# Generative Adversarial Network
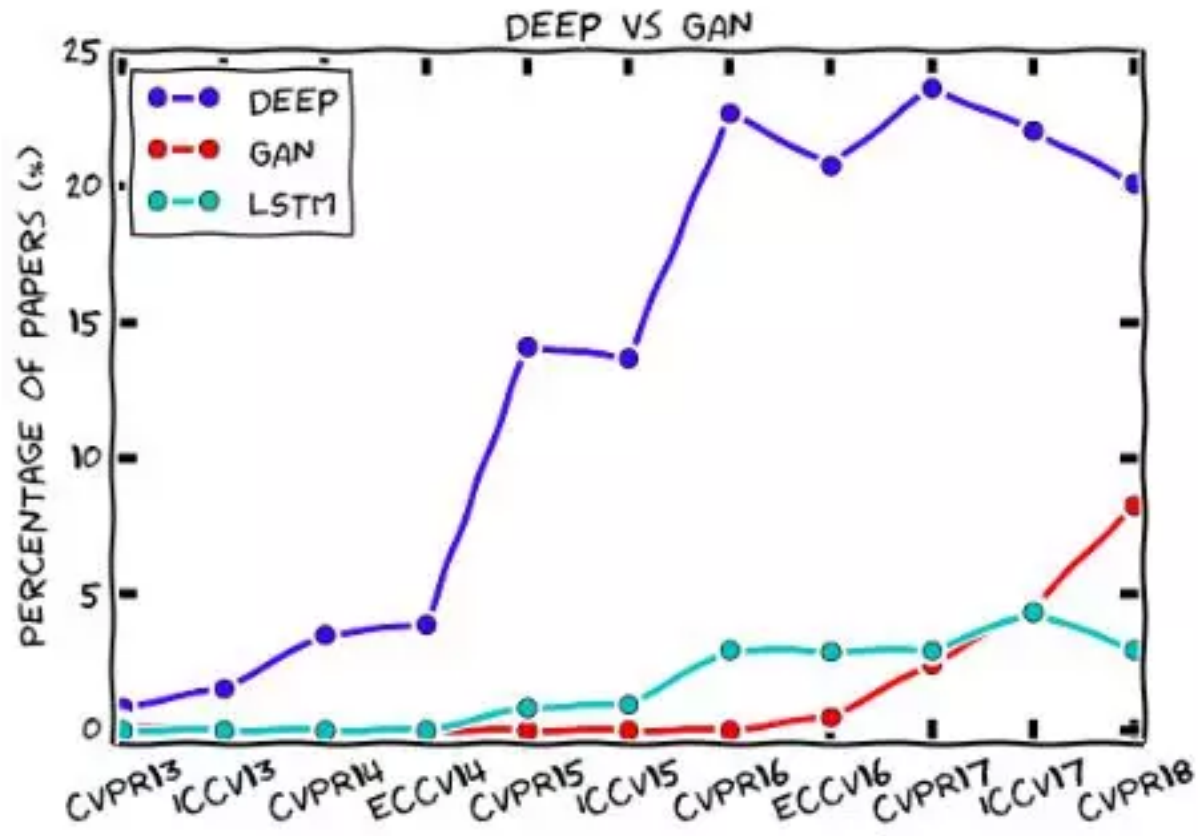
## Liangjie

# Outline

- Basic theory (GAN, DCGAN)

- Recent improvement (WGAN, AAE, WAE)
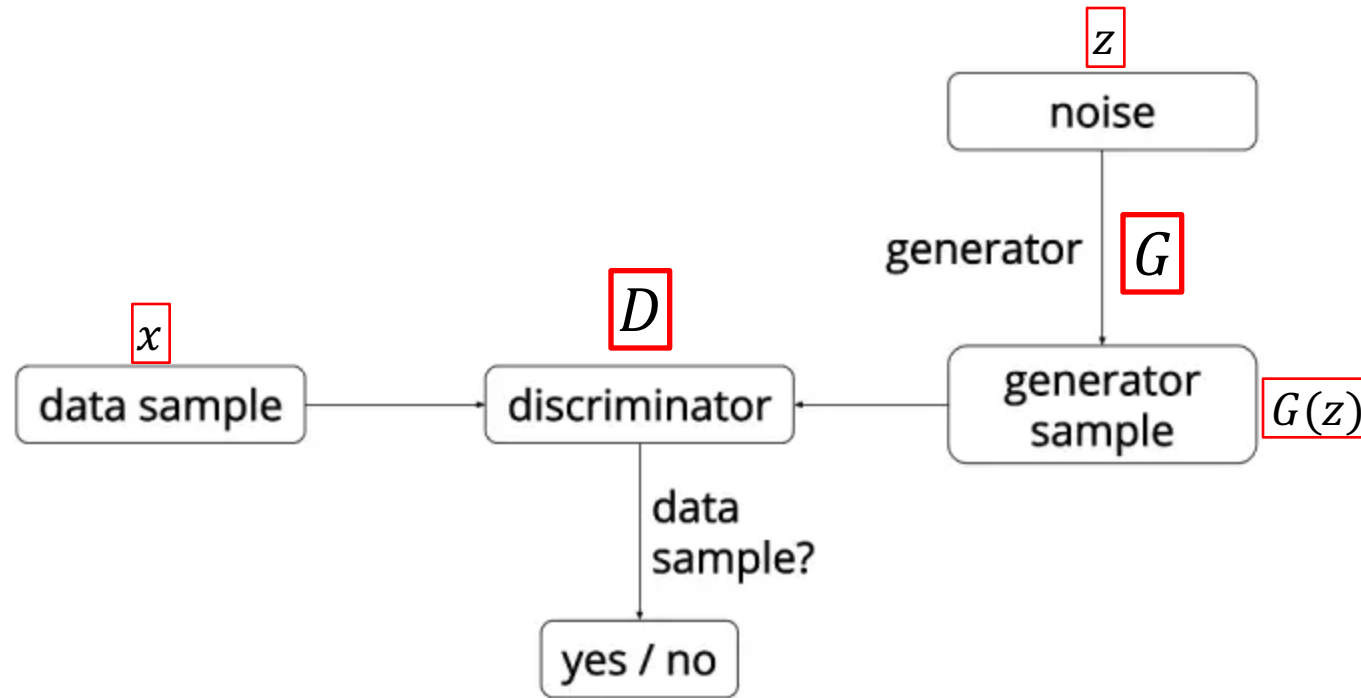
- Applications (CGAN, Text2Img, DeblurGAN)

# A hot topic



GAN related papers per community in CVPR 2018:

- style transfer/cycle GAN/domain adaption: **13**
- Photo enhancement/deblur/high-resolution reconstruction/..: **7**
- Optimizing GAN theory: **6**
- Image synthesis: **10**
- human face related: **7**
- human pose related: **4**
- Person Re-ID: **3**
- Others …

- **TOTAL**: **75+, around 8% of 979 papers.**

# The idea of GAN



Competition in this game drives both teams (**G** and **D**) to improve their capacity until the counterfeits are indistinguishable from the genuine articles.

I. Goodfellow, et.al. Generative Adversarial Nets, NIPS2014.

# The Algorithm

$$\min_{G} \max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} \left[ \log \boxed{D(\boldsymbol{x})} \right] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z})))].$$

The probability that $x$ is a real sample.

**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 $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.
        • Sample minibatch of $m$ examples $\{\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\boldsymbol{x})$.
        • Update the discriminator by ascending its stochastic gradient:

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

    **end for**
    • Sample minibatch of $m$ noise samples $\{\boldsymbol{z}^{(1)}, \ldots, \boldsymbol{z}^{(m)}\}$ from noise prior $p_g(\boldsymbol{z})$.
    • Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right).$$
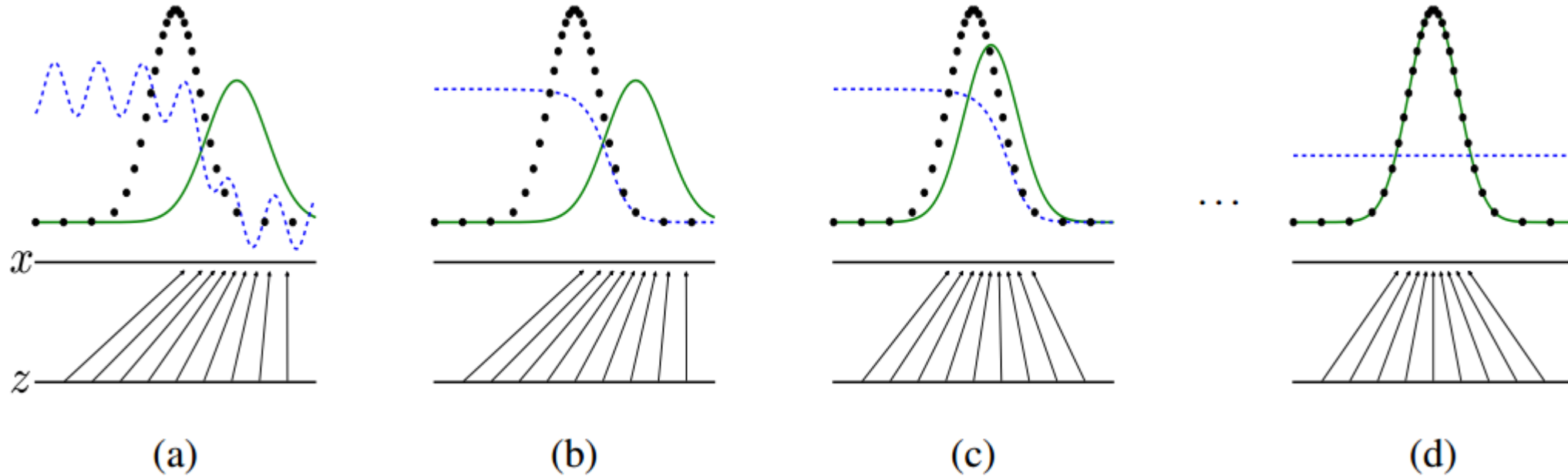
**end for**
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

First train **D** for several times,
Then train **G**.

k is hard to control!

# The training process of the original GAN



(a)

(b)

(c)

(d)

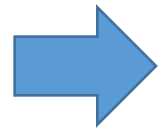The original state. Both G and D are not powerful.

After training D. D can tell which sample is real.

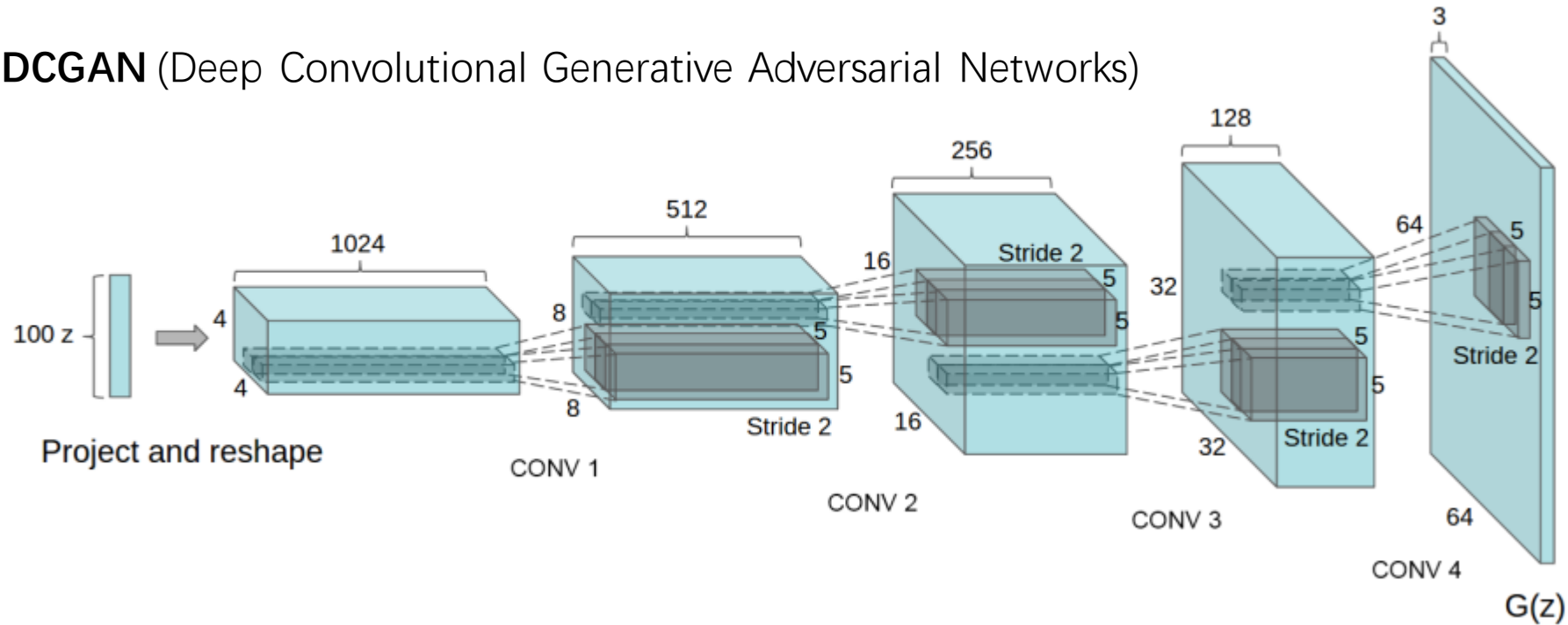After training G. G can generate 'real' samples.

The final state. Both G and D are powerful.

# The weakness

1) The structure is too simple thus lacks ability.
2) The generative model and the discriminative model should utilize the deep convolutional network.

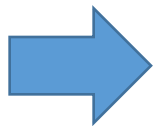**DCGAN** (Deep Convolutional Generative Adversarial Networks)



A. Radford, et.al. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks, ICLR2016.

# Still have weakness

1. The training process is unstable.

2. Hardly to control the ability of the **G** and **D**.

3. The gradient always disappears.

4. Do not have an index to show the performance of the model.

5. The obtained samples lack of diversity.

   ➡️  **WGAN** (Wasserstein GAN)

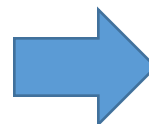Martin. Arjovsky, et.al. Wasserstein GAN. arxiv.org/abs/1701.07875

# WGAN (Wasserstein GAN)

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

---

**Require:** : $\alpha$, the learning rate. $c$, the clipping parameter. $m$, the batch size. $n_{\text{critic}}$, the number of iterations of the critic per generator iteration.

**Require:** : $w_0$, initial critic parameters. $\theta_0$, initial generator's parameters.

1: **while** $\theta$ has not converged **do**
2:      **for** $t = 0, ..., n_{\text{critic}}$ **do**
3:         Sample $\{x^{(i)}\}_{i=1}^{m} \sim \mathbb{P}_r$ a batch from the real data.
4:         Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a batch of prior samples.
5:         $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^{m} f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^{m} f_w(g_\theta(z^{(i)})) \right]$
6:         $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$
7:         $w \leftarrow \text{clip}(w, -c, c)$
8:      **end for**
9:      Sample $\{z^{(i)}\}_{i=1}^{m} \sim p(z)$ a batch of prior samples.
10:      $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^{m} f_w(g_\theta(z^{(i)}))$
11:      $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$
12: **end while**

**The modifications:**

1. No sigmoid in last layer of **D**.

2. No $\log(a)$ in the loss of the **G** and **D**.

3. For **D**, clip the updated parameter to [-c , c].

Martin Arjovsky, et.al. Wasserstein GAN. arxiv.org/abs/1701.07875

# The problem in GAN
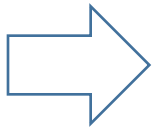
1. If the **D** has strong ability to make decision, the gradient of **G** will disappear

When training the **D**: we maximize the equation:

$$L(D, g_\theta) = \mathbb{E}_{x \sim P_r} \log D(x) + \mathbb{E}_{x \sim P_g} \log[1 - D(x)]$$

$$= P_r(x) \log D(x) + P_g(x) \log[1 - D(x)]$$

The derivative of $D(x)$ equals 0:

$$\frac{P_r(x)}{D(x)} - \frac{P_g(x)}{1 - D(x)} = 0.$$

The optimal discriminator :

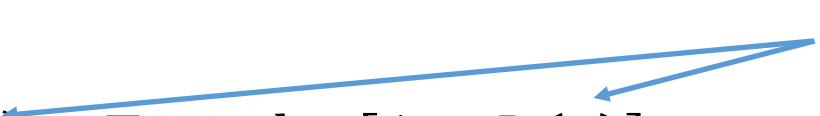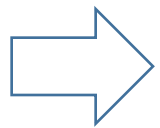$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)}$$

M. Arjovsky, et.al. Wasserstein GAN. arxiv.org/abs/1701.07875

# The problem in GAN

When training the **G**: we minimize the equation:

$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)}$$

$$\mathbb{E}_{x \sim P_r} \log D(x) + \mathbb{E}_{x \sim P_g} \log[1 - D(x)]$$

$$\Rightarrow \quad \mathbb{E}_{x \sim P_r} \log \frac{P_r(x)}{\frac{1}{2}(P_r(x) + P_g(x))} + \mathbb{E}_{x \sim P_g} \log[1 - \frac{P_r(x)}{\frac{1}{2}(P_r(x) + P_g(x))}] - 2\log 2$$

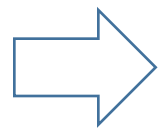M. Arjovsky, et.al. Towards Principled Methods for Training Generative Adversarial Networks. ICLR 2017

# The problem in GAN

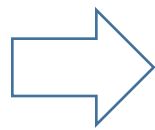1. If the **D** has strong ability to make decision, the gradient of **G** will disappear

KL divergence :
$$KL(P_1||P_2) = \mathbb{E}_{x \sim P_r} \log \frac{P_1}{P_2}$$

JS divergence :
$$JS(P_1||P_2) = \frac{1}{2} KL(P_1||(P_1 + P_2)/2) + \frac{1}{2} KL(P_2||(P_1 + P_2)/2)$$

$$\mathbb{E}_{x \sim P_r} \log \frac{P_r(x)}{\frac{1}{2}(P_r(x) + P_g(x))} + \mathbb{E}_{x \sim P_g} \log[1 - \frac{P_r(x)}{\frac{1}{2}(P_r(x) + P_g(x))}] - 2\log 2$$
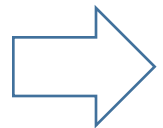
$$\implies = 2JS(P_r||P_g) - 2log2 \implies$$

When training the **G**, we minimize the JS divergence between real distribution and generative distribution.

M. Arjovsky, et.al. Towards Principled Methods for Training Generative Adversarial Networks. ICLR 2017

# The problem in GAN

1. If the **D** has strong ability to make decision, the gradient of **G** will disappear

$$JS(P_r||P_g) = \begin{cases} log2 & P_r = 0, P_g \neq 0; \; or \; P_r \neq 0, P_g = 0 \\ 0 & P_r = 0, P_g = 0; \; or \; \boxed{P_r \neq 0, P_g \neq 0} \end{cases}$$

The two distribution can hardly have combinations since $z$ is obtained randomly and the intrinsic dimension is much lower than the real space.

➡ The gradient of **G** disappears.

M. Arjovsky, et.al. Towards Principled Methods for Training Generative Adversarial Networks. ICLR 2017

# The problem in GAN

The training process is unstable.

For **G**, it minimizes the object function:

$$\mathbb{E}_{x \sim P_g}[-\log D^*(x)] = KL(P_r||P_g) - \mathbb{E}_{x \sim P_g}\log[1 - D^*(x)]$$

$$= KL(P_r||P_g) - 2JS(P_r||P_g) + 2\log 2 + \mathbb{E}_{x \sim P_r}\log D^*(x)$$

$$\Longrightarrow \quad KL(P_r||P_g) - 2JS(P_r||P_g)$$

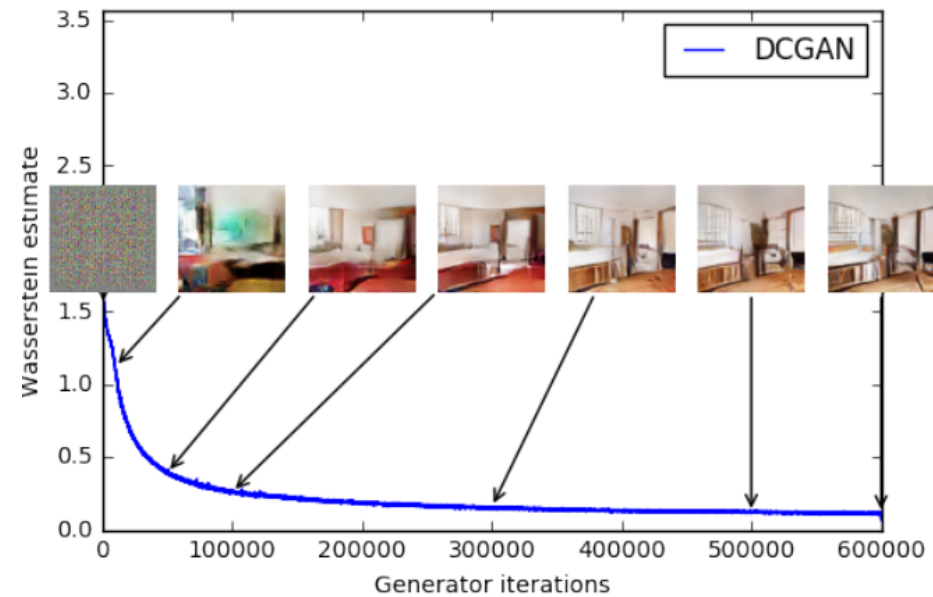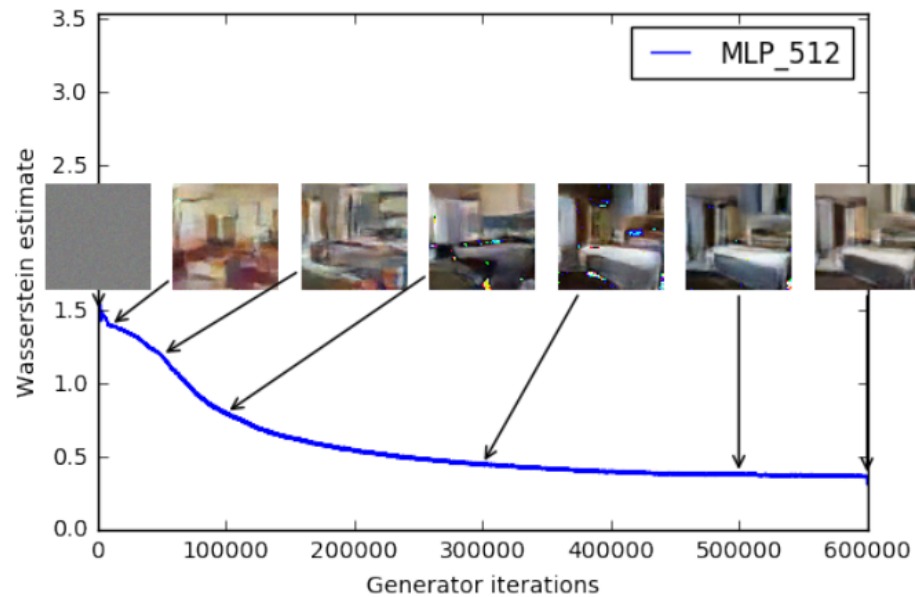One to close while another to far! $\Longrightarrow$ unstable

M. Arjovsky, et.al. Towards Principled Methods for Training Generative Adversarial Networks. ICLR 2017

# The problem in GAN

Wasserstein Distance (Earth-Mover distance):

$$W(P_r, P_g) = \inf_{\gamma \sim \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma}[||x - y||]$$



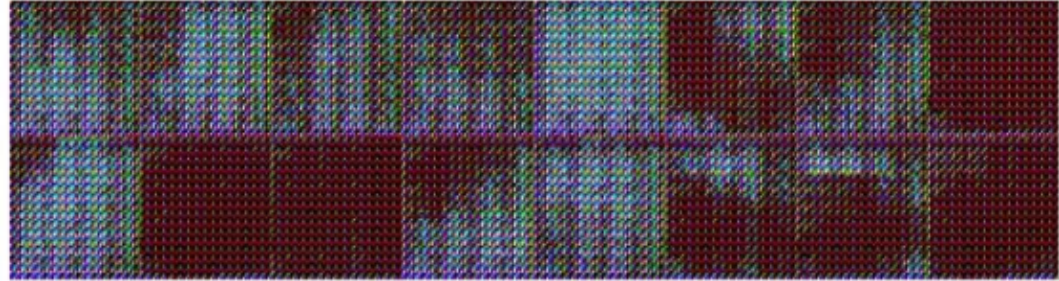M. Arjovsky, et.al. Wasserstein GAN. arxiv.org/abs/1701.07875

# Wasserstein GAN

**Advantages:**

- does not need to carefully control the capacity of G and D

- alleviates the mode collapse problem

- provides a metric of model capacity during training

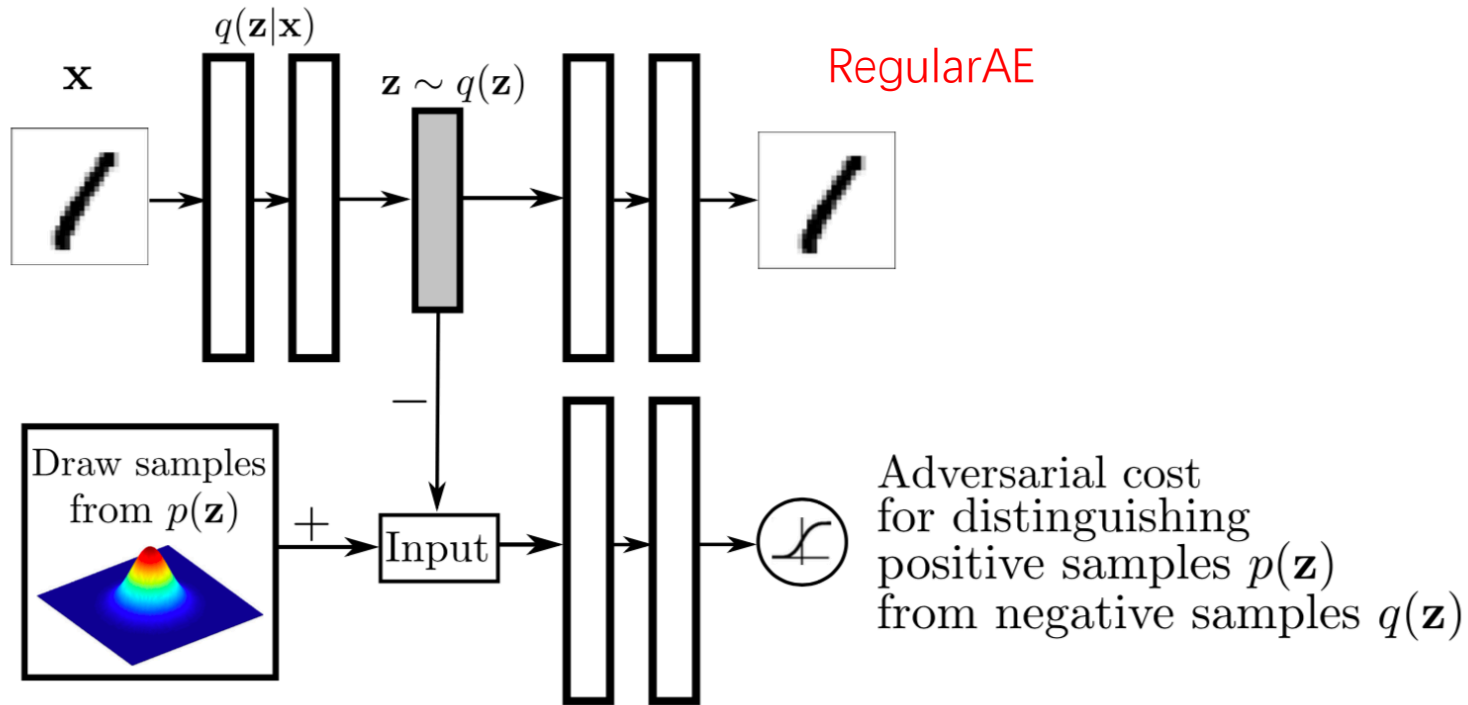- does not require well-designed network architectures

# Some Results

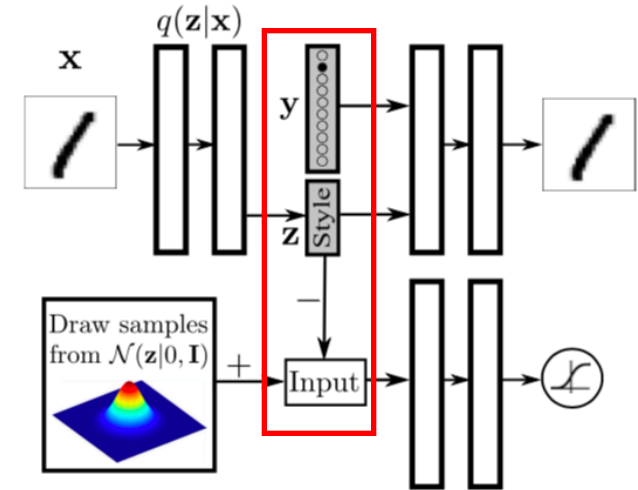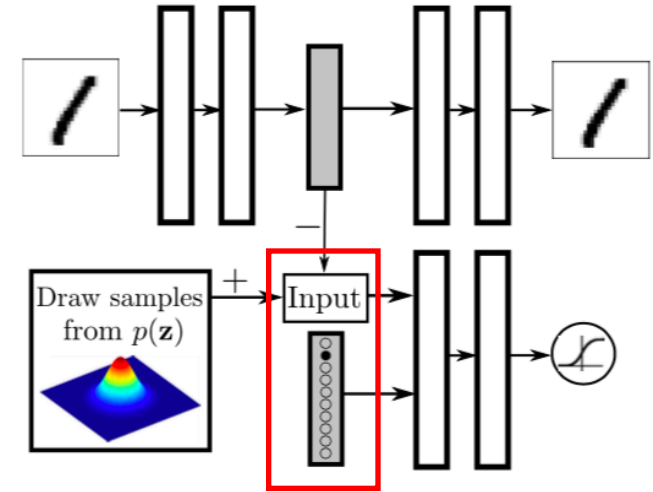## WGANs vs DCGANs (without BN)



## WGANs vs GANs (without CNN)



M. Arjovsky, et.al. Wasserstein GAN. arxiv.org/abs/1701.07875

# Adversarial Auto-Encoders



RegularAE

$q(\mathbf{z}|\mathbf{x})$

$\mathbf{x}$

$\mathbf{z} \sim q(\mathbf{z})$

Draw samples from $p(\mathbf{z})$

Input

Adversarial cost for distinguishing positive samples $p(\mathbf{z})$ from negative samples $q(\mathbf{z})$

D：z is derived from the encoder of the random sampled from prior distribution

Finally: $q_z$ has same distribution with $p_z$, we can feed the random sampled latent code to the decoder, to generate new samples.

Draw samples from $p(\mathbf{z})$

Input

$q(\mathbf{z}|\mathbf{x})$

$\mathbf{x}$

y

Style

z

Draw samples from $\mathcal{N}(\mathbf{z}|0, \mathbf{I})$

Input

A Makhzani, et.al. Adversarial Auto-encoders. arXiv preprint arXiv:1511.05644

# Adversarial Auto-Encoders



(a) MNIST

(b) SVHN

A Makhzani, et.al. Adversarial Auto-encoders. arXiv preprint arXiv:1511.05644

# Wasserstein Auto-Encoders

**Encoder:**
- Based on a specified divergence, <span style="color:red">matches</span> the encoded distribution $Q_Z$ of training samples to the prior $P_Z$.

- Ensures that the latent codes fed to the decoders are <span style="color:red">informative to reconstruct</span> the training sample.

**Decoder:**
- According to the cost function, <span style="color:red">reconstruct</span> the encoded training sample.
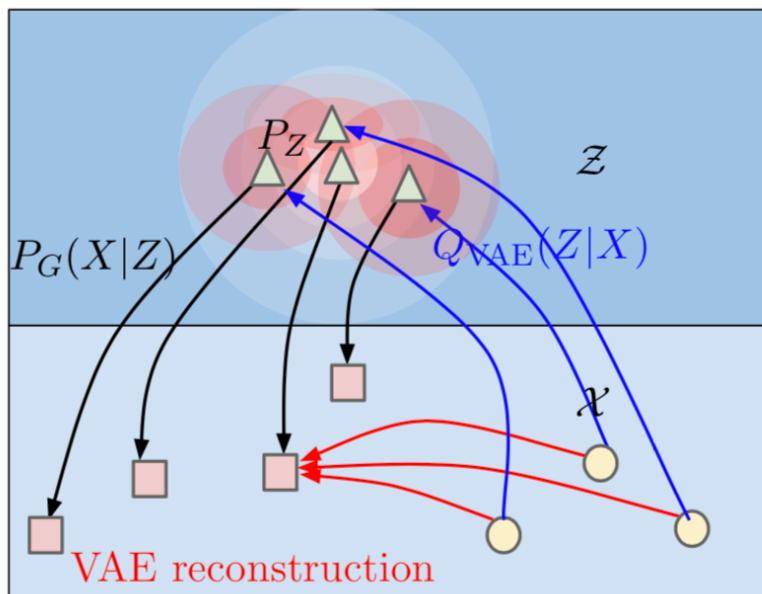
**Loss Function:** Optimal transport cost $W_c(P_X, P_G)$, a family of Wasserstein distance

$$D_{\mathrm{WAE}}(P_X, P_G) := \inf_{Q(Z|X) \in \mathcal{Q}} \mathbb{E}_{P_X} \mathbb{E}_{Q(Z|X)} \left[ \boxed{c(X, G(Z))} \right] + \lambda \cdot \boxed{\mathcal{D}_Z(Q_Z, P_Z)},$$

<div style="color:red">Reconstruction loss        Penalty</div>

Ilya Tolstikhin, et.al. Wasserstein Auto-Encoders. ICLR 2018 (8.0)

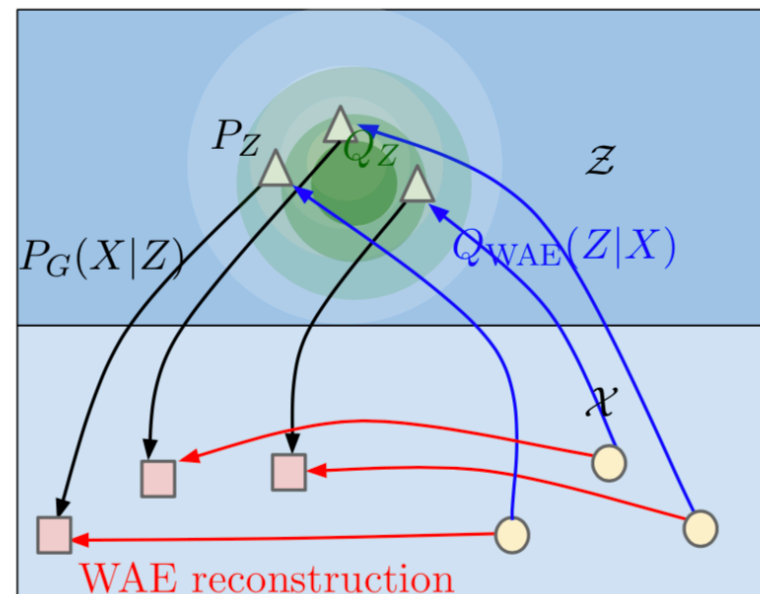# Wasserstein Auto-Encoders



For each samples in $X$ (yellow circle), forces $Q(Z|X)$ (triangle) to match $P(Z)$ (white shape).

Forces $Q_z = \int Q(Z|X) dP_X$ to match $P(Z)$ (green ball).

Ilya Tolstikhin, et.al. Wasserstein Auto-Encoders. ICLR 2018

# Wasserstein Auto-Encoders

**Algorithm 1** Wasserstein Auto-Encoder with GAN-based penalty (WAE-GAN).

**Require:** Regularization coefficient $\lambda > 0$.

   Initialize the parameters of the encoder $Q_\phi$, decoder $G_\theta$, and latent discriminator $D_\gamma$.

   **while** $(\phi, \theta)$ not converged **do**

      Sample $\{x_1, \ldots, x_n\}$ from the training set

      Sample $\{z_1, \ldots, z_n\}$ from the prior $P_Z$

      Sample $\tilde{z}_i$ from $Q_\phi(Z|x_i)$ for $i = 1, \ldots, n$

      Update $D_\gamma$ by ascending:

$$\frac{\lambda}{n} \sum_{i=1}^{n} \log D_\gamma(z_i) + \log\big(1 - D_\gamma(\tilde{z}_i)\big)$$

      Update $Q_\phi$ and $G_\theta$ by descending:

$$\frac{1}{n} \sum_{i=1}^{n} c\big(x_i, G_\theta(\tilde{z}_i)\big) - \lambda \cdot \log D_\gamma(\tilde{z}_i)$$

**end while**

---

**Algorithm 2** Wasserstein Auto-Encoder with MMD-based penalty (WAE-MMD).

**Require:** Regularization coefficient $\lambda > 0$, characteristic positive-definite kernel $k$.

   Initialize the parameters of the encoder $Q_\phi$, decoder $G_\theta$, and latent discriminator $D_\gamma$.

   **while** $(\phi, \theta)$ not converged **do**

      Sample $\{x_1, \ldots, x_n\}$ from the training set
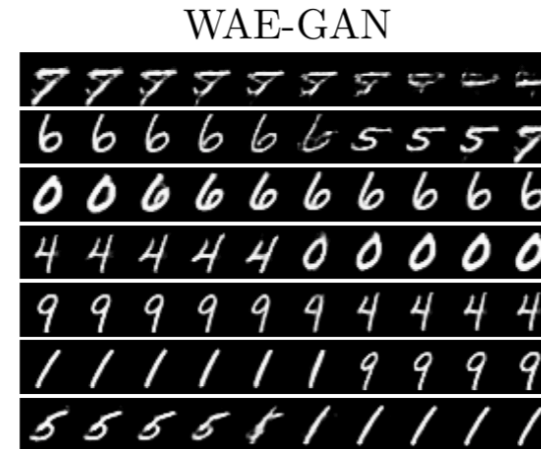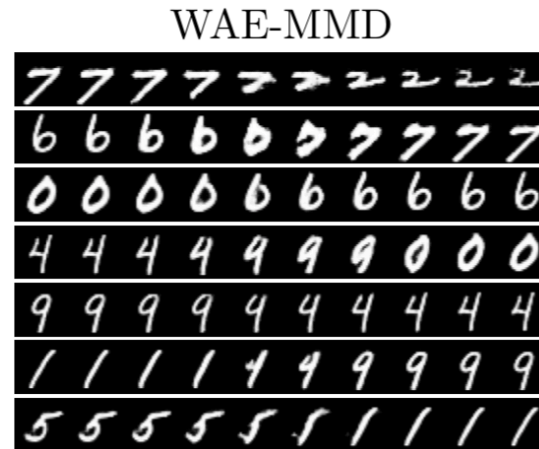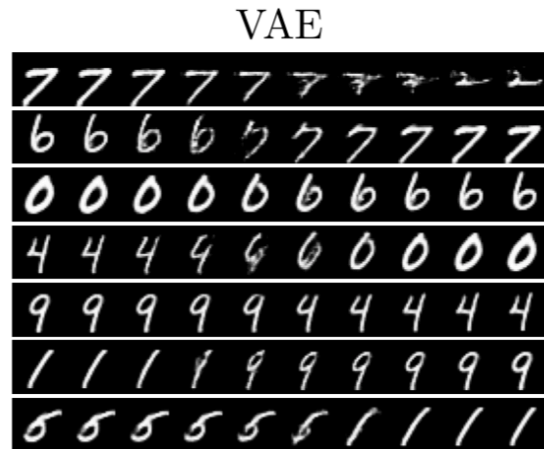
      Sample $\{z_1, \ldots, z_n\}$ from the prior $P_Z$

      Sample $\tilde{z}_i$ from $Q_\phi(Z|x_i)$ for $i = 1, \ldots, n$
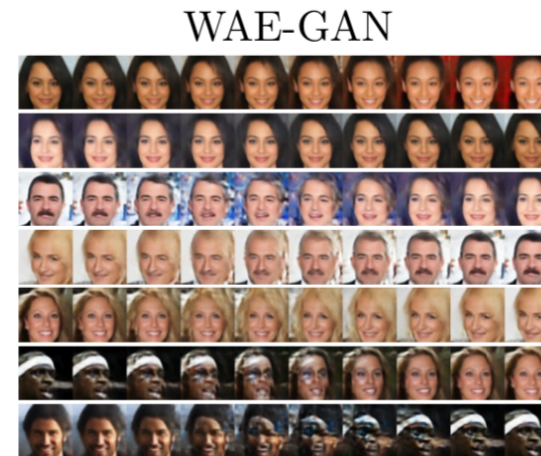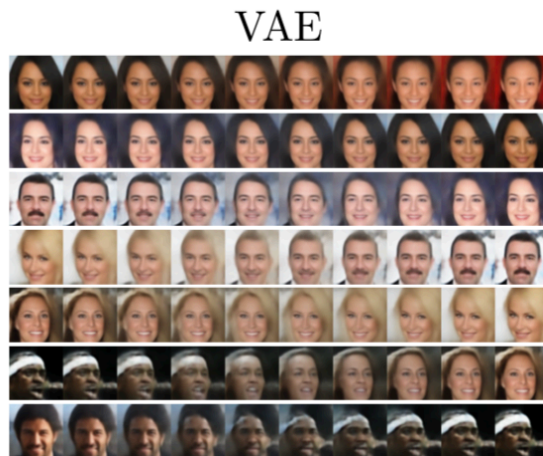
      Update $Q_\phi$ and $G_\theta$ by descending:

$$\frac{1}{n} \sum_{i=1}^{n} c\big(x_i, G_\theta(\tilde{z}_i)\big) + \frac{\lambda}{n(n-1)} \sum_{\ell \neq j} k(z_\ell, z_j)$$
$$+ \frac{\lambda}{n(n-1)} \sum_{\ell \neq j} k(\tilde{z}_\ell, \tilde{z}_j) - \frac{2\lambda}{n^2} \sum_{\ell, j} k(z_\ell, \tilde{z}_j)$$

**end while**

Ilya Tolstikhin, et.al. Wasserstein Auto-Encoders. ICLR 2018

# Wasserstein Auto-Encoders



VAE   WAE-MMD   WAE-GAN

**MNIST:**
**28*28**

VAE   WAE-MMD   WAE-GAN

**celebA:**
**64*64**

Ilya Tolstikhin, et.al. Wasserstein Auto-Encoders. ICLR 2018
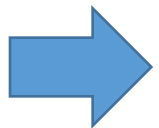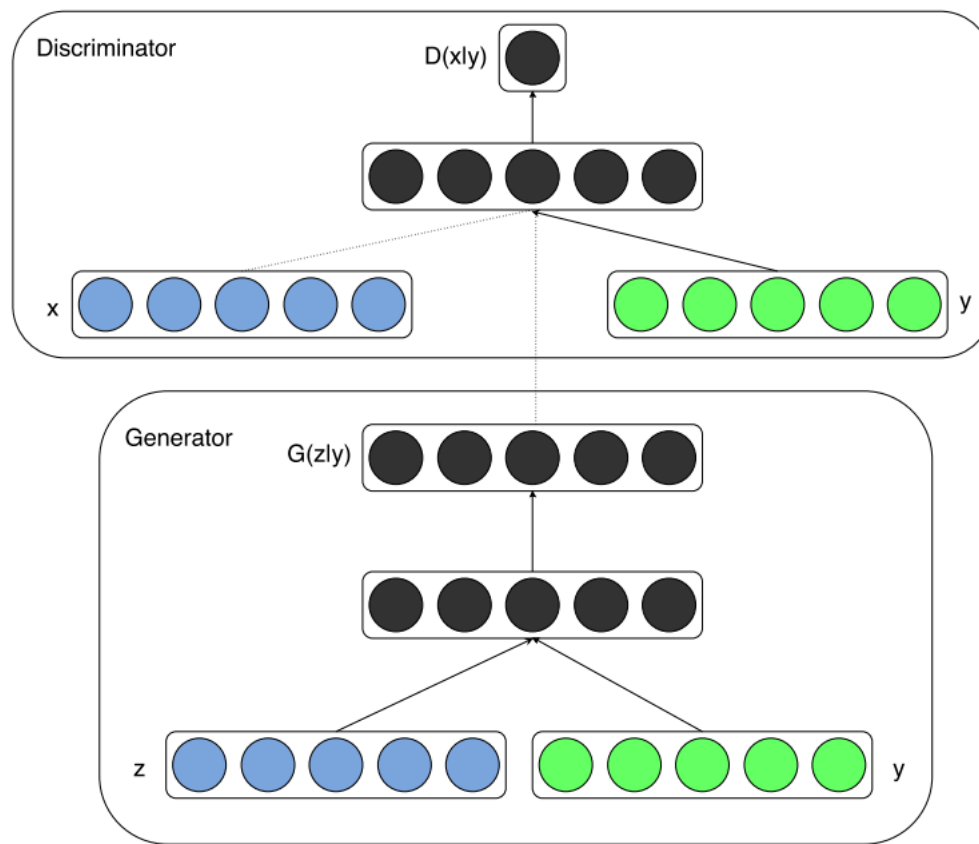
# Applications

1. Generate new samples for training, especially those hard to collect.

2. Combine with various low-level vision tasks such as segmentation, etc.

3. Complete the broken image (inpainting).

4. Generate high resolution image from lower one.

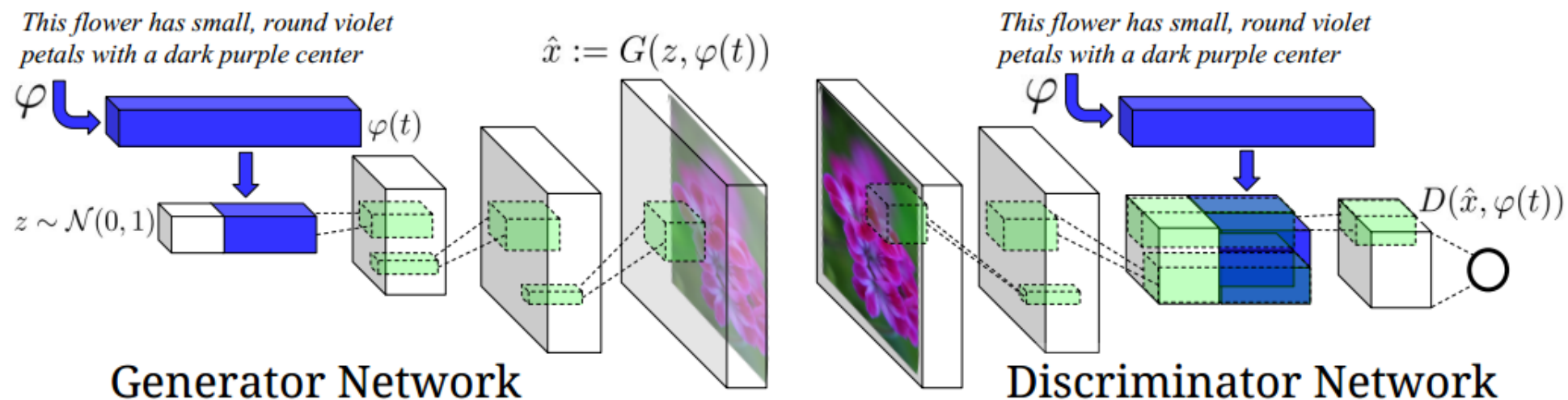5. Generate images from the text descriptions.

...

**GAN can be combined with any task that generates some new things, e.g., the mask in segmentation task, the broken part in image inpainting task, or the high-resolution image, etc.**

# Conditional GAN



$$\min_{G} \max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} [\log D(\boldsymbol{x}|\boldsymbol{y})] + \mathbb{E}_{\boldsymbol{z} \sim p_{z}(\boldsymbol{z})} [\log(1 - D(G(\boldsymbol{z}|\boldsymbol{y})))].$$

Mehdi Mirza, et.al. Conditional Generative Adversarial Nets. arXiv preprint arXiv:1411.1784

# Text2Img



Generator Network

Discriminator Network

$$\hat{x} := G(z, \varphi(t))$$

$$z \sim \mathcal{N}(0,1)$$

$$\varphi(t)$$

$$D(\hat{x}, \varphi(t))$$

This flower has small, round violet petals with a dark purple center

This flower has small, round violet petals with a dark purple center

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

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

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

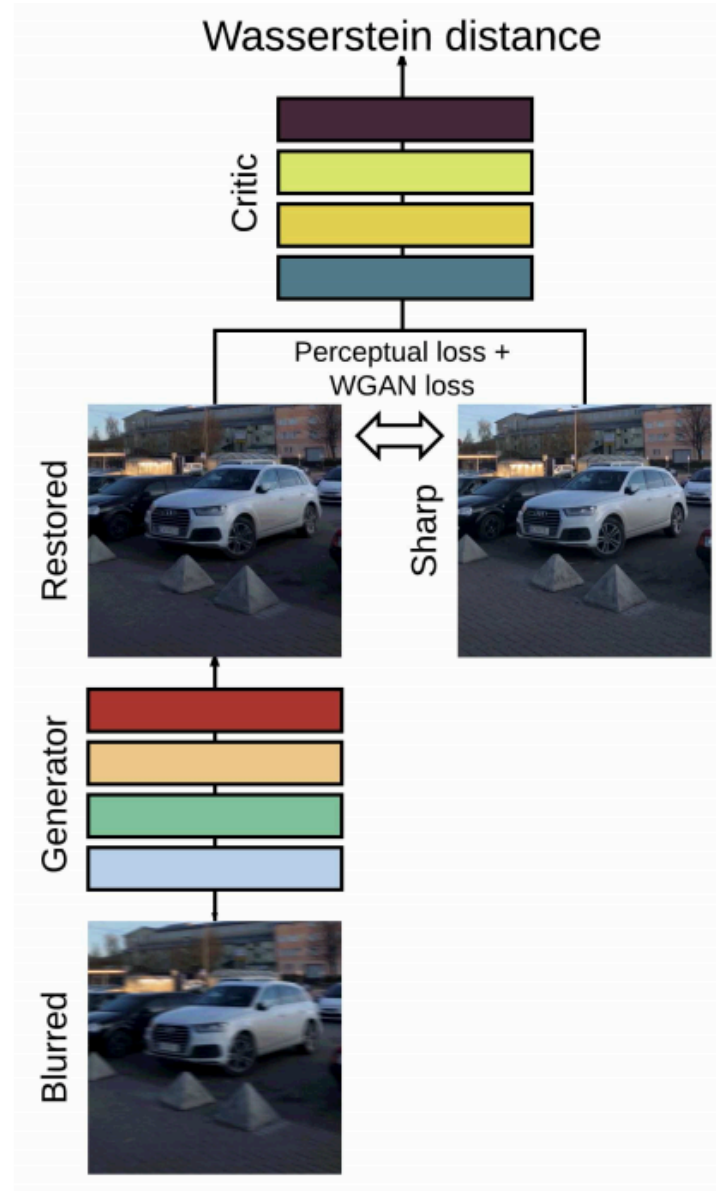this white and yellow flower have thin white petals and a round yellow stamen

S.Reed, et.al. Generative Adversarial Text to Image Synthesis. ICML2016

# DeblurGAN

**WGAN loss:**

Wasserstein distance

**Perceptual loss:**

the **difference** between the VGG-19 conv3.3 feature maps of the sharp and restored images.



Orest Kupyn, et.al. DeblurGAN: Blind Motion Deblurring Using Conditional Adversarial Networks. CVPR 2018

# Thanks