

# CS-E4890: Deep Learning

## Generative adversarial networks

---

Alexander Ilin

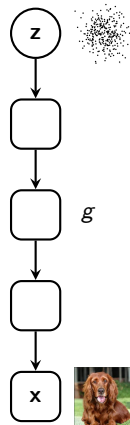
- Recall the generative model model of the VAE:
  - Hidden variables  $\mathbf{z}$  are normally distributed:  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$
  - Data vectors  $\mathbf{x}$  are nonlinear transformations of the latent variables with possible noise  $\epsilon$ :

$$\mathbf{x} = g(\mathbf{z}, \boldsymbol{\theta}) + \epsilon$$

- VAE is an explicit density model because we define an explicit parametric for of  $p(\mathbf{x})$ , for example:

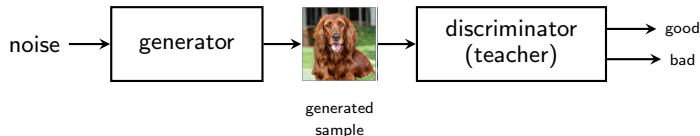
$$p(\mathbf{x}) = \prod_i \int p(\mathbf{z}_i) p(\mathbf{x}_i | \mathbf{z}_i, \boldsymbol{\theta}) d\mathbf{z}_i = \prod_i \int \mathcal{N}(\mathbf{z}_i | 0, \mathbf{I}) \mathcal{N}(\mathbf{x}_i | g(\mathbf{z}_i, \boldsymbol{\theta}), \sigma^2 \mathbf{I}) d\mathbf{z}_i$$

- In this lecture, we consider generative adversarial networks (GAN) which is an implicit density model. We can draw samples from the model but we do not explicitly define  $p(\mathbf{x})$ .



# Generative adversarial networks (GAN)

- GANs consist ([Goodfellow et al. 2014](#)) of a generator which generates samples and a discriminator which tells whether the generated samples are good or bad.



- The generator can be any parametric differentiable model (a deep neural network) that generates random samples.
- The discriminator is a classifier (a deep neural network) that classifies the generated samples into two classes.

- A popular choice for the GAN generation process:

- sample from an isotropic Gaussian distribution

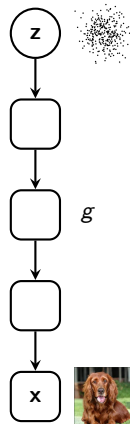
$$\mathbf{z} \sim \mathcal{N}(0, I)$$

- transform  $\mathbf{z}$  into the data space by a deep neural network

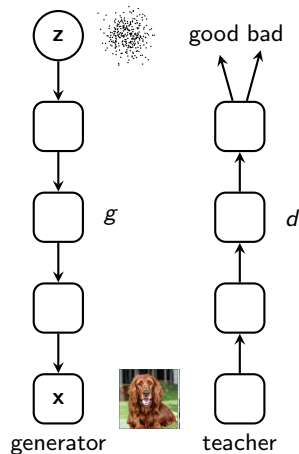
$$\mathbf{x} = g(\mathbf{z}, \theta)$$

- Other popular strategies:

- Apply additive or multiplicative noise to hidden layers.
  - Concatenate noise to hidden layers.



- The generator is guided by the teacher network that assesses the quality of the generated samples.
  - Teacher: a classifier that separates samples into classes “good” and “bad”.
- How to train the teacher  $d(\mathbf{x}, \theta_d)$ ?
  - Class “good”: samples from the training set.
  - Class “bad”: samples generated by the generator.
- The teacher is more traditionally called *discriminator*.



- The discriminator  $d(\mathbf{x})$  learns to separate generated samples from training data:

$$(\text{bad}=\text{generated}=\text{fake}) \quad 0 < d(\mathbf{x}) < 1 \quad (\text{good}=\text{real})$$

- The discriminator can be trained by maximizing the following log-likelihood function:

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log d(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim p_g} \log(1 - d(\mathbf{x})) \rightarrow \max_d$$

or equivalently:

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log d(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p_z} \log(1 - d(g(\mathbf{z}))) \rightarrow \max_d$$

- In practice, if we have  $N$  real examples  $\mathbf{x}_i$  and  $N$  generated examples  $g(\mathbf{z}_i)$ , we minimize the standard binary cross-entropy loss:

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

- The generator  $g(\mathbf{z})$  is trained to produce samples that are classified as real by the discriminator:

$$(\text{bad}=\text{generated}=\text{fake}) \quad 0 < d(\mathbf{x}) < 1 \quad (\text{good}=\text{real})$$

- The generator can be trained by maximizing the following function:

$$\mathbb{E}_{\mathbf{x} \sim p_g} \log d(\mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim p_z} \log d(g(\mathbf{z})) \rightarrow \max_g$$

- In practice, if we have  $N$  generated examples  $g(\mathbf{z}_i)$ , we can minimize the following loss:

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

## A non-saturating objective for the generator

- In principle, there are two ways to train the generator:

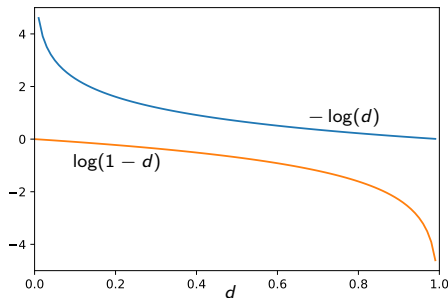
- To minimize the probability of being fake:

$$\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} \log(1 - d(g(\mathbf{z}))) \rightarrow \min_g$$

- To maximize the probability of being real:

$$\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} \log(d(g(\mathbf{z}))) \rightarrow \max_g$$

Both formulations result in the same fixed point.



- The latter formulation provides much stronger gradients. In the beginning of training, the discriminator can reject samples produced by the generator with high confidence ( $d \approx 0$ ). In this regime, the generator receives almost no gradient information, which slows down the convergence.



## 1. Update the discriminator:

- Sample  $N$  examples  $\mathbf{x}_i$  from the training set.
- Generate  $N$  samples  $g(\mathbf{z}_i)$  using the generator.
- 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  $\theta_d$  by stochastic gradient descent:  $\theta_d \leftarrow \theta_d - \nabla_{\theta_d} \mathcal{L}_d$

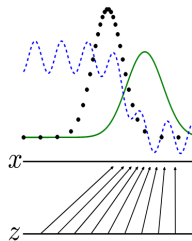
## 2. Update the generator:

- Generate  $N$  samples  $g(\mathbf{z}_i)$  using the generator.
- Compute the loss function

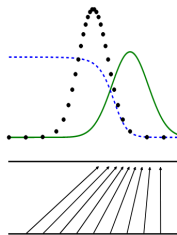
$$\mathcal{L}_g = -\frac{1}{N} \sum_{i=1}^N \log d(g(\mathbf{z}_i))$$

- Update  $\theta_g$  by stochastic gradient descent:  $\theta_g \leftarrow \theta_g - \nabla_{\theta_g} \mathcal{L}_g$  (gradients flow through the discriminator).

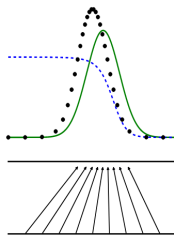
# Illustration of GAN training



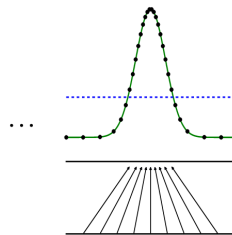
$d(x)$  is a partially accurate classifier.



$d(x)$  is trained to discriminate samples from data.



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.



After several steps of training,  $g$  and  $d$  reach a point at which both cannot improve because  $p_g = p_{\text{data}}$ . The discriminator is unable to differentiate between the two distributions  $d(x) = 0.5$ .

Images from (Goodfellow et al., 2014)

- A popular view at GANs: It is a two-player minimax game in which the generator tries to fool the discriminator and the discriminator tries to catch the fakes.
- The game can be described with one objective:

$$v(g, d) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log d(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} \log(1 - d(g(\mathbf{z})))$$
$$g^* = \arg \min_g \max_d v(g, d)$$

- The discriminator is trained to maximize  $v(g, d)$ . The generator tries to minimize  $v(g, d)$ .
- The equilibrium (also known as Nash equilibrium) is a saddle point of  $v$ .
- Nash equilibrium: No player has anything to gain by changing only their own strategy.

- For fixed generator  $g$ , the optimal discriminator is given by

$$d_g^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$$

- If we use the optimal discriminator  $d_g^*(\mathbf{x})$  to tune the generator, minimization of the GAN loss  $\mathcal{L}_g$  is equivalent to minimization of the Jensen-Shannon divergence between the model's distribution  $p_g$  and the data distribution  $p_{\text{data}}$ :

$$JSD(p_{\text{data}} \parallel p_g) = KL\left(p_{\text{data}} \parallel \frac{p_{\text{data}} + p_g}{2}\right) + KL\left(p_g \parallel \frac{p_{\text{data}} + p_g}{2}\right)$$

The global minimum of  $\mathcal{L}_g$  is achieved if and only if  $p_g = p_{\text{data}}$ .

- Training of GANs is often unstable: the convergence may be slow or difficult to achieve.
- A typical problem is so-called *mode collapse*: the generator produces the same output point (or slight variations of the same output, e.g., different views of the same dog) that the discriminator believes is most likely to be real rather than fake.
- There has been a lot of progress in GAN research and the results obtained with modern GANs look very impressive.

Zero-centered gradient penalties  
(Mescheder et al., 2018)

## A toy problem for studying the convergence of GANs

- Mescheder et al. (2018) studied the convergence in this optimization problem using a simple example:

- The true data distribution is a Dirac-distribution concentrated at 0.
- The generator distribution is  $p_\theta = \delta_\theta$ .
- The discriminator is linear  $d_\phi(x) = \phi x$ .

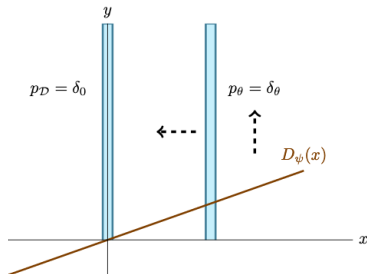
- We can write different variants of the GAN objective function as

$$v(\theta, \phi) = \mathbb{E}_{p_{\text{data}}(x)}[f(-d_\phi(x))] + \mathbb{E}_{p(z)}[f(d_\phi(g_\theta(z)))]$$

$$v(\theta, \phi) \rightarrow \min_{\theta} \max_{\phi}$$

where  $f(t) = -\log(1 + \exp(-t))$  yields the conventional GAN objective:

$$v(g, d) = \mathbb{E}_{x \sim p_{\text{data}}} \log d(x) + \mathbb{E}_{z \sim p_z(z)} \log(1 - d(g(z)))$$



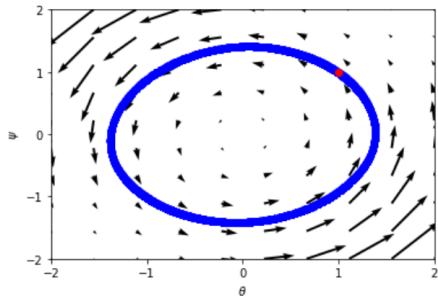
## Convergence of the original GANs

- Parameters  $\theta$ ,  $\phi$  can be optimized by simultaneous or alternating gradient descent.
- The optimization trajectories can be visualized on the 2d plane (visualizations are for alternating gradient descent).

- GANs in which the generator is trained as

$$\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} \log(1 - d(g(\mathbf{z}))) \rightarrow \min_g$$

- Training does not always converge to the Nash equilibrium.



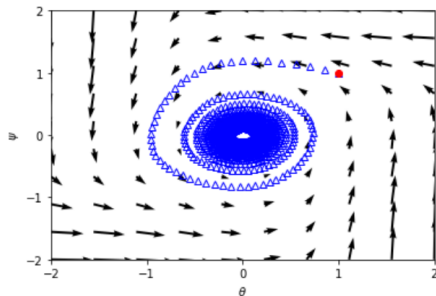


## Convergence of the original GANs with a non-saturating objective

- GANs in which the generator is trained as

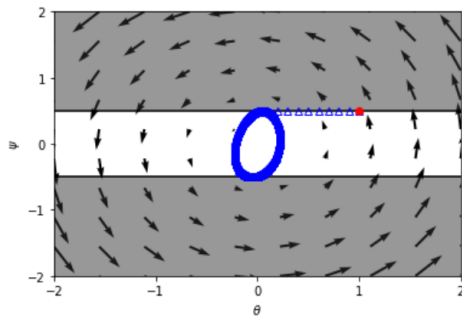
$$\mathbb{E}_{z \sim p_z(z)} \log(d(g(z))) \rightarrow \max_g$$

- Training converges but the convergence rate is extremely slow.

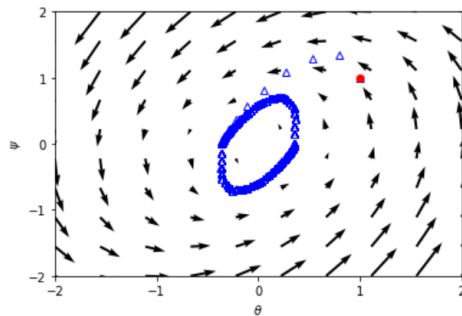


# Convergence of Wasserstein GAN

- WGAN (Arjovsky et al., 2017) and WGAN-GP (Gulrajani et al., 2017) with a finite number of discriminator updates per generator update do not always converge to the equilibrium point.



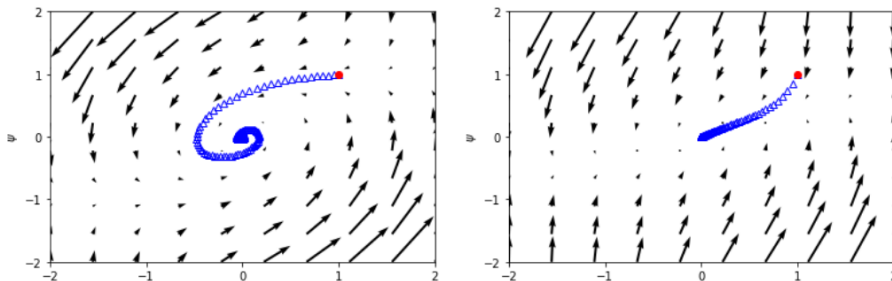
WGAN ( $n_d = 5$ )



WGAN-GP ( $n_d = 5$ )

## Zero-centered gradient penalties

- Penalizing the gradients of the discriminator has a positive effect on convergence: the discriminator is penalized for deviating from the Nash-equilibrium.
- The authors proposed *zero-centered* gradient penalty:
  - Penalize gradients on real data  $R_1 = \frac{\gamma}{2} \mathbb{E}_{x \sim p_{\text{data}}} (\|\nabla d(x)\|^2)$
  - Penalize gradients on generated data  $R_2 = \frac{\gamma}{2} \mathbb{E}_{x \sim p_g} (\|\nabla d(x)\|^2)$

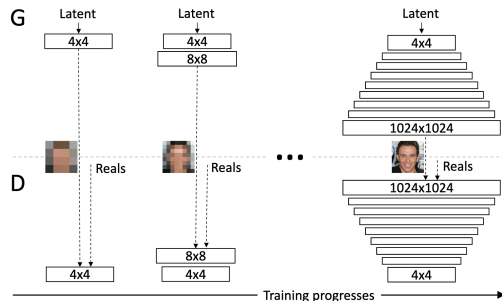


Training behavior with different  $\gamma$  ( $R_1$  and  $R_2$  are equivalent for the toy problem).

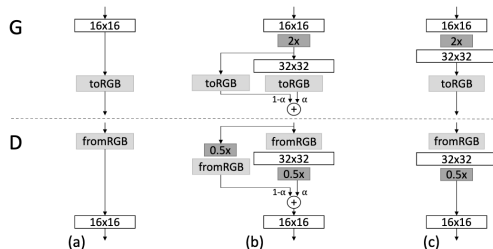
Progressive growing (ProGAN)  
(Karras et al., 2018)

## Progressive growing of GANs (Karras et al., 2018)

- The idea: start by building a generative model for low-resolution images, then progressively increase the resolution by adding layers to the networks.
- Generation of smaller low-resolution images is more stable because the problem is much simpler compared to the end goal.
- The training time is reduced because many iterations are done at lower resolutions.



- When doubling the resolution, the new layers are faded in smoothly: During the transition (b) the layers that operate on the higher resolution are treated like a residual block, whose weight  $\alpha$  increases linearly from 0 to 1.
- When training the discriminator, real images are downsampled to match the current resolution of the network.
- During the resolution transition, the authors interpolate between two resolutions (both in the generator and the discriminator) to mitigate possible negative effects caused by the change of the training scenario.



## ProGAN: Generated samples



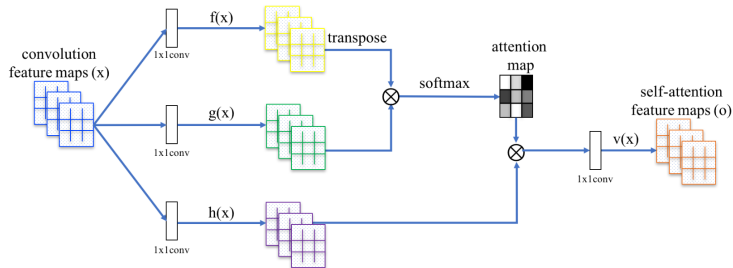
# Self-Attention GAN (SAGAN)

(Zhang et al., 2018)

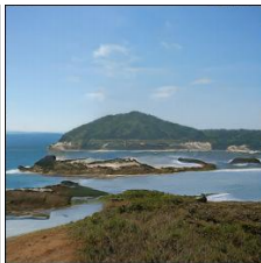
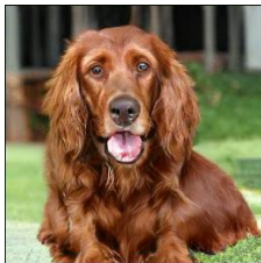


## Self-Attention GAN (SAGAN) (Zhang et al., 2018)

- SAGAN: Using a self-attention module (inspired by the transformers) in the generator.



- The results suggest that a self-attention module is beneficial for improving the quality of the generated images.



# Style-Based Generators (StyleGAN)

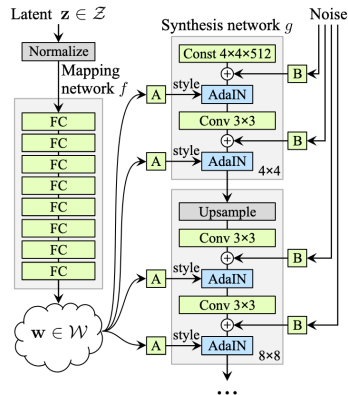
(Karras et al., 2018)

- Motivated by the style-transfer literature (Huang and Belongie, 2017), the authors re-design the generator architecture.
- The generator starts from a learned constant input.
- The “style” of the image at each convolution layer is adjusted using adaptive instance normalization:

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

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

- The style vectors ( $\mathbf{y}_s, \mathbf{y}_b$ ) are produced from latent code  $z$  by an MLP.
- Additional noise is injected directly into the network.

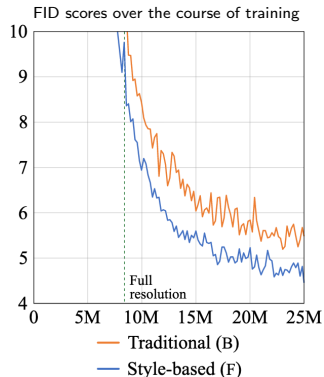


# StyleGAN improves the quality of generated images

- Automatic evaluation of the quality of generated images is not a trivial task. One popular metric is called Fréchet Inception distance (FID) (Heusel et al., 2017). It is a Fréchet distance between two Gaussian distributions  $\mathcal{N}(\mathbf{m}_r, \mathbf{C}_r)$  and  $\mathcal{N}(\mathbf{m}_g, \mathbf{C}_g)$ :

$$\text{FID} = \|\mathbf{m}_r - \mathbf{m}_g\|_2^2 + \text{Tr}(\mathbf{C}_r + \mathbf{C}_g - 2(\mathbf{C}_r \mathbf{C}_g)^{1/2})$$

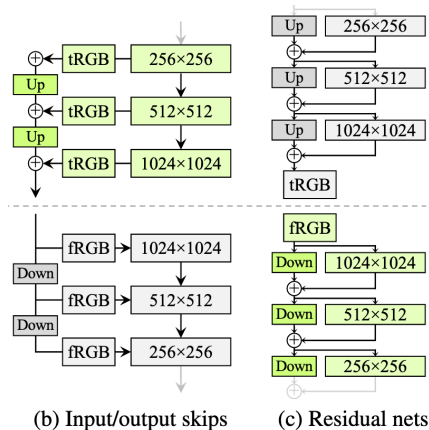
- Statistics  $\mathbf{m}_r$  and  $\mathbf{C}_r$  are computed in the following way:
  - Propagate real images through an Inception-v3 classifier pre-trained on natural images.
  - Compute mean  $\mathbf{m}_r$  and covariance matrix  $\mathbf{C}_r$  of the outputs of one of the layers.
- $\mathbf{m}_g$  and  $\mathbf{C}_g$  are computed similarly on generated images.



Horizontal axis denotes the number of training images seen by the discriminator. The dashed vertical line marks the point when training has progressed to full  $1024^2$  resolution.

- Replaced adaptive instance normalization with weight demodulation.
- New architectures of the generator and discriminator.
- Removed progressive growing.

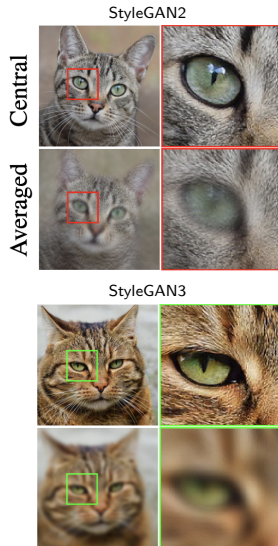
Configuration	FID ↓
A Baseline StyleGAN [24]	4.40
B + Weight demodulation	4.39
C + Lazy regularization	4.38
D + Path length regularization	4.34
E + No growing, new G & D arch.	3.31
F + Large networks (StyleGAN2)	<b>2.84</b>
Config A with large networks	3.98



New architectures of the generator and the discriminator in StyleGAN2.

## StyleGAN3: Alias-Free GAN (Karras et al., 2021)

- StyleGAN2 has a “texture sticking” problem: unwanted dependence of the synthesis process on absolute pixel coordinates.
  - Above: an image generated from a latent code (I guess the input of the synthesis network).
  - Below: The average of images generated from a small neighborhood around the same code.
  - The intended result is uniformly blurry because all details should move together. However, with StyleGAN2 many details stick to the same pixel coordinates, showing unwanted sharpness.
- StyleGAN3 fixes this problem:
  - Interpret all signals in the network as continuous
  - Re-design the architecture of the synthesis network to make it fully equivariant to translation and rotation.



Home assignment



- You need to implement and train on MNIST:
  - DCGAN

- Papers cited in the lecture slides.