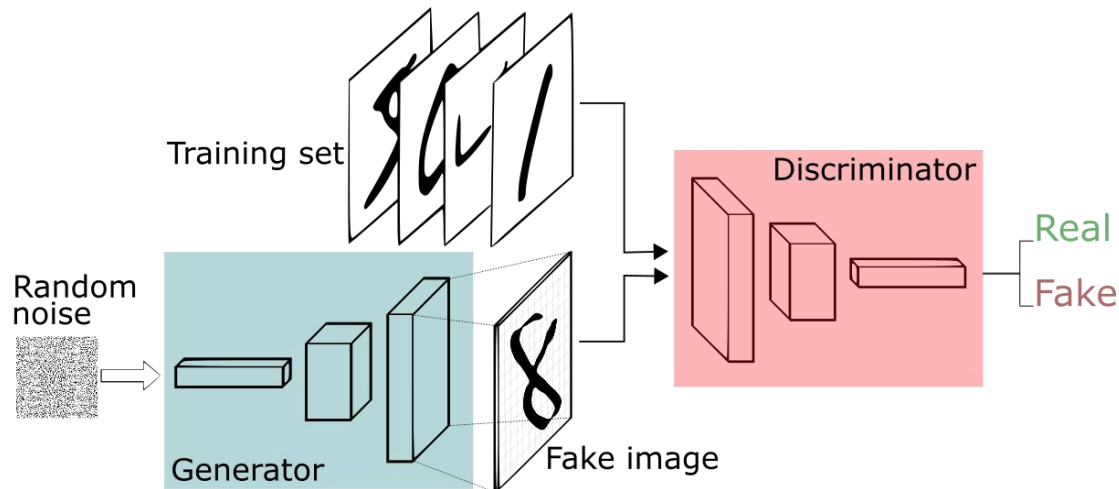
Generative adversarial networks

- When the discriminator spots a fake the generator adjusts its parameters, until at the end the generator reproduces the true data distribution and the discriminator is unable to find differences



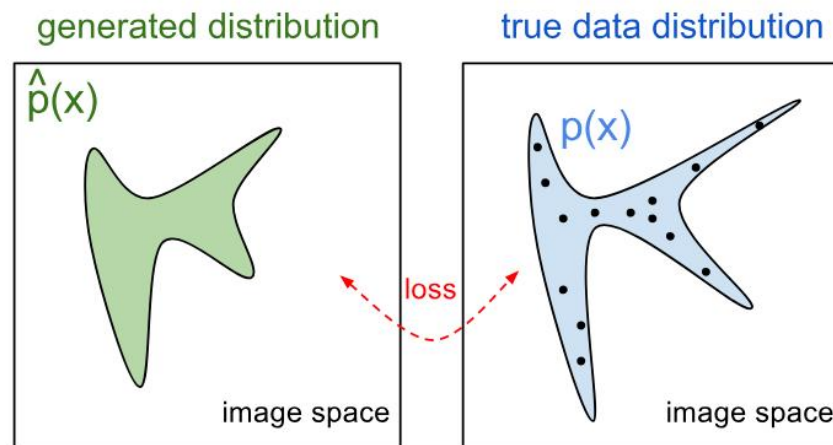
Generative adversarial networks

- **Note:** both the generator and the discriminator need to be differentiable
- Typically both implemented as (deep) neural **networks**, so we can use backpropagation

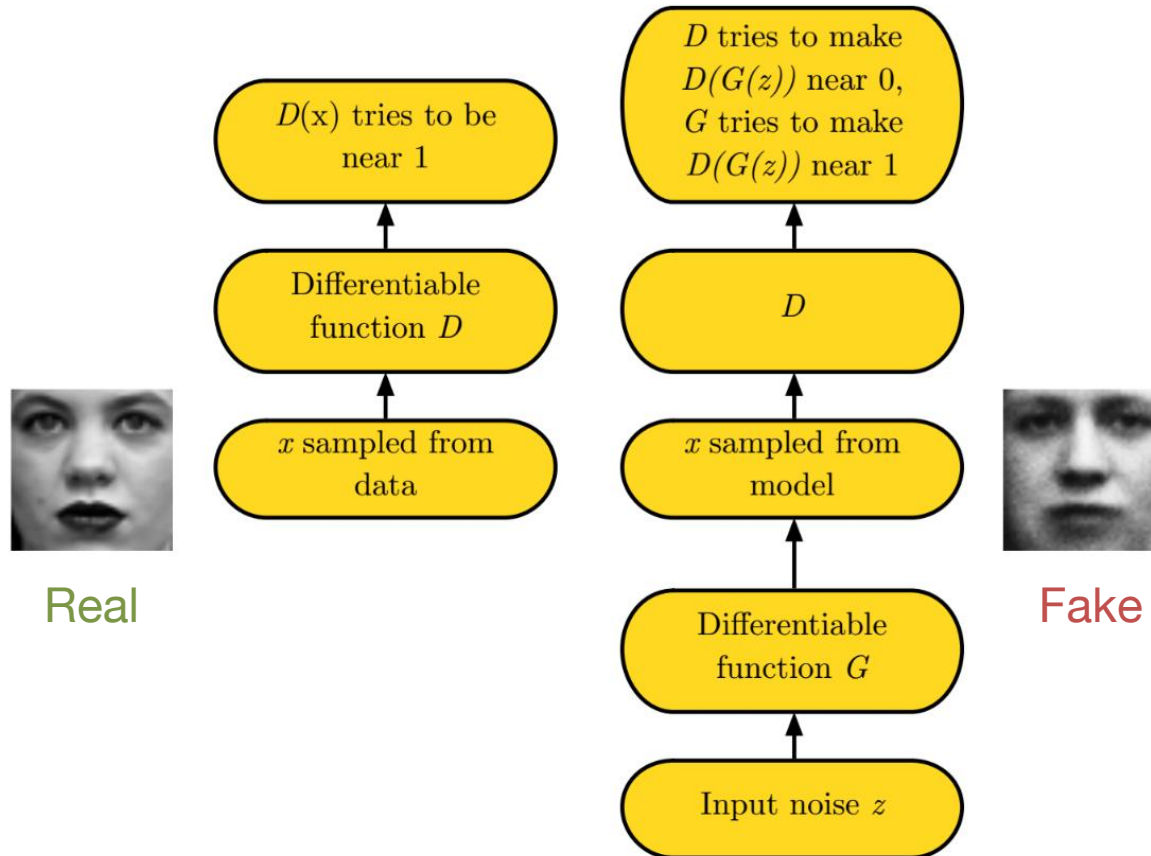


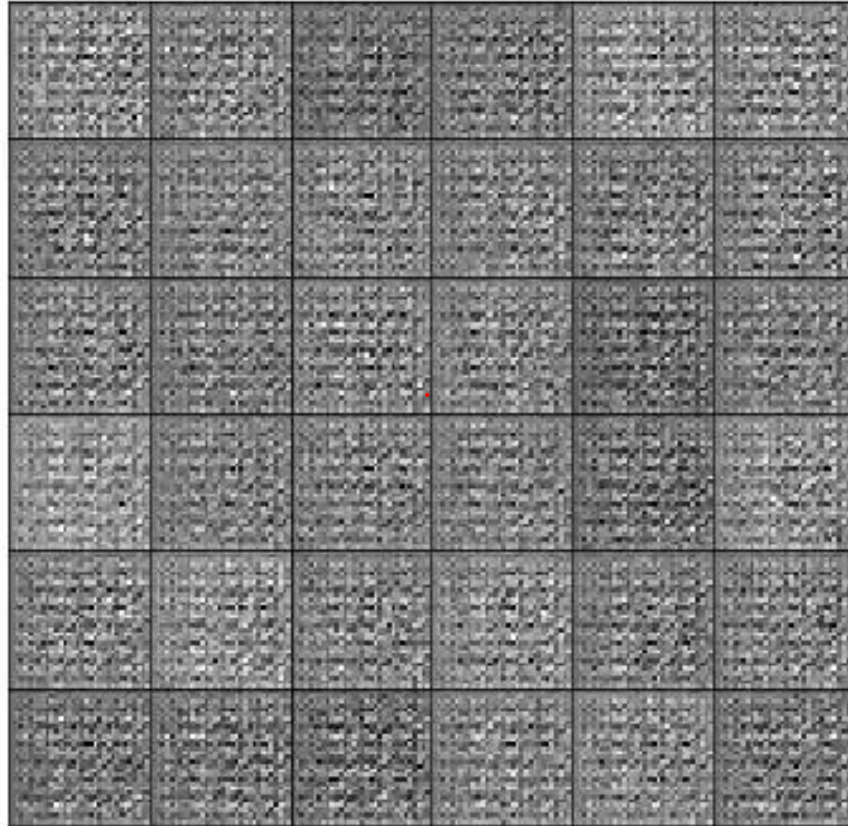
Generative adversarial networks

- That's like saying that the discriminator at the end is unable to tell the difference between the generated data distribution and the true data distribution (where training data comes from)



GANs pipeline







Real



Fake

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Prior distribution on input noise variables

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

The generator is a differentiable function (e.g., MLP) mapping noise to data space

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

The discriminator is a differentiable function (e.g., MLP) with single scalar output, i.e., the probability that \mathbf{x} comes from the true data distribution

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

$D(\mathbf{x}) = 1$ when the discriminator thinks \mathbf{x} is a real image

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

$D(G(\mathbf{z})) = 1$ when the discriminator thinks $G(\mathbf{z})$ is a real image

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

D is trained to maximise the probability of giving the correct label to samples from the true data distribution (the label should be 1) as well as to samples from G (the label should be 0)

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Large when $D(x)$ close to 1

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Large when $D(G(\mathbf{z}))$ close to 0

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Simultaneously G is trained to minimise $\log(1 - D(G(\mathbf{z})))$, i.e., fool the discriminator

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

- For a fixed G, the loss is effectively the binary cross-entropy
- The minimax strategy is used for **zero-sum games**, i.e., loss of one player = gain of the adversary
- The minimax solution is the **Nash equilibrium**
Nash equilibrium is a state in which D and G don't have any incentive to change because doing so would make the value of the objective function worse

Minimax and zero-sum games

- G and D follow a **minimax strategy**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

- For a fixed G, the loss is effectively the binary cross-entropy
- The minimax strategy is used for **zero-sum games**, i.e., loss of one player = gain of the adversary
- The minimax solution is the **Nash equilibrium**
In machine learning terms it means we reached convergence of the gradient descent for the joint optimisation problem

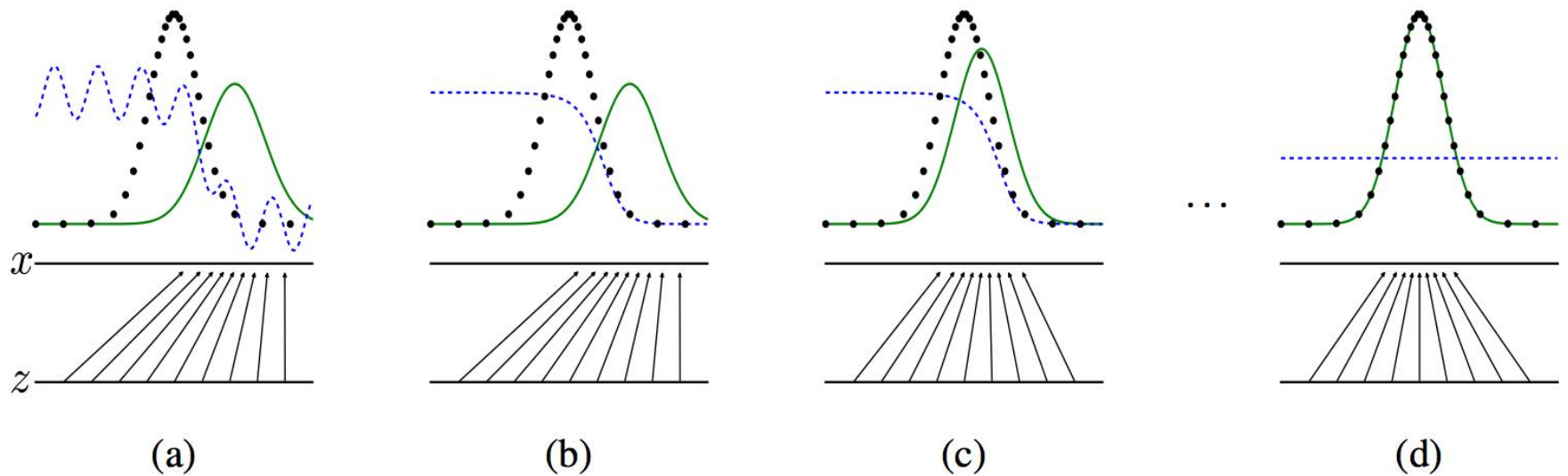


Figure 1: Generative adversarial nets are trained by simultaneously updating the discriminative distribution (D , blue, dashed line) so that it discriminates between samples from the data generating distribution (black, dotted line) p_x from those of the generative distribution p_g (G) (green, solid line). The lower horizontal line is the domain from which z is sampled, in this case uniformly. The horizontal line above is part of the domain of x . The upward arrows show how the mapping $x = G(z)$ 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 $D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$. (c) After an update to G , gradient of D has guided $G(z)$ 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. $D(x) = \frac{1}{2}$.

GAN training

- $V(G, D) = \mathbb{E}_{x \sim p_{\text{data}}} \log D(x) + \mathbb{E}_{z \sim p} \log(1 - D(G(z)))$

- Alternate between

- *Gradient ascent* on discriminator:

$$D^* = \arg \max_D V(G, D)$$

- *Gradient descent* on generator (minimize log-probability of discriminator being right):

$$\begin{aligned} G^* &= \arg \min_G V(G, D) \\ &= \arg \min_G \mathbb{E}_{z \sim p} \log(1 - D(G(z))) \end{aligned}$$

- In practice, do *gradient ascent* on generator (maximize log-probability of discriminator being wrong):

$$G^* = \arg \max_G \mathbb{E}_{z \sim p} \log(D(G(z)))$$

Training GANs

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:

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

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:

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

end for

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

Training GANs

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:

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

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:

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

In practice this saturates in the beginning since D rejects the samples with high confidence because they are clearly different from the training data

end for

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

Training GANs

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:

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

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:

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

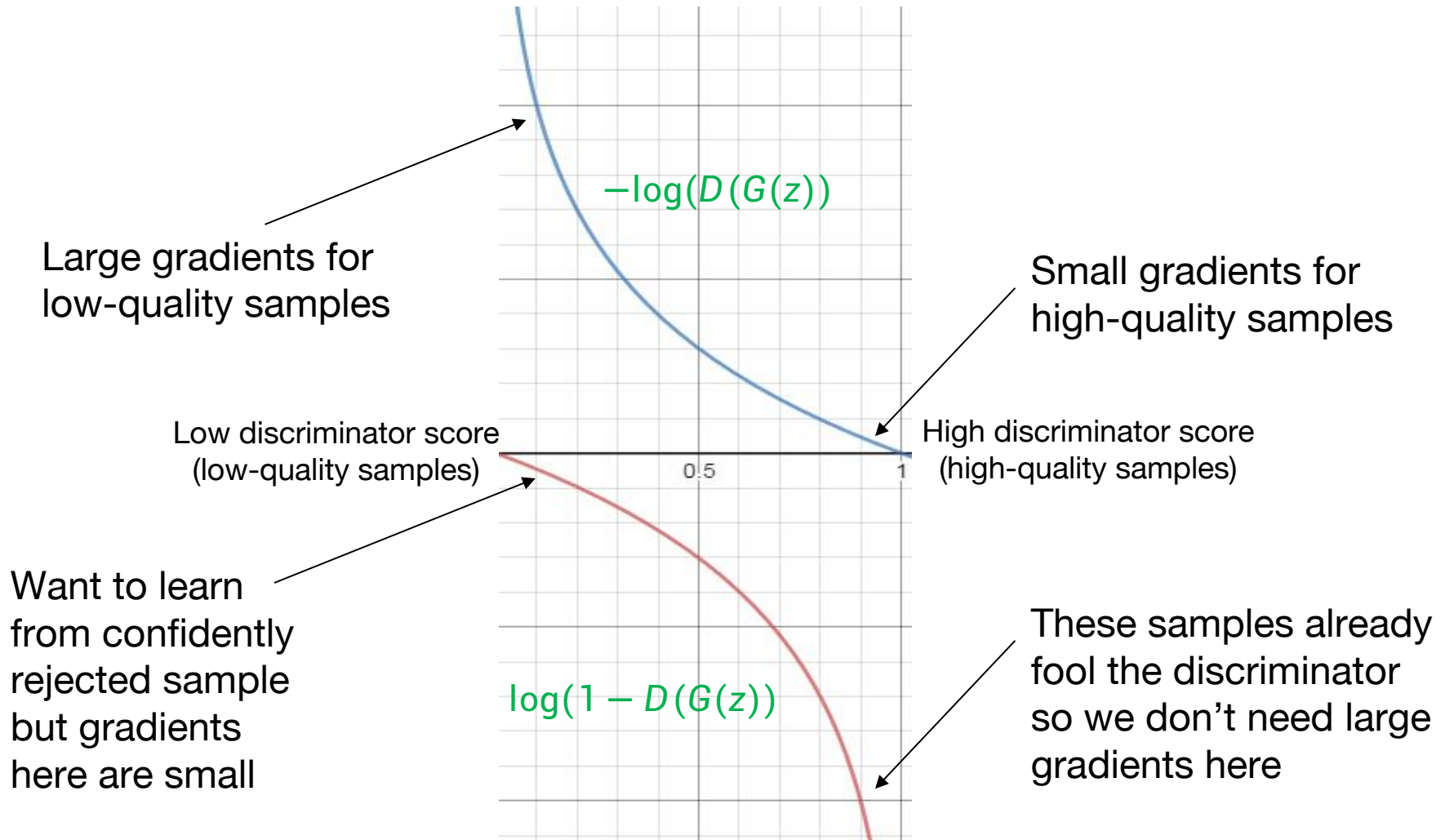
So we use this loss function for the generator instead, i.e., the generator maximises the log-prob of the discriminator being mistaken

end for

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

GAN training

$$\min_{w_G} \mathbb{E}_{z \sim p} \log(1 - D(G(z))) \text{ vs. } \max_{w_G} \mathbb{E}_{z \sim p} \log(D(G(z)))$$



Example generator: DCGAN

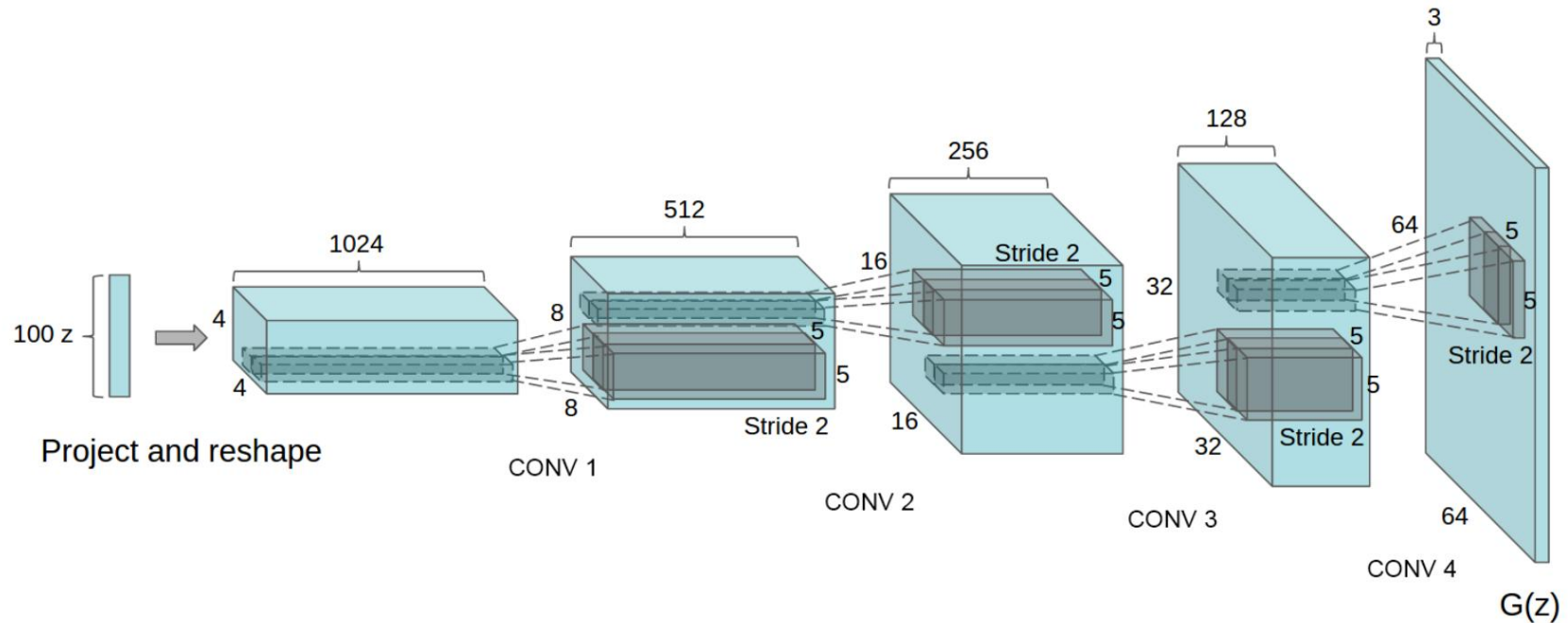
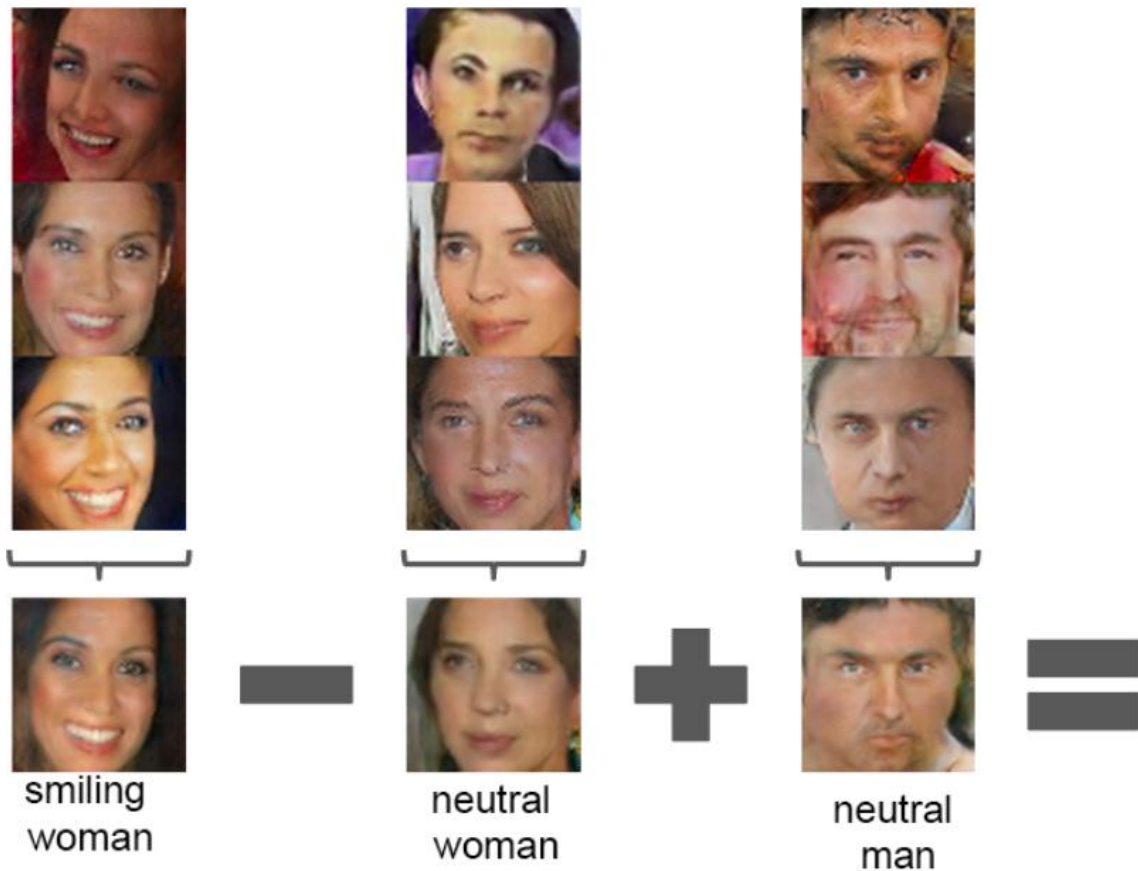


Figure 1: DCGAN generator used for LSUN scene modeling. A 100 dimensional uniform distribution Z is projected to a small spatial extent convolutional representation with many feature maps. A series of four fractionally-strided convolutions (in some recent papers, these are wrongly called deconvolutions)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

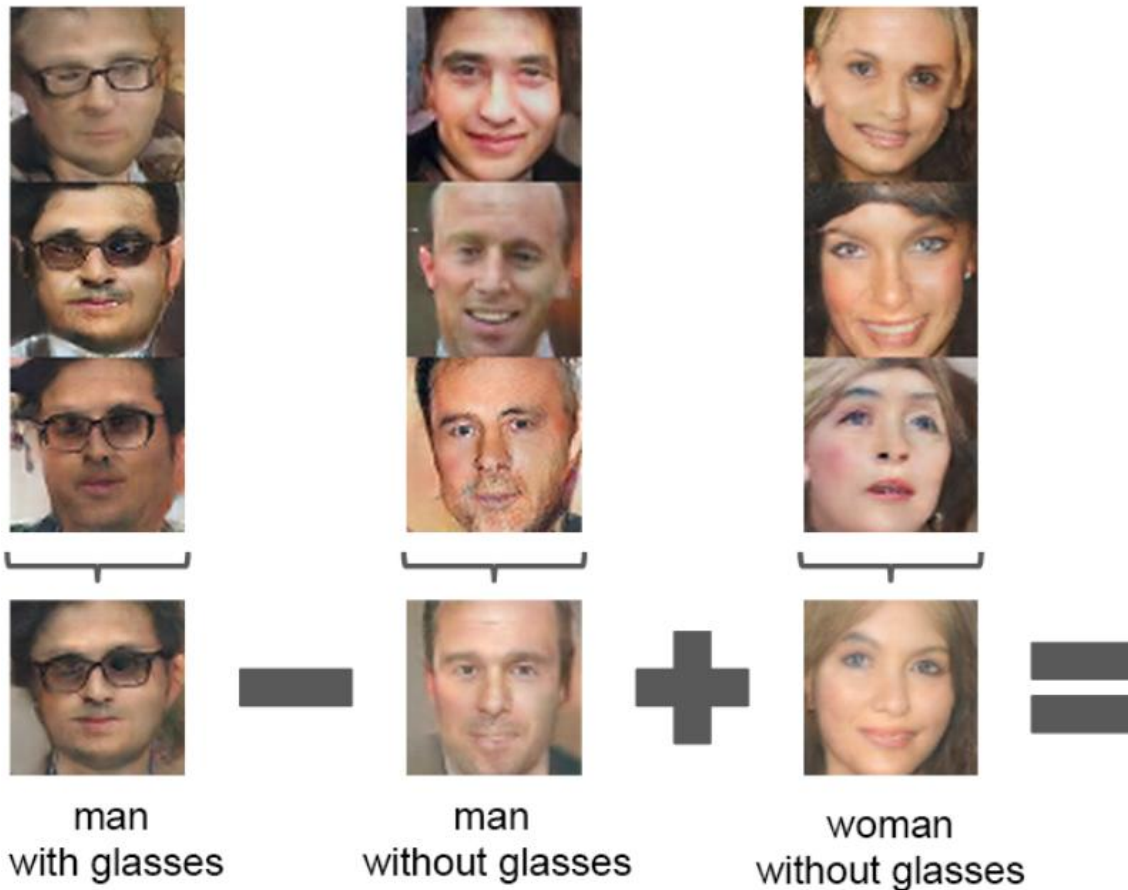
DCGAN results

- Vector arithmetic in the z space



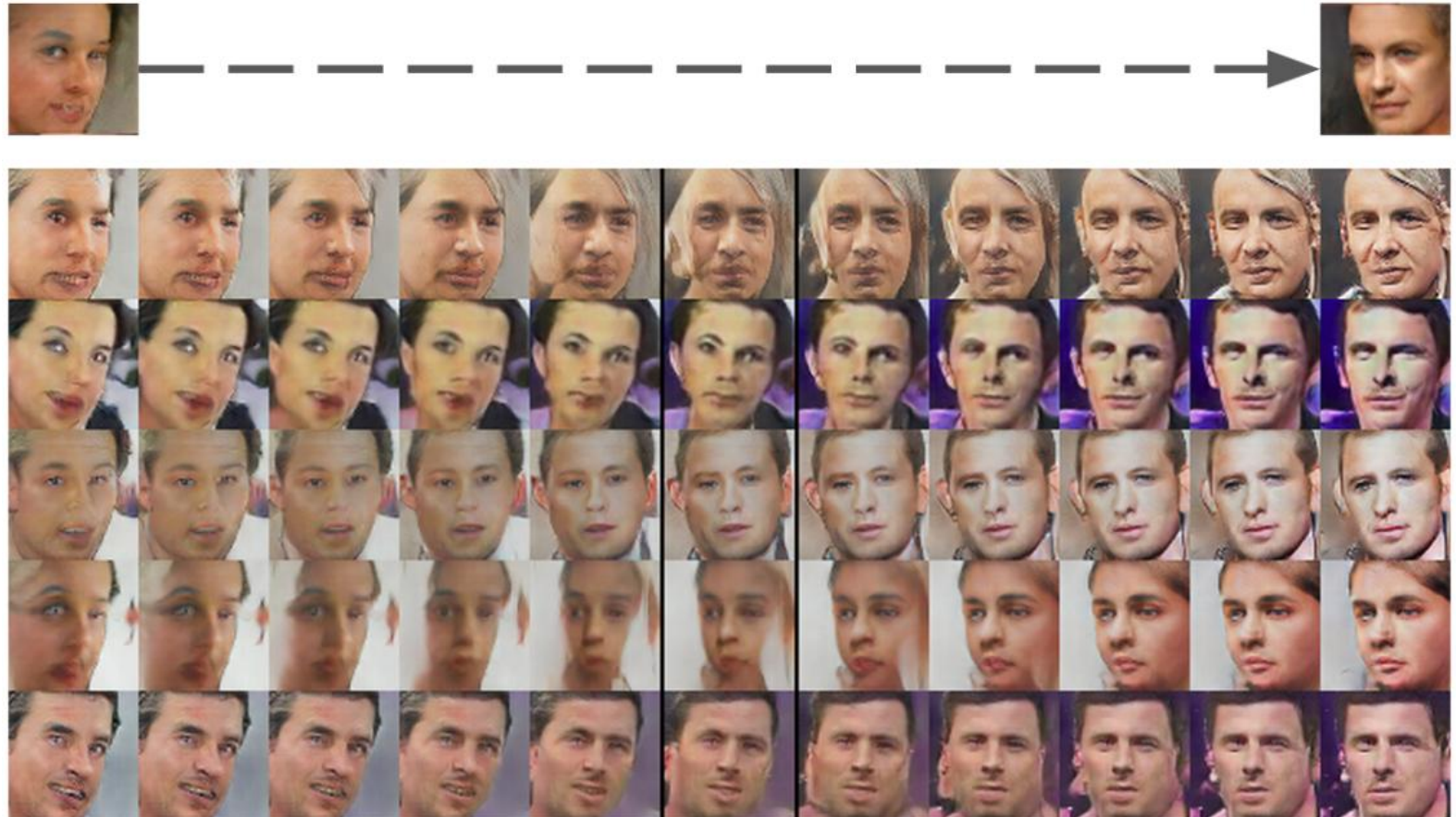
DCGAN results

- Vector arithmetic in the z space



DCGAN results

- Pose transformation by adding a “turn” vector



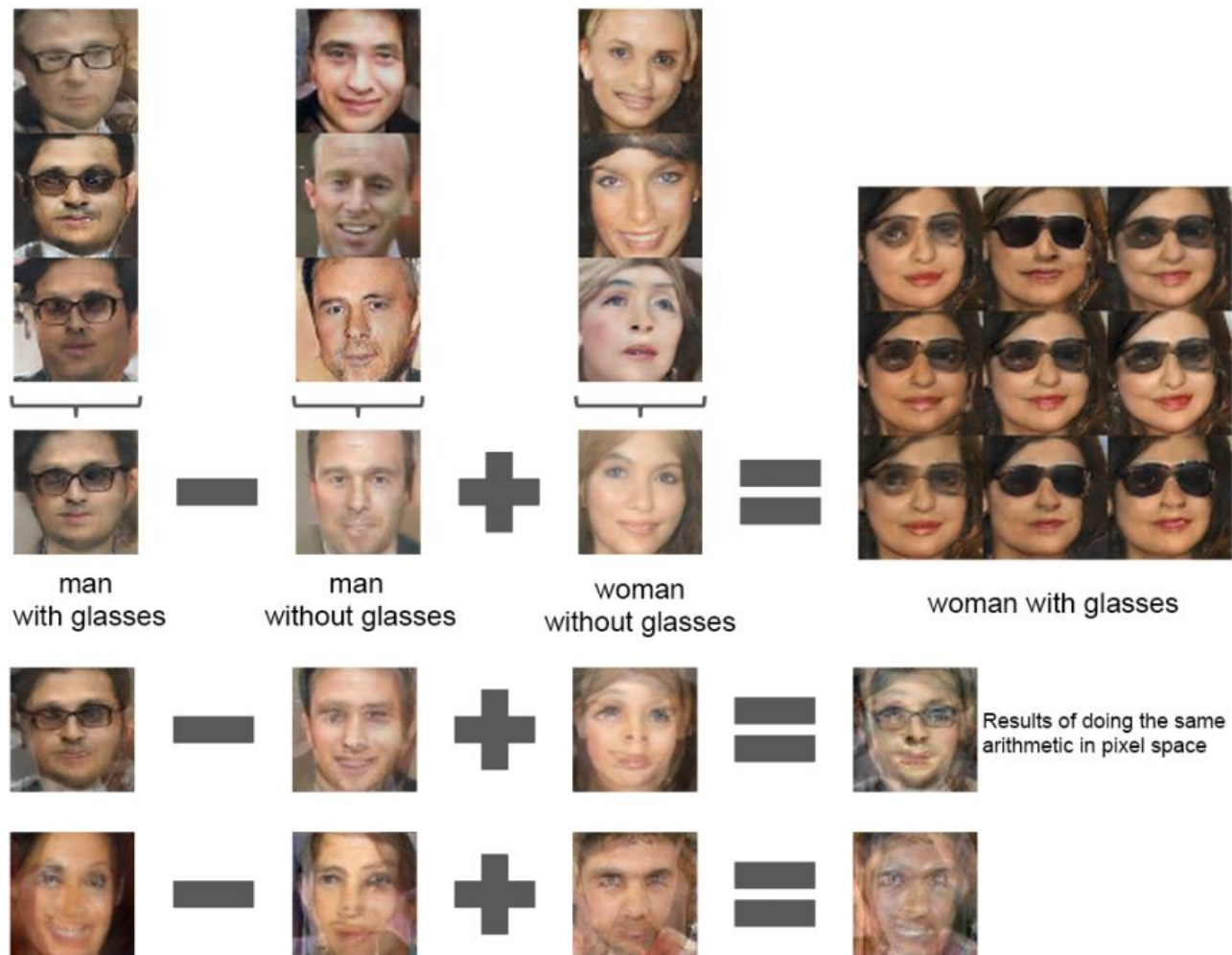


Figure 7: Vector arithmetic for visual concepts. For each column, the Z vectors of samples are averaged. Arithmetic was then performed on the mean vectors creating a new vector Y . The center sample on the right hand side is produce by feeding Y as input to the generator. To demonstrate the interpolation capabilities of the generator, uniform noise sampled with scale ± 0.25 was added to Y to produce the 8 other samples. Applying arithmetic in the input space (bottom two examples) results in noisy overlap due to misalignment.

Problems with GANs

- In the real world we face multiple problems when training GANs:
 - **Finite sample size**: training set is finite, not the full distribution
 - **Limited capacity**: the generator has limit capacity, i.e. cannot perfectly represent any distribution
 - **Optimisation errors**: optimisers can get stuck in local optima or never exactly converge to global optima

Problems with GANs

- In the real world we face multiple problems when training GANs:
 - **Saddle point problem**: harder than finding a maximum or minimum
 - **Balancing updates**: D too weak/strong means no gradient for G to improve

M. Arjovsky, S. Chintala, L. Bottou, [Wasserstein generative adversarial networks](#)

Conditional generation (cGAN)

goldfish



indigo
bunting



redshank



saint
bernard

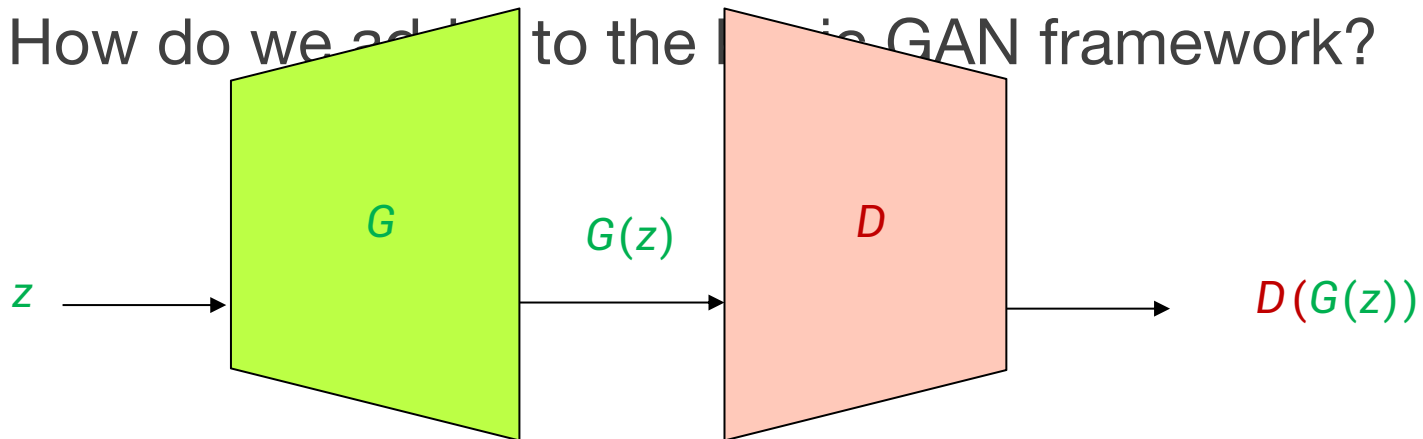


tiger
cat



Conditional generation (cGAN)

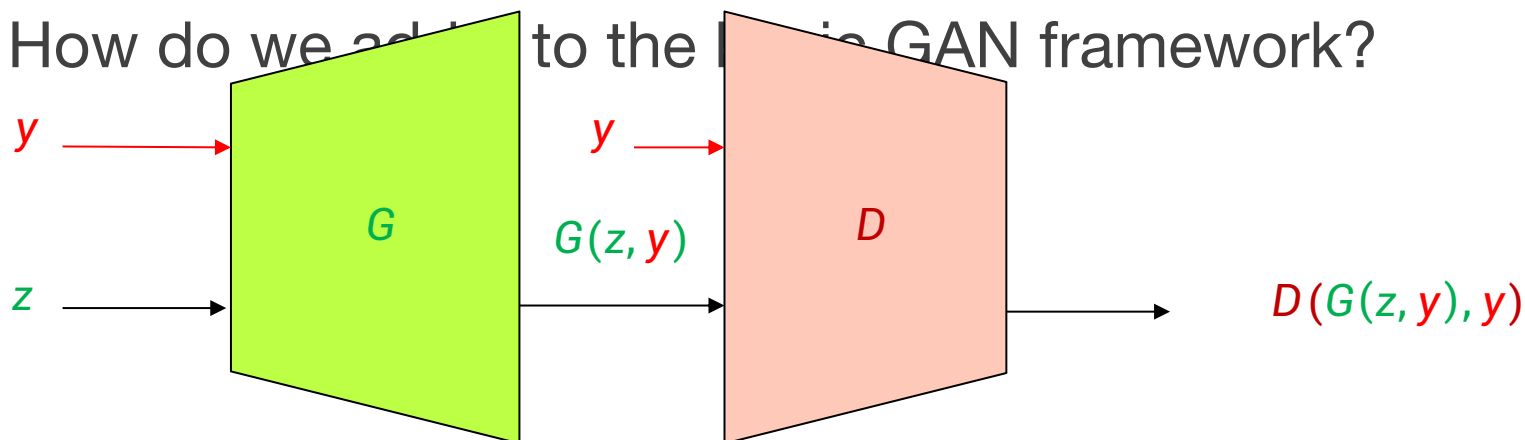
- Suppose we want to condition the generation of samples on discrete side information (label) y
- How do we add this to the basic GAN framework?



Conditional generation

- Suppose we want to condition the generation of samples on discrete side information (label) y

- How do we add this to the basic GAN framework?



Conditional generation

- Example: simple network for generating 28 x 28 MNIST digits

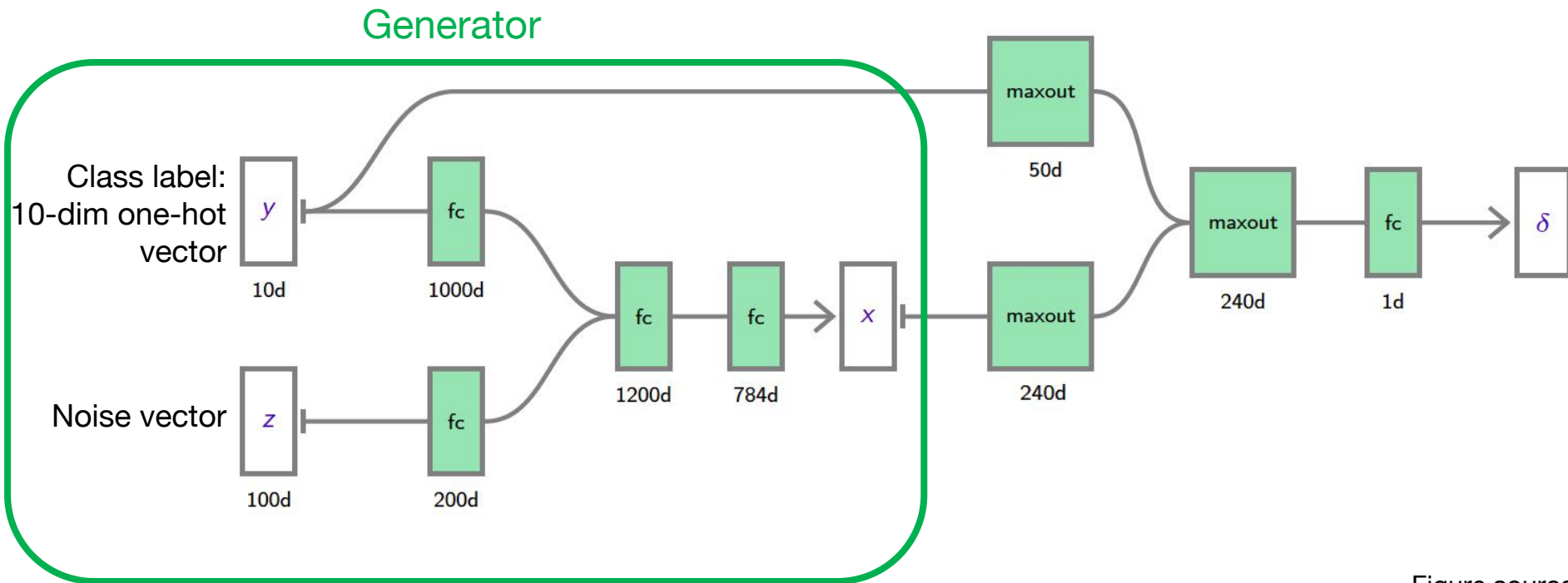


Figure source:
[F. Fleuret](#)

Conditional generation

- Example: simple network for generating 28 x 28 MNIST digits

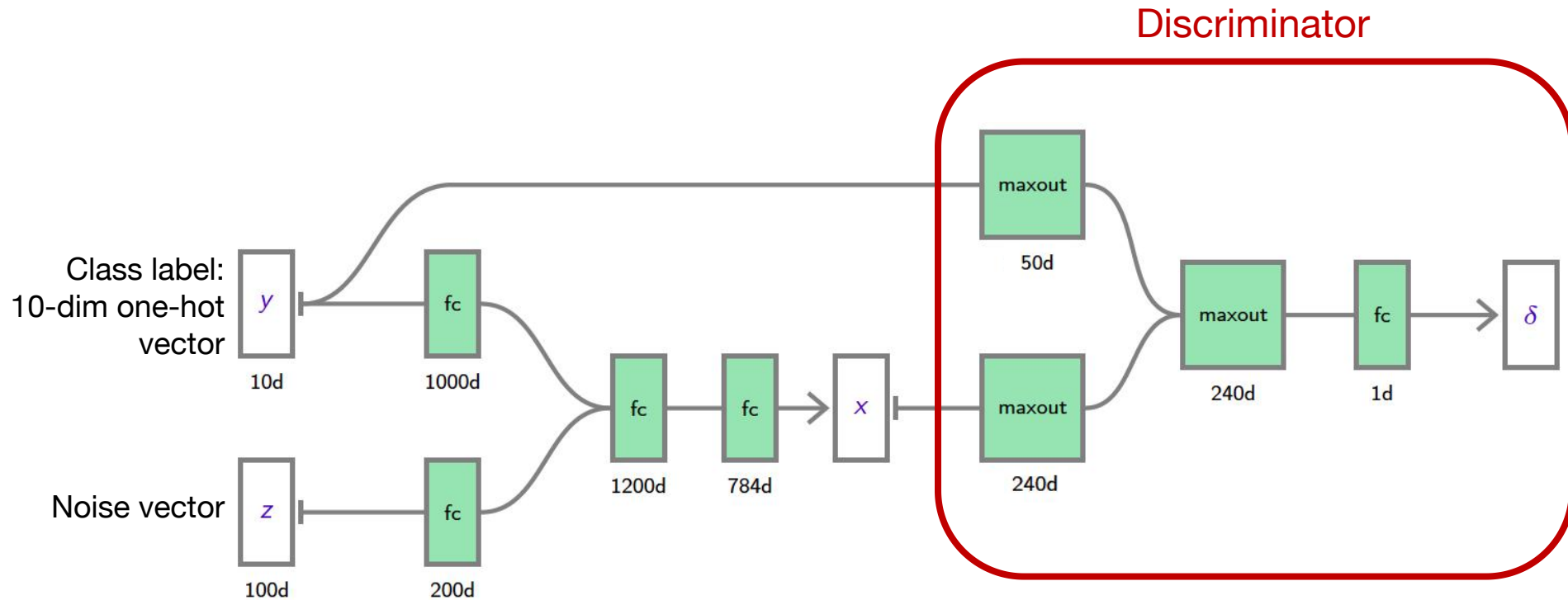


Figure source:
[F. Fleuret](#)

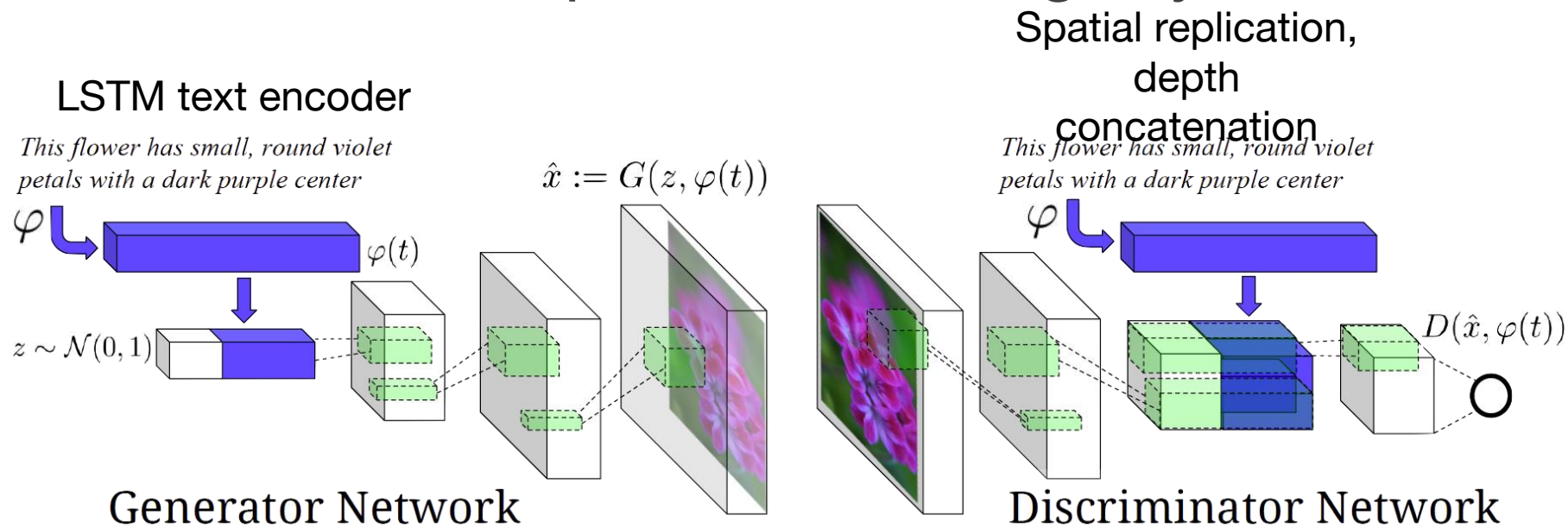
Conditional generation

- Example: simple network for generating 28 x 28 MNIST digits



Conditional generation

- Another example: text-to-image synthesis



Conditional generation

- Another example: text-to-image synthesis

Previously unseen
captions (*zero-shot*
setting)

this small bird has a pink
breast and crown, and black
primaries and secondaries.



the flower has petals that
are bright pinkish purple
with white stigma



this magnificent fellow is
almost all black with a red
crest, and white cheek patch.

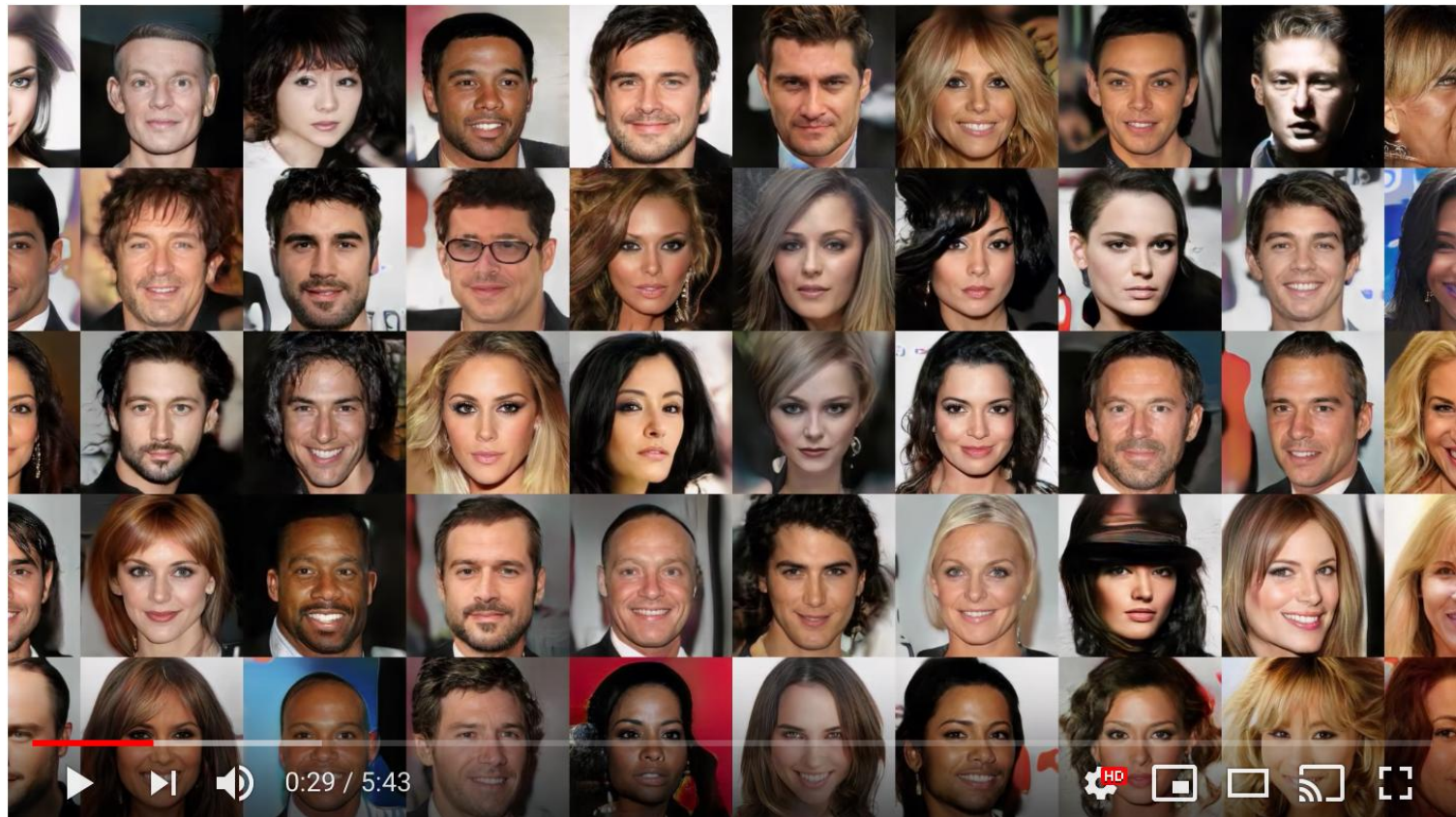


this white and yellow flower
have thin white petals and a
round yellow stamen



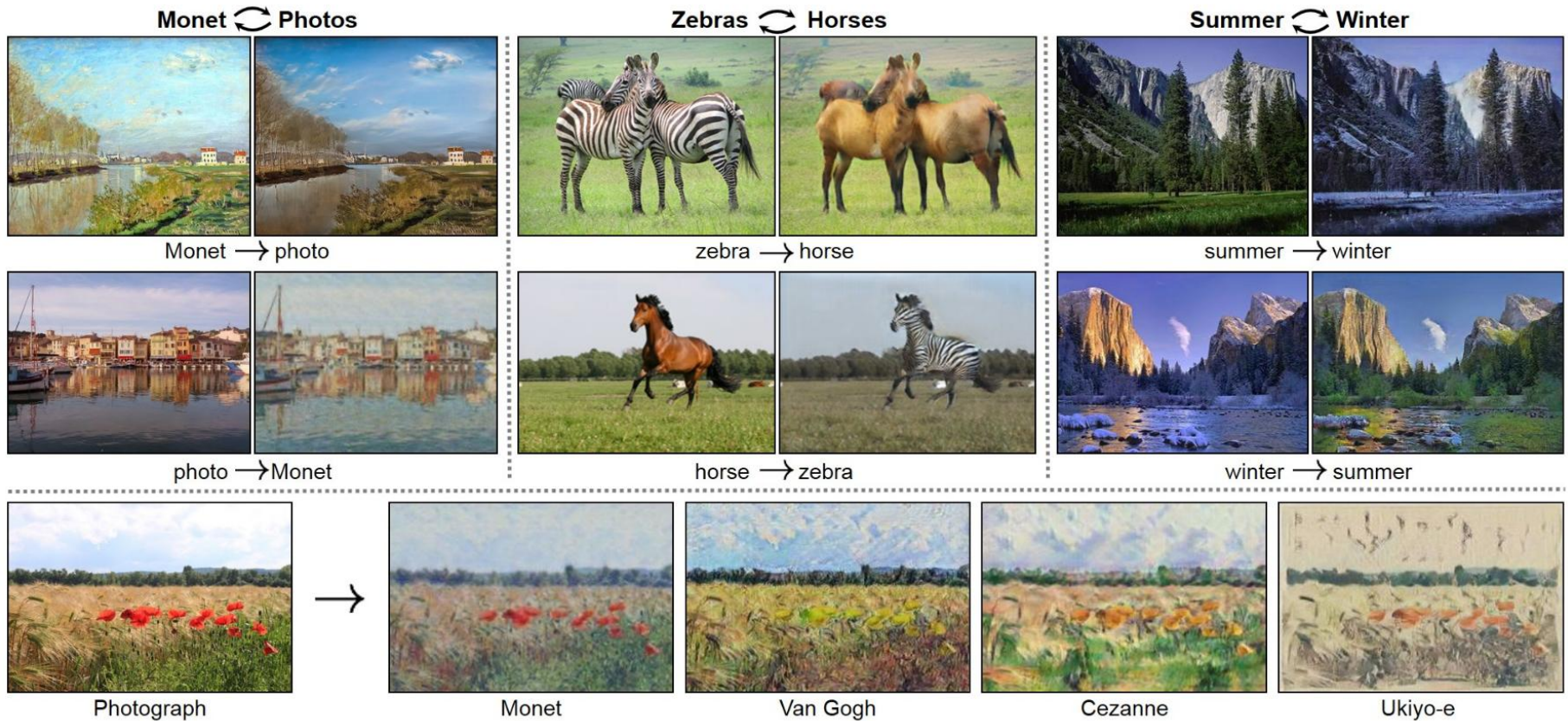
Captions seen in
the training set

Application: face generation

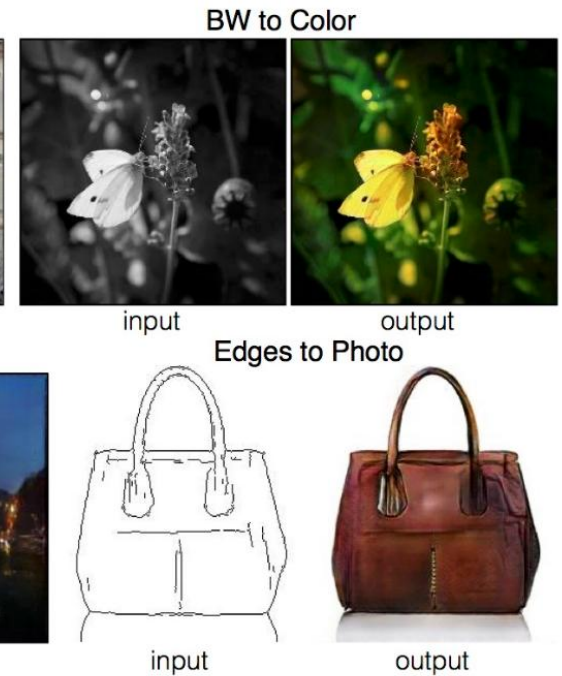
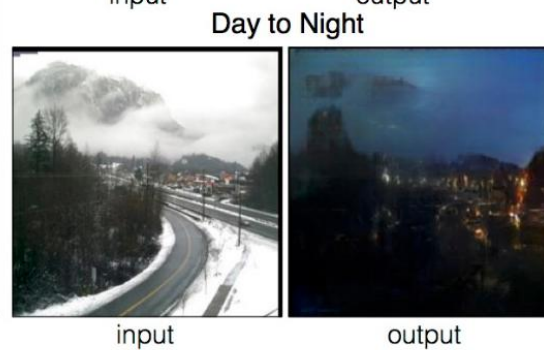
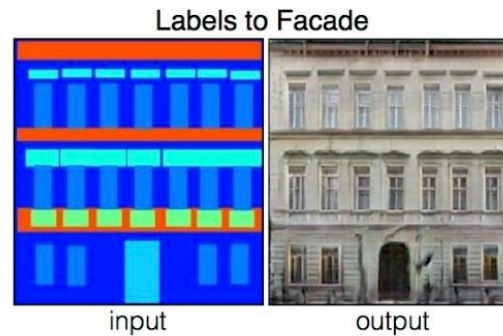


Source: <https://www.youtube.com/watch?v=G06dEcZ-QTg>

Application: style transfer



Application: image translation



Conclusion

- GANs provide an attractive way to create a generative model without using an explicit encode (as in variational auto-encoders) but using supervised learning
- However the training is far from trivial and in general finding Nash equilibria in high-dimensional non-convex games is still an important open research problem