

# Deep Generative Models

## Lecture 11

Roman Isachenko



Ozon Masters

Spring, 2021

## Recap of previous lecture

### Kantorovich-Rubinstein duality

$$W(\pi || p) = \frac{1}{K} \max_{\|f\|_L \leq K} [\mathbb{E}_{\pi(x)} f(x) - \mathbb{E}_{p(x)} f(x)],$$

where  $\|f\|_L \leq K$  are  $K$ -Lipschitz continuous functions  
( $f : \mathcal{X} \rightarrow \mathbb{R}$ ).

### Gradient penalty

$$W(\pi || p) = \underbrace{\mathbb{E}_{\pi(x)} f(x) - \mathbb{E}_{p(x)} f(x)}_{\text{original critic loss}} + \lambda \underbrace{\mathbb{E}_{U[0,1]} \left[ (\|\nabla_{\hat{x}} f(\hat{x})\|_2 - 1)^2 \right]}_{\text{gradient penalty}}.$$

Samples  $\hat{x}_t = t\mathbf{x} + (1-t)\mathbf{y}$  with  $t \in [0, 1]$  are uniformly sampled along straight lines between pairs of points:  $\mathbf{x}$  from the data distribution  $\pi(\mathbf{x})$  and  $\mathbf{y}$  from the generator distribution  $p(\mathbf{x}|\theta)$ .

## Recap of previous lecture

$$f(\mathbf{x}, \phi) = \mathbf{W}_{K+1} a_K (\mathbf{W}_K a_{K-1} (\dots a_1 (\mathbf{W}_1 \mathbf{x}) \dots)).$$

- ▶  $a_k$  is a pointwise nonlinearities. We assume that  $\|a_k\|_L = 1$  (it holds for ReLU).
- ▶  $\mathbf{g}(\mathbf{x}) = \mathbf{W}\mathbf{x}$  is a linear transformation ( $\nabla \mathbf{g}(\mathbf{x}) = \mathbf{W}$ ).

$$\|\mathbf{g}\|_L = \sup_{\mathbf{x}} \sigma(\nabla \mathbf{g}(\mathbf{x})) = \sigma(\mathbf{W}).$$

### Critic spectral norm

$$\|f\|_L \leq \|\mathbf{W}_{K+1}\| \cdot \prod_{k=1}^K \|a_k\|_L \cdot \|\mathbf{W}_k\| = \prod_{k=1}^{K+1} \sigma(\mathbf{W}_k).$$

## Spectral Normalization GAN

If we replace the weights in the critic  $f(\mathbf{x}, \phi)$  by  $\mathbf{W}_k^{SN} = \mathbf{W}_k / \sigma(\mathbf{W}_k)$ , we will get  $\|f\|_L \leq 1$ .

Power iteration approximates the value of  $\sigma(\mathbf{W})$ .

## Recap of previous lecture

### What is a divergence?

Let  $\mathcal{S}$  be the set of all possible probability distributions. Then  $D : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$  is a divergence if

- ▶  $D(\pi || p) \geq 0$  for all  $\pi, p \in \mathcal{S}$ ;
- ▶  $D(\pi || p) = 0$  if and only if  $\pi \equiv p$ .

### f-divergence minimization

$$D_f(\pi || p) = \mathbb{E}_{p(\mathbf{x})} f\left(\frac{\pi(\mathbf{x})}{p(\mathbf{x})}\right) \rightarrow \min_p .$$

### Variational divergence estimation

$$D_f(\pi || p) \geq \sup_{T \in \mathcal{T}} [\mathbb{E}_\pi T(\mathbf{x}) - \mathbb{E}_p f^*(T(\mathbf{x}))],$$

where  $f^*(t) = \sup_{u \in \text{dom}_f} (ut - f(u))$  – Fenchel conjugate.

## Recap of previous lecture

Let take some pretrained image classification model to get the conditional label distribution  $p(y|\mathbf{x})$  (e.g. ImageNet classifier).

### Evaluation of likelihood-free models

- ▶ Sharpness  $\Rightarrow$  low  $H(y|\mathbf{x}) = -\sum_y \int_{\mathbf{x}} p(y, \mathbf{x}) \log p(y|\mathbf{x}) d\mathbf{x}$ .
- ▶ Diversity  $\Rightarrow$  high  $H(y) = -\sum_y p(y) \log p(y)$ .

### Inception Score

$$IS = \exp(H(y) - H(y|\mathbf{x})) = \exp(\mathbb{E}_{\mathbf{x}} KL(p(y|\mathbf{x}) || p(y)))$$

### Frechet Inception Distance

$$D^2(\pi, p) = \|\mathbf{m}_\pi - \mathbf{m}_p\|_2^2 + \text{Tr} \left( \mathbf{C}_\pi + \mathbf{C}_p - 2\sqrt{\mathbf{C}_\pi \mathbf{C}_p} \right).$$

FID is related to moment matching.

---

Salimans T. et al. *Improved Techniques for Training GANs*, 2016

Heusel M. et al. *GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium*, 2017

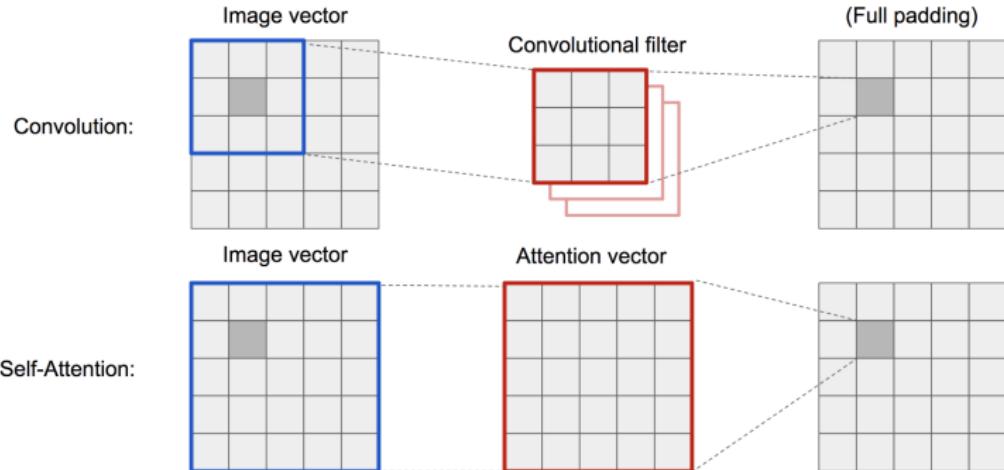
# Evolution of GANs



- ▶ **Vanilla GAN** <https://arxiv.org/abs/1406.2661>
- ▶ **DCGAN** <https://arxiv.org/abs/1511.06434>
- ▶ **CoGAN** <https://arxiv.org/abs/1606.07536>
- ▶ **ProGAN** <https://arxiv.org/abs/1710.10196>
- ▶ **StyleGAN** <https://arxiv.org/abs/1812.04948>

# Self-Attention GAN

- ▶ Convolutional layers process the information in a local neighborhood.
- ▶ Using convolutional layers alone is computationally inefficient for modeling long-range dependencies in images.

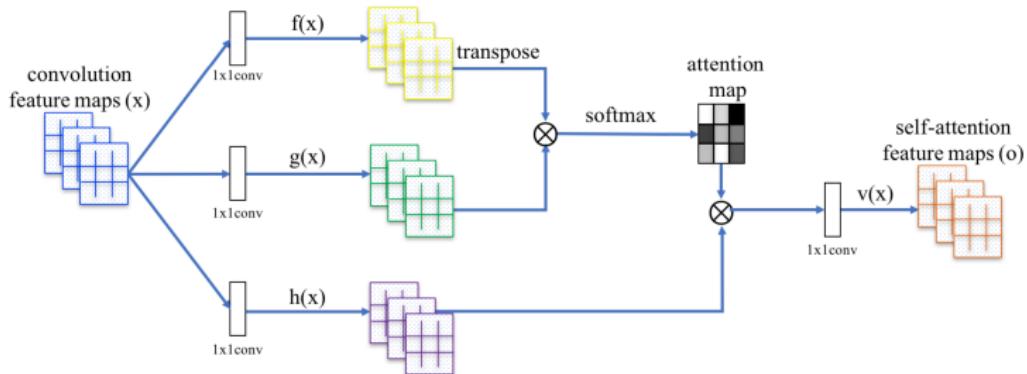


---

*image credit:*

<https://lilianweng.github.io/lil-log/2018/06/24/attention-attention.html>

# Self-Attention GAN



- ▶  $x$  – feature vector for one feature location.
- ▶  $N$  – number of feature locations.

$$\mathbf{f}(\mathbf{x}) = \mathbf{W}_f \mathbf{x}, \quad \mathbf{g}\mathbf{x} = \mathbf{W}_g \mathbf{x}, \quad \mathbf{h}\mathbf{x} = \mathbf{W}_h \mathbf{x}, \quad \mathbf{v}\mathbf{x} = \mathbf{W}_v \mathbf{x}$$

$$s_{ij} = \mathbf{f}(\mathbf{x}_i)^T \mathbf{g}(\mathbf{x}_j), \quad a_{ij} = \frac{\exp s_{ij}}{\sum_{i=1}^N \exp s_{ij}}, \quad \mathbf{o}_j = \mathbf{v} \left( \sum_{i=1}^N a_{ij} \mathbf{h}(\mathbf{x}_i) \right)$$

# Self-Attention GAN

## Technical Details

- ▶ Hinge loss for training.
- ▶ SpectralNorm in both the generator and the discriminator.
- ▶ Separate learning rates for the generator and the discriminator.

Model	Inception Score	Intra FID	FID
AC-GAN ( <a href="#">Odena et al., 2017</a> )	28.5	260.0	/
SNGAN-projection ( <a href="#">Miyato &amp; Koyama, 2018</a> )	36.8	92.4	27.62*
SAGAN	<b>52.52</b>	<b>83.7</b>	<b>18.65</b>

## Visualization of attention maps



# BigGAN

## Technical Details

- ▶ Hinge loss.
- ▶ Self-Attention GAN baseline.
- ▶ **Orthogonal regularization**

$$\|\mathbf{W}^T \mathbf{W} - \mathbf{I}\|^2 \Rightarrow \|\mathbf{W}^T \mathbf{W} - \text{diag}(\mathbf{W}^T \mathbf{W})\|^2$$

- ▶ **Truncation trick.** Components of  $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$  which fall outside a predefined range are resampled.

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

# BigGAN

Samples (512x512)



Interpolations



---

Brock A., Donahue J., Simonyan K. *Large Scale GAN Training for High Fidelity Natural Image Synthesis*, 2018

# Progressive Growing GAN

## Problems with HR image generation

- ▶ Disjoint manifolds  $\Rightarrow$  gradient problem.
- ▶ Small minibatch  $\Rightarrow$  training instability.

## Samples (1024x1024)



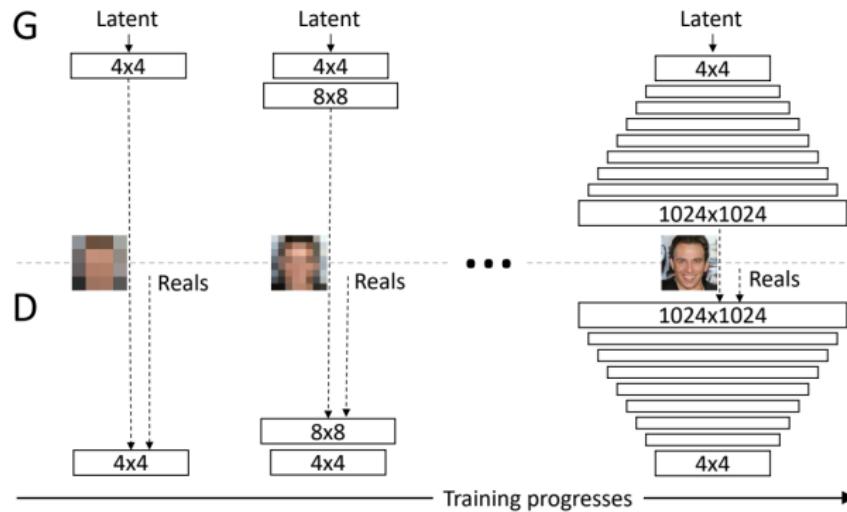
---

Karras T. et al. *Progressive Growing of GANs for Improved Quality, Stability, and Variation*, 2017

# Progressive Growing GAN

Grow both the generator and discriminator progressively, new layers will introduce higher-resolution details as the training progresses.

- ▶ Train GAN which generate 4x4 images (2 convs for G and D).
- ▶ Add upsampling layers to G, downsampling layers to D.
- ▶ Train GAN which generate 8x8 images.
- ▶ etc.



# StyleGAN

- ▶ Generating of HR images is hard.
- ▶ Progressive growing greatly simplifies the task.
- ▶ The ability to control specific features of the generated image is very limited.

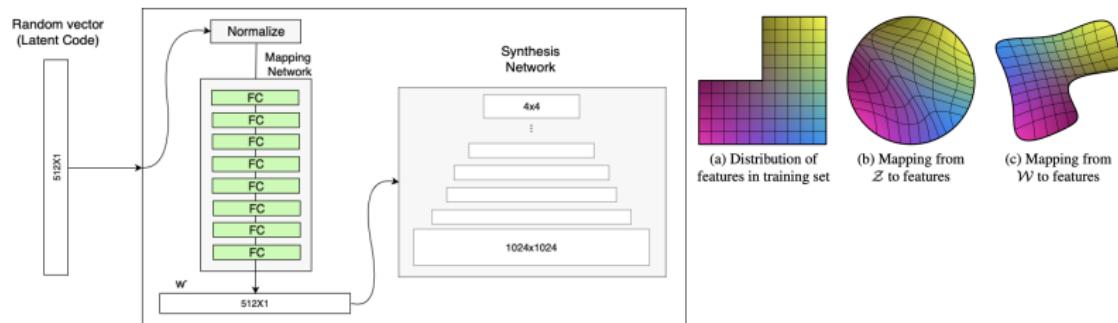
## Face image features

- ▶ Coarse (pose, general hair style, face shape). Resolution  $4^2 - 8^2$ .
- ▶ Middle (finer facial features, hair style, eyes open/closed). Resolution  $16^2 - 32^2$ .
- ▶ Fine (color scheme (eye, hair and skin) and micro features). Resolution  $64^2 - 1024^2$ .

# StyleGAN

## Mapping Network

- ▶ Generator input is likely to be **disentangled**. Each component of input vector  $\mathbf{z}$  should be responsible for one generative factor.
- ▶ Mapping network  $f : \mathcal{Z} \rightarrow \mathcal{W}$  is used to reduce correlations between components of  $\mathbf{z}$ .

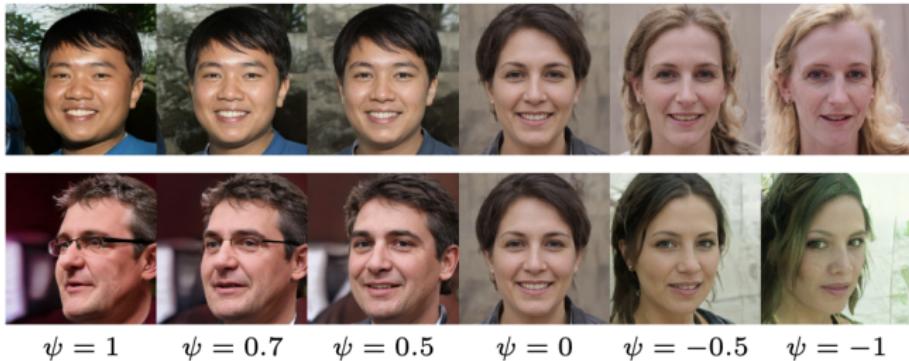


# StyleGAN

## Truncation trick

$$\mathbf{w}' = \hat{\mathbf{w}} + \psi \cdot (\mathbf{w} - \hat{\mathbf{w}}), \quad \hat{\mathbf{w}} = \mathbb{E}_{\mathbf{z}} p(f(\mathbf{z}))$$

- ▶ Constant  $\psi$  is a tradeoff between diversity and fidelity.
- ▶  $\psi = 0.7$  is used for most of the results.
- ▶ Truncation is done only at the low-resolution layers.



# StyleGAN

## Results

Method	CelebA-HQ	FFHQ
A Baseline Progressive GAN [30]	7.79	8.04
B + Tuning (incl. bilinear up/down)	6.11	5.25
C + Add mapping and styles	5.34	4.85
D + Remove traditional input	5.07	4.88
E + Add noise inputs	<b>5.06</b>	4.42
F + Mixing regularization	5.17	<b>4.40</b>

## Samples (1024x1024)



Karras T., Laine S., Aila T. A Style-Based Generator Architecture for Generative Adversarial Networks, 2018

# Summary