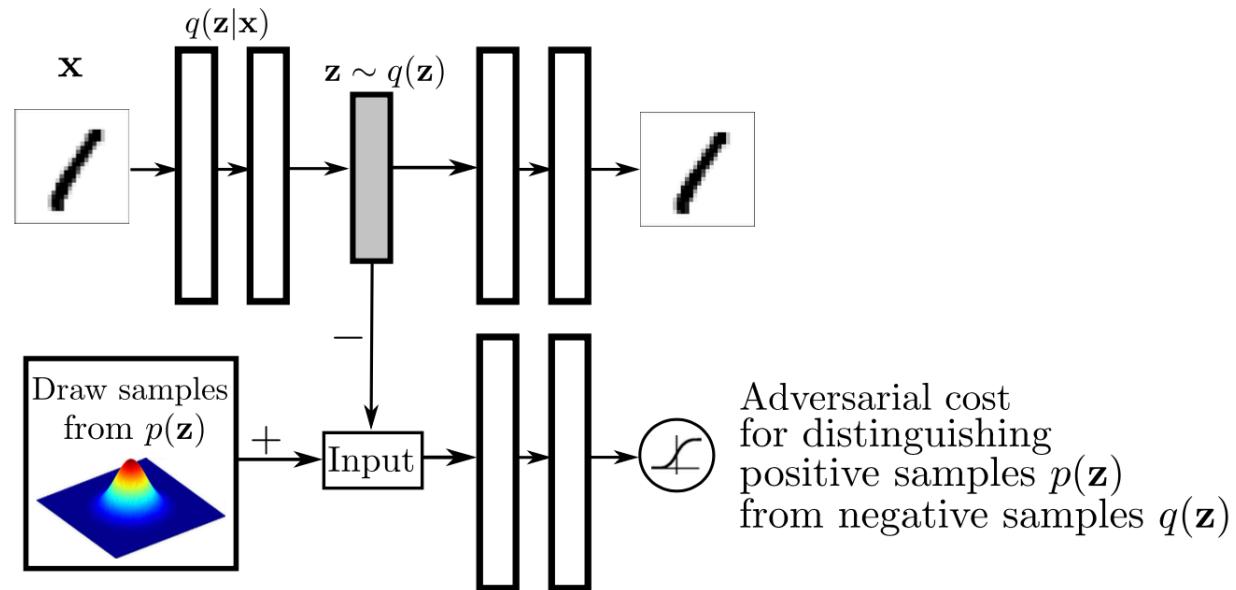
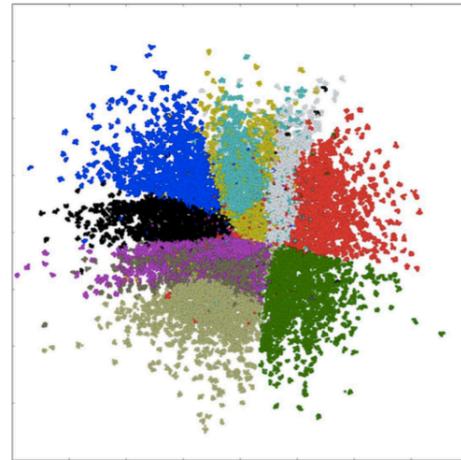


Part II

$$q(\mathbf{z}) = \int_{\mathbf{x}} q(\mathbf{z}|\mathbf{x}) p_{\text{data}}(\mathbf{x}) d\mathbf{x}$$

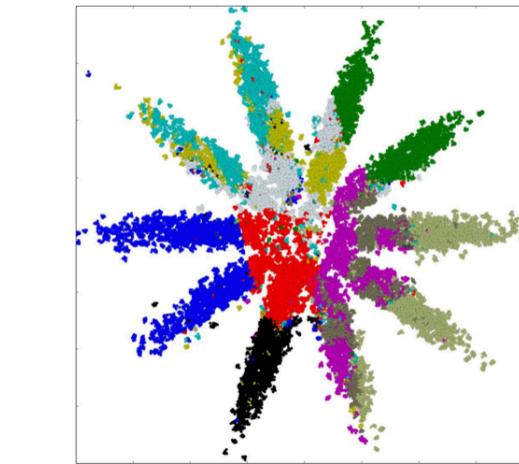


Code Space
of MNIST:



Gaussian Prior

Makhzani et al., 2015

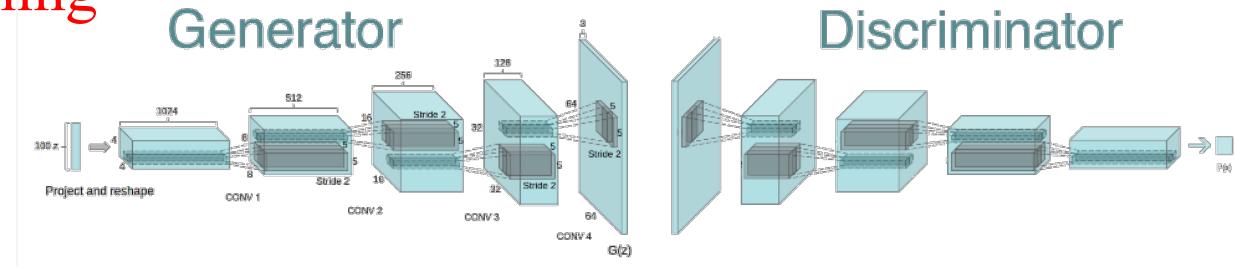


Mixture of Gaussians

Wasserstein GAN

Difficulty of Generative Adversarial Networks

- The gradient vanishing



$$J^{(D)} = -E_{x \sim P_r} [\log D(x)] - E_{x \sim P_g} [\log(1 - D(x))]$$

$$J^{(G)} = E_{x \sim P_g} [\log(1 - D(x))]$$

- At the very early phase of training, D is very easy to be confident in detecting G, so D will output almost always 0

In GAN, better discriminator leads to worse vanishing gradient in its generator!

Proof Sketch:

[Martin Arjovsky et al., Wasserstein GAN, 2017]

- Minimizing generator yields minimizing the JS divergence when the discriminator is optimal:

$$J^{(D)} = -E_{x \sim P_r} [\log D(x)] - E_{x \sim P_g} [\log(1 - D(x))]$$

- The optimal D for any P_r and P_g is always: $-P_r(x) \log D(x) - P_g(x) \log[1 - D(x)]$

$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)} \quad -\frac{P_r(x)}{D(x)} + \frac{P_g(x)}{1 - D(x)} = 0$$

- The generative loss (by adding a term independent of P_g) is:

$$J^{(G)} = E_{x \sim P_g} [\log(1 - D(x))]$$

Plug D^* into $J^{(G)}$:

Jensen–Shannon divergence

$$J^{(D^*)} = 2JS(P_r || P_g) - 2 \log 2$$

So, when D is *optimal*, minimizing the loss is equal to minimizing the *JS divergence*.

JS divergence:

$$\text{JS}(p \parallel q) = \frac{1}{2} \text{KL}(p \parallel m) + \frac{1}{2} \text{KL}(q \parallel m)$$

where

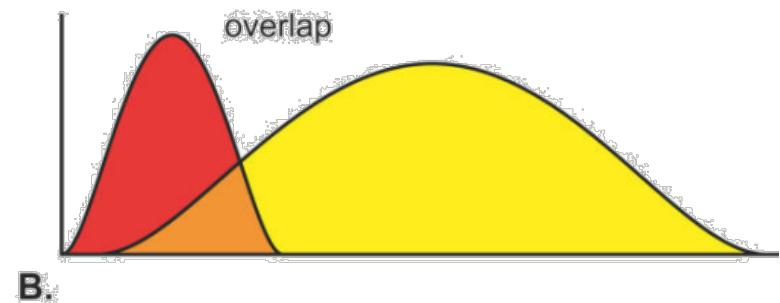
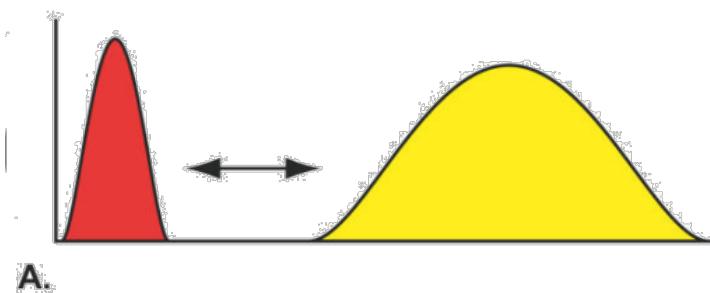
$$m(x) = \frac{1}{2} (p(x) + q(x))$$

$$\text{KL}(p \parallel q) = - \sum_x p(x) \log \left(\frac{q(x)}{p(x)} \right)$$

Proof Sketch:

$$J^{(D^*)} = 2JS(P_r || P_g) - 2 \log 2$$

- If the supports (underlying low-dimension manifolds) of P_r and P_g (almost) have no overlap, then $JS(P_r || P_g) = \log 2$, and thus the gradient of $J^{(G)}$ w.r.t. P_g vanishes.
- The probability that the support of P_r and P_g have almost zero overlap is 1.



Preliminaries: distance measures for distributions

- JS

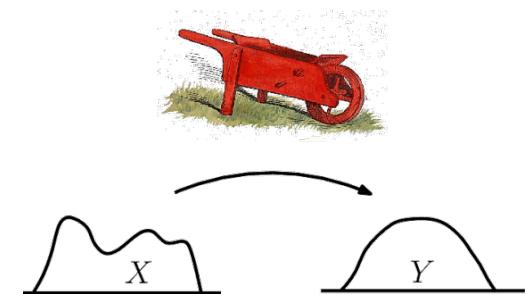
$$JS(P||Q) = \frac{1}{2}KL(P\| \frac{P+Q}{2}) + \frac{1}{2}KL(Q\| \frac{P+Q}{2})$$

- Wasserstein

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

the set of all joint γ distributions whose marginals are P and Q , respectively.

γ indicates a plan to transport “mass” from x to y , when performing P into Q .



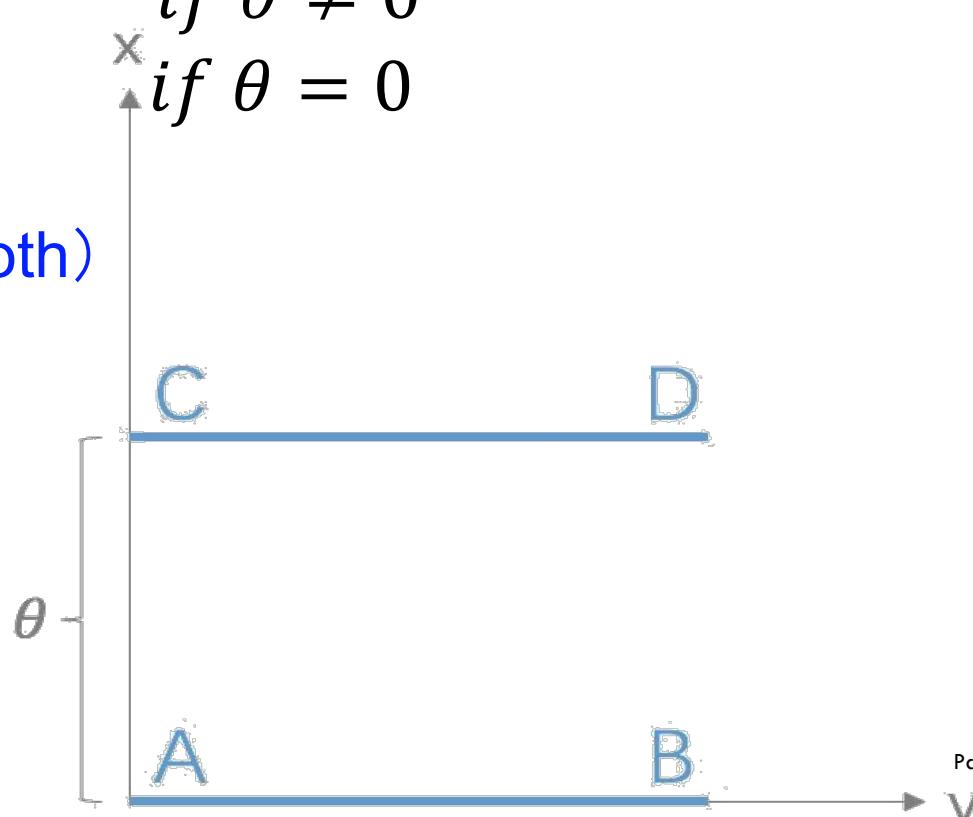
The Wasserstein (or Earth-Mover) distance is then the “cost” of the **optimal** transport plan

Examples

$$KL(P_1 || P_2) = \begin{cases} +\infty & \text{if } \theta \neq 0 \\ 0 & \text{if } \theta = 0 \end{cases}$$

$$JS(P_1 || P_2) = \begin{cases} \log 2 & \text{if } \theta \neq 0 \\ 0 & \text{if } \theta = 0 \end{cases}$$

$$W(P_1, P_2) = |\theta| \text{ (Smooth)}$$



Wasserstein GANs

- The Earth-Mover (EM) distance or Wasserstein-1:

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

have nicer properties when optimized than JS

- However, the infimum is highly intractable.
- Wasserstein distance has a duality form

$$\begin{aligned} W(P_r, P_g) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim P_r} [f(x)] - \mathbb{E}_{x \sim P_g} [f(x)] \\ &= \frac{1}{K} \sup_{\|f\|_L \leq K} \mathbb{E}_{x \sim P_r} [f(x)] - \mathbb{E}_{x \sim P_g} [f(x)] \end{aligned}$$

where supremum is over all the K-Lipschitz functions

Wasserstein GANs

- Wasserstein distance has a duality form

$$\begin{aligned} W(P_r, P_g) &= \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim P_r}[f(x)] - \mathbb{E}_{x \sim P_g}[f(x)] \\ &= \frac{1}{K} \sup_{\|f\|_L \leq K} \mathbb{E}_{x \sim P_r}[f(x)] - \mathbb{E}_{x \sim P_g}[f(x)] \end{aligned}$$

where supremum is over all the K-Lipschitz functions

- Consider a w -parameterized family of functions $\{f_w\}_{w \in W}$ that are all K -Lipschitz

$$W(P_r, P_g) = \max_{w \in W} \mathbb{E}_{x \sim P_r}[f_w(x)] - \mathbb{E}_{x \sim P_g}[f_w(x)]$$

For example, $W = [-c, c]^l$. To satisfy this requirement, WGAN enforces the weights of D lie within a compact space $[-c, c]$ by applying **weight clipping**

Wasserstein GANs

- The loss for **discriminator/critic** f_w

$$\mathbb{E}_{x \sim P_r} [f_w(x)] - \mathbb{E}_{x \sim P_g} [f_w(x)]$$

- f_w 's goal: maximize Wasserstein distance between real data distribution and generative distribution
- The loss for **generator**
 - $-\mathbb{E}_{x \sim P_g} [f_w(x)] = -\mathbb{E}_{z \sim p(z)} [f_w(g_\theta(z))]$
- g_θ 's goal: minimize Wasserstein distance between real data distribution and generative distribution

Algorithm

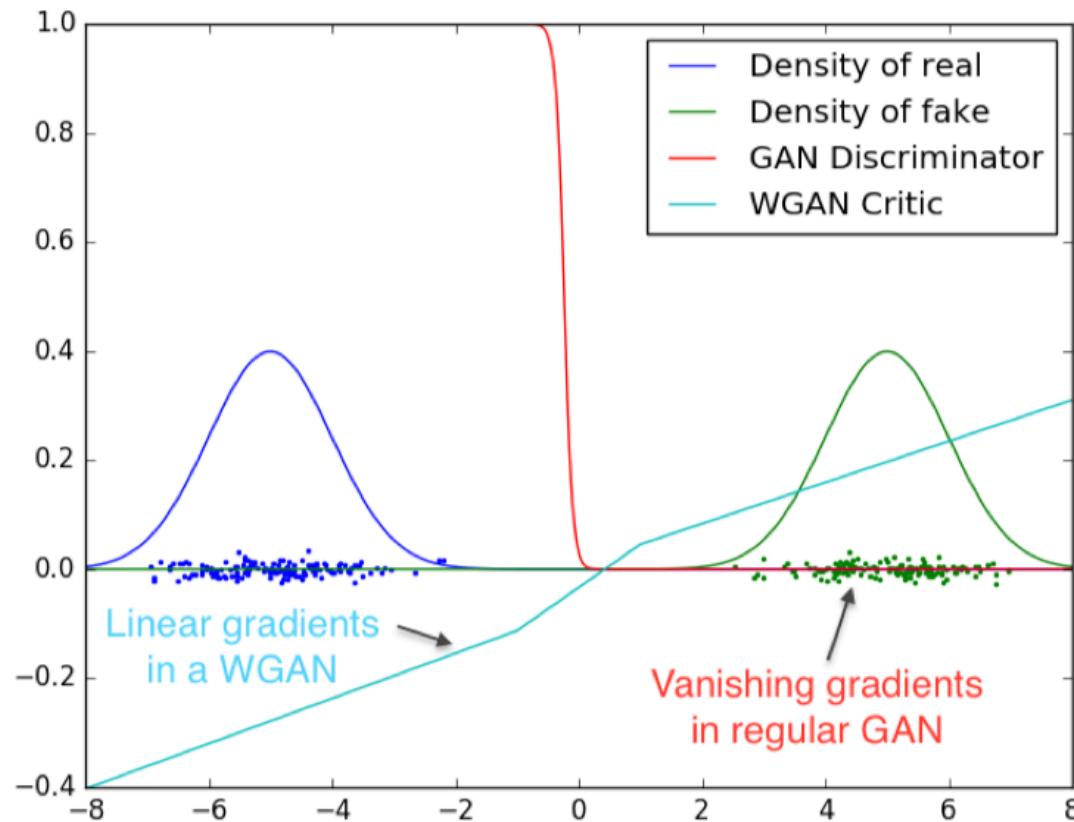
Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

Require: : α , the learning rate. c , the clipping parameter. m , the batch size. n_{critic} , the number of iterations of the critic per generator iteration.

Require: : w_0 , initial critic parameters. θ_0 , initial generator's parameters.

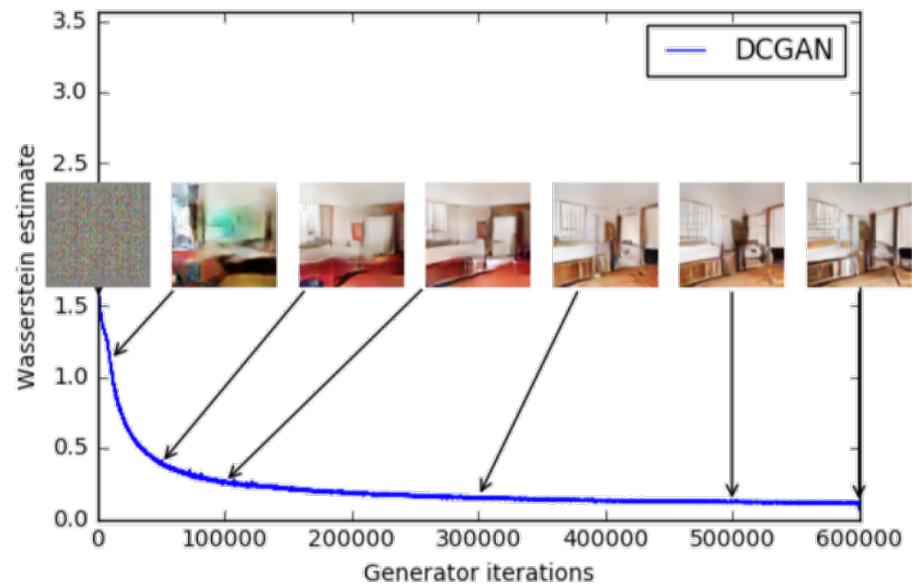
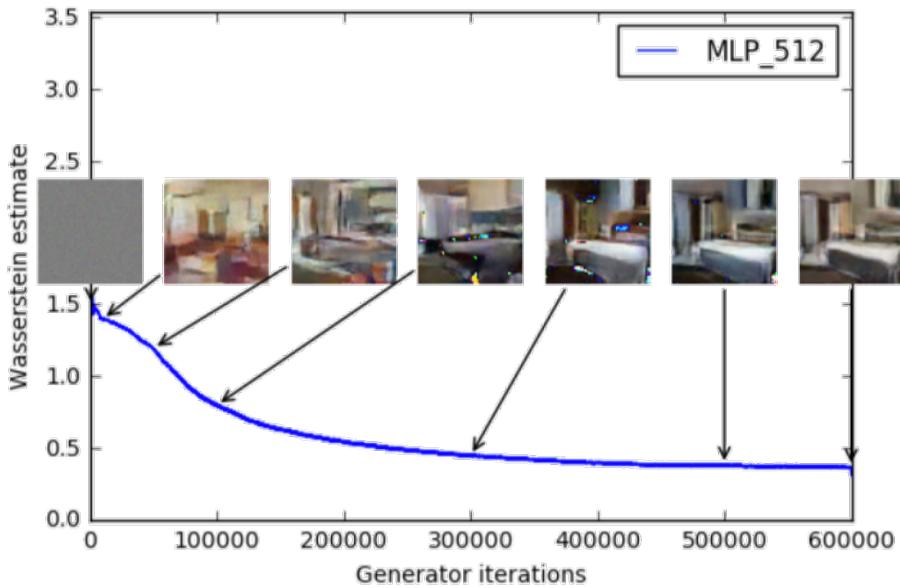
```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$  → compute gradient of  $f_w$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$  → gradient ascent on  $f_w$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$  → weight clipping to satisfy that functions  $\{f_w\}_{w \in W}$  are all K - Lipschitz
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$  → compute gradient of  $g_\theta$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$  → gradient decent on  $g_\theta$ 
12: end while
```

Algorithm



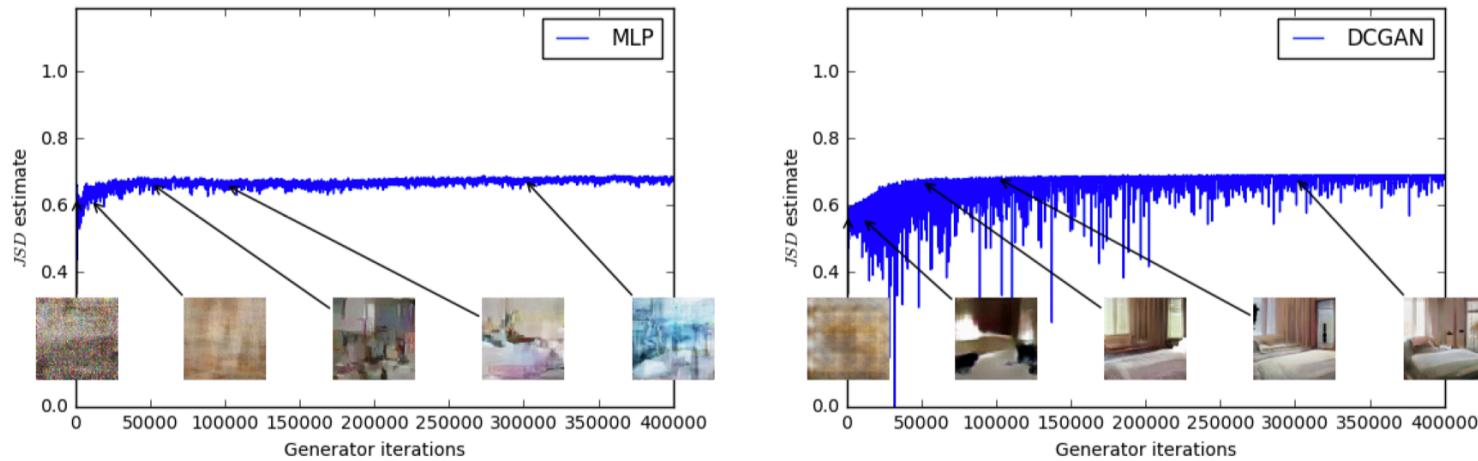
Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the discriminator of a minimax GAN saturates and results in vanishing gradients. WGAN critic provides very clean gradients on all parts of the space.

Meaningful loss metric



Meaningful loss metric

- Vanilla GANs



JS estimates for an MLP generator (upper left) and a DCGAN generator (upper right) trained with the standard GAN procedure. Both had a DCGAN discriminator. Both curves have increasing error. Samples get better for the DCGAN but the JS estimate increases or stays constant, pointing towards no significant correlation between sample quality and loss.

Image-to-Image Translation



Image de-raining



Image inpainting



Sketch -> Image



Image Super-resolution

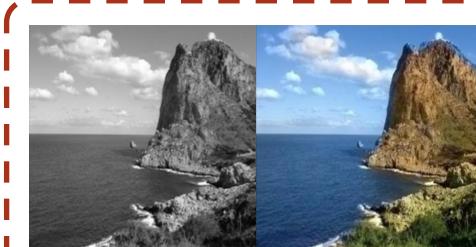
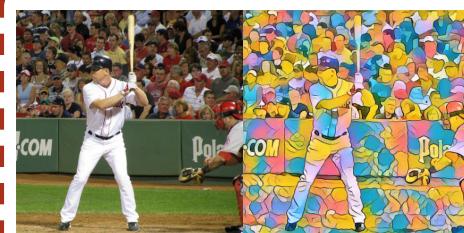


Image coloration



Rotated images



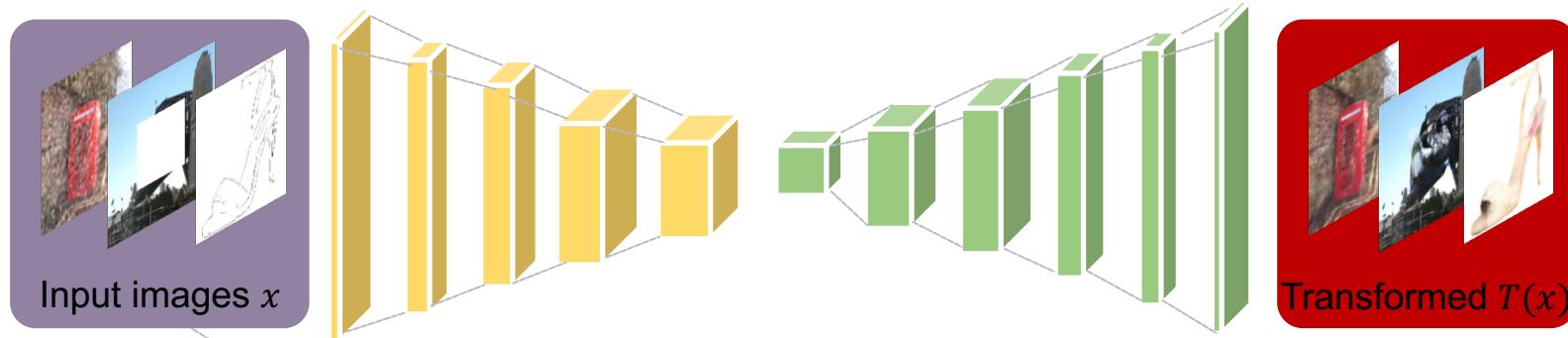
Style transfer



Day -> Night

Image-to-Image Translation

Image transformation network T



Pixel-wise loss



GANs loss

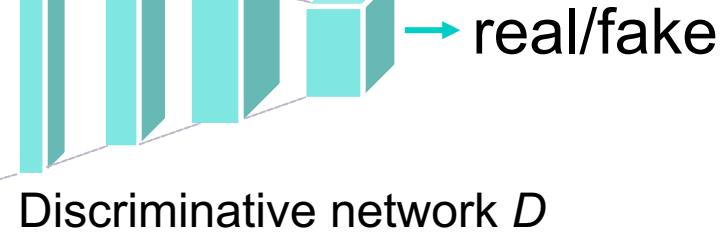
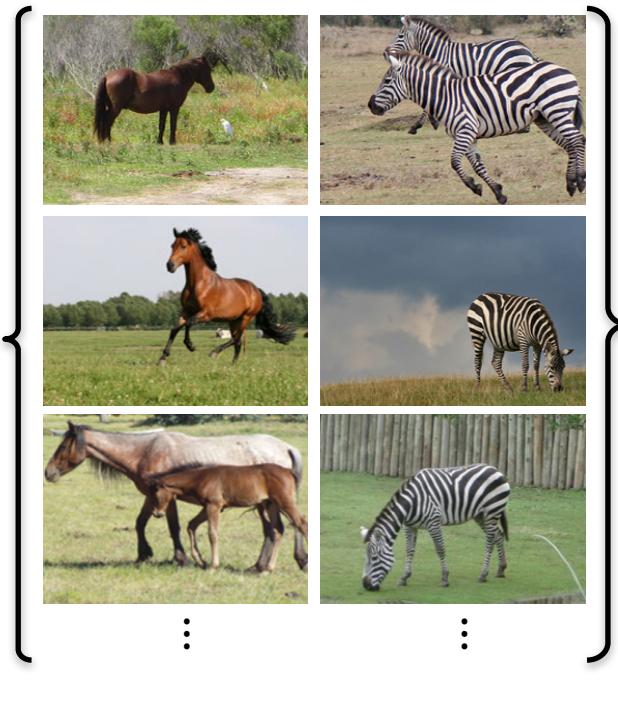
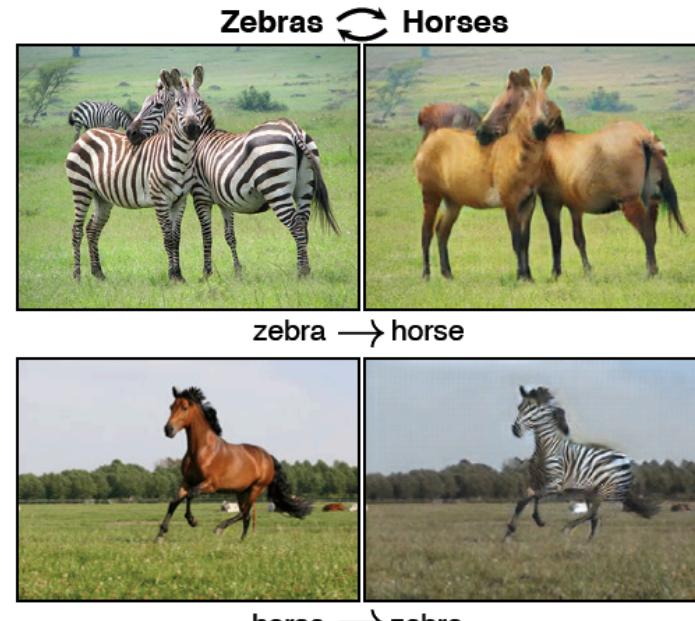


Image-to-Image Translation

- Unpaired Image-to-image Translation (CycleGAN)



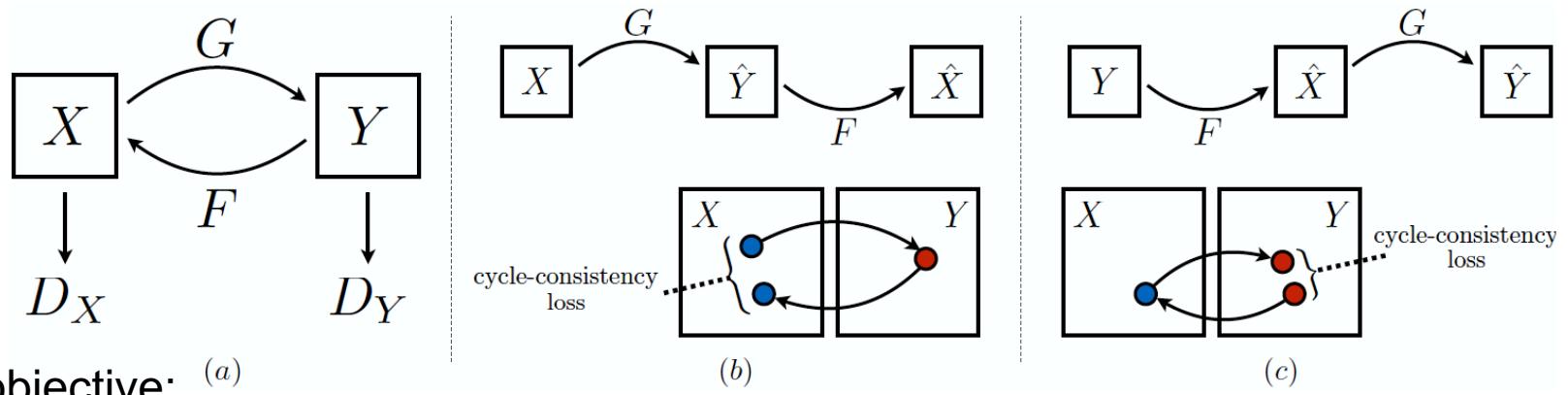
Unpaired Datasets



Object Transfiguration

CycleGAN

- Unpaired Image-to-image Translation (CycleGAN)



Full objective: (a)

$$\begin{aligned}\mathcal{L}(\mathcal{G}, \mathcal{F}, D_X, D_Y) = & \mathcal{L}_{GAN}(\mathcal{G}, D_Y, X, Y) \\ & + \mathcal{L}_{GAN}(\mathcal{F}, D_X, Y, X) + \lambda \mathcal{L}_{cyc}(\mathcal{G}, \mathcal{F}),\end{aligned}$$

$$\mathcal{L}_{cyc}(G, F) = \mathbb{E}_{x \in X} [\|\mathcal{F}(\mathcal{G}(x)) - x\|_1] + \mathbb{E}_{y \in Y} [\|\mathcal{G}(\mathcal{F}(y)) - y\|_1].$$

$$\mathcal{L}_{GAN}(\mathcal{G}, D_Y, X, Y) = \mathbb{E}_{y \in Y} [D_Y^2(y)] + \mathbb{E}_{x \in X} [(D_Y(\mathcal{G}(x)) - 1)^2],$$

Attention-GAN for object transfiguration in wild images

[ECCV 2018]

- Two responsibilities for generator G:
 1. detect the objects of interests
 2. convert the object from source domain to another domain.

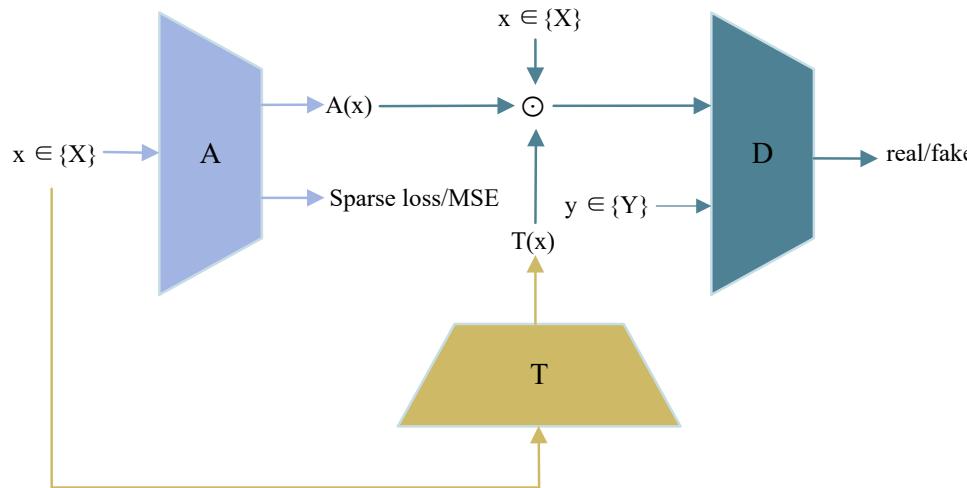


Figure 3 An illustration of Attention-GAN. A , T , D respectively represent the attention network, the transformation network and the discriminative network. $A(x)$ denotes the attention map predicted by the attention network. $T(x)$ denotes the transformed images. \odot denotes the layered operation.

Attention-GAN

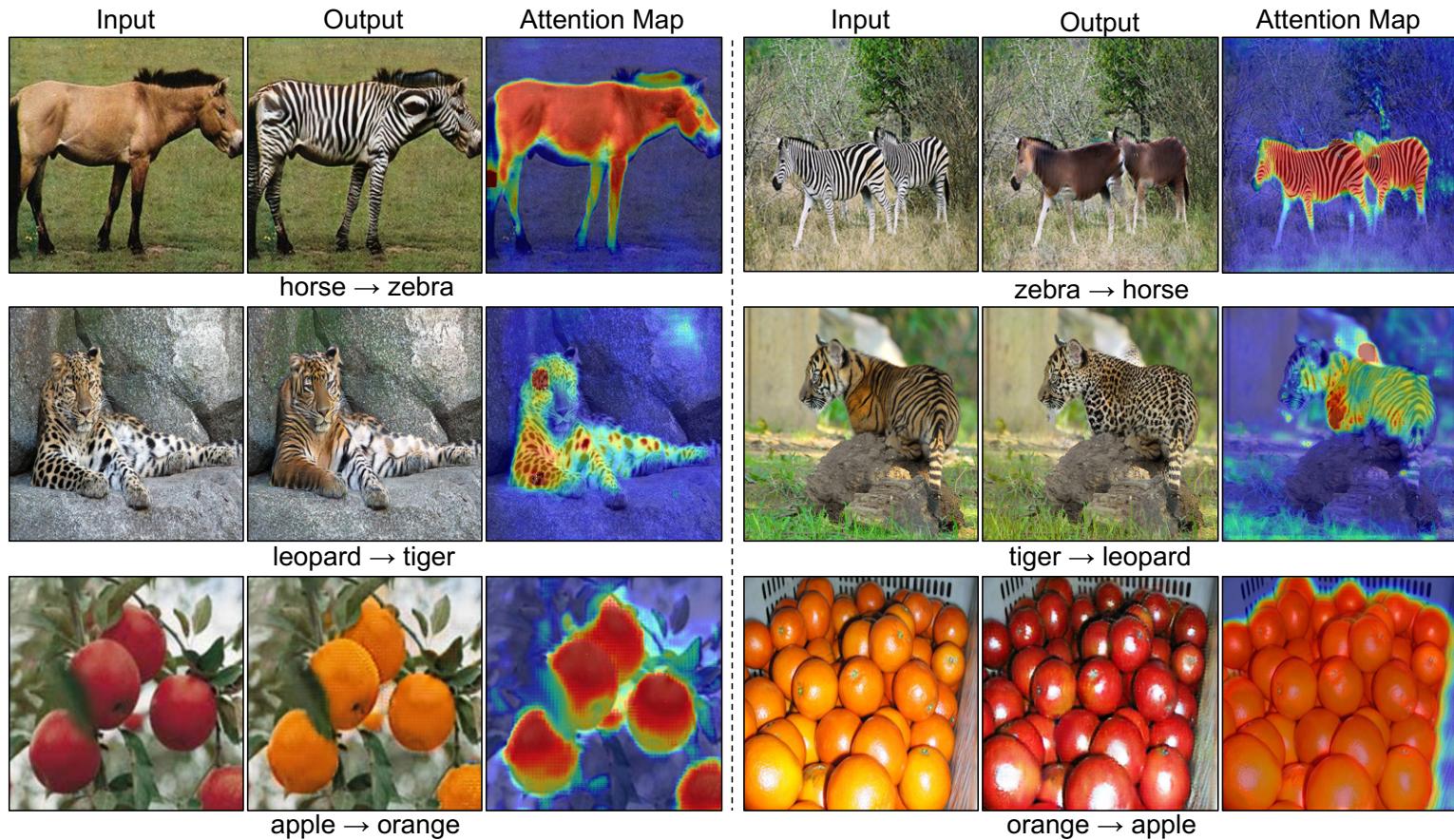


Figure 1 Results of object transfiguration on tasks: horse ↔ zebra, leopard ↔ tiger and apple↔orange. In each case, the first image is the original images, the second image is the synthesized image, and the third image is the predicted attention map. Our proposed model only manipulates the attention parts of image and preserves the background.

Architecture

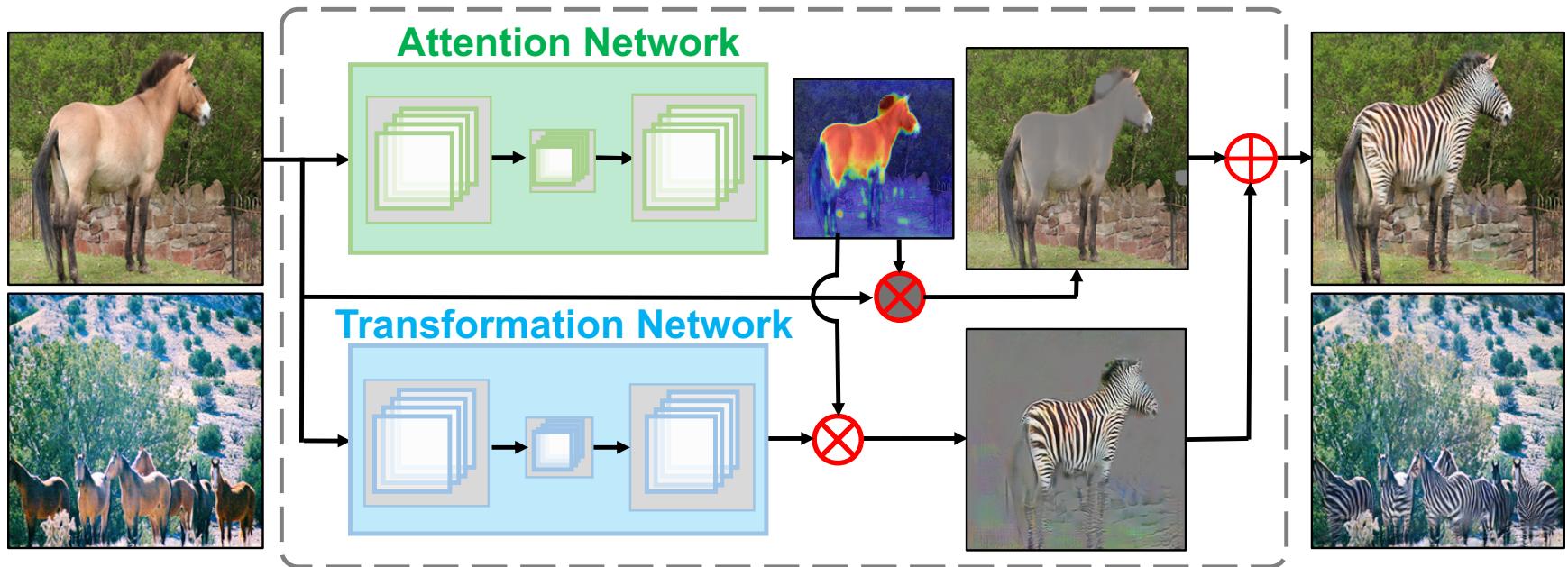
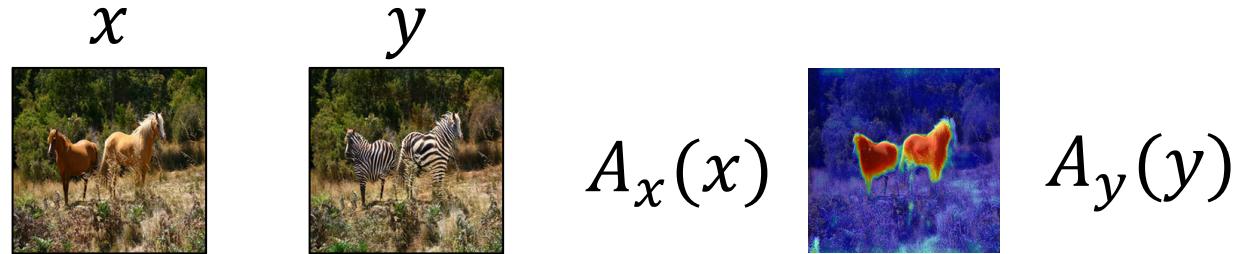


Figure 4 The proposed Attention-GAN for object transfiguration from one class to another. The attention network predicts the attention maps. The transformation network synthesizes the target object. A layered operation is applied on the background and transformed images to output the resulting image.

· Formula: $\mathcal{G}(x) \equiv A_X(x) \odot T_X(x) + (1 - A_X(x)) \odot x$

$$\mathcal{F}(y) \equiv A_Y(y) \odot T_Y(y) + (1 - A_Y(y)) \odot y.$$

Attention Loss



The regions of interest in the original image and the transformed image should be consistent,

$$\begin{aligned} & \mathcal{L}_{A_{cyc}}(A_X, A_Y) \\ &= \mathbb{E}_{x \in X} \left[\|A_X(x) - A_Y(\mathcal{G}(x))\|_1 \right] + \mathbb{E}_{y \in Y} \left[\|A_Y(y) - A_X(\mathcal{F}(y))\|_1 \right] \end{aligned}$$

Sparse loss encourages the attention network to pay attention to a small region related to the object:

$$\mathcal{L}_{sparse}(A_X, A_Y) = \mathbb{E}_{x \in X} [\|A_X(x)\|_1] + \mathbb{E}_{y \in Y} [\|A_Y(y)\|_1]$$

Attention GAN

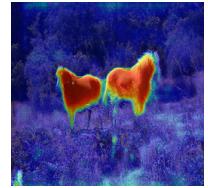
x



y



$A_x(x)$



$A_y(y)$

➤ Full objective Loss

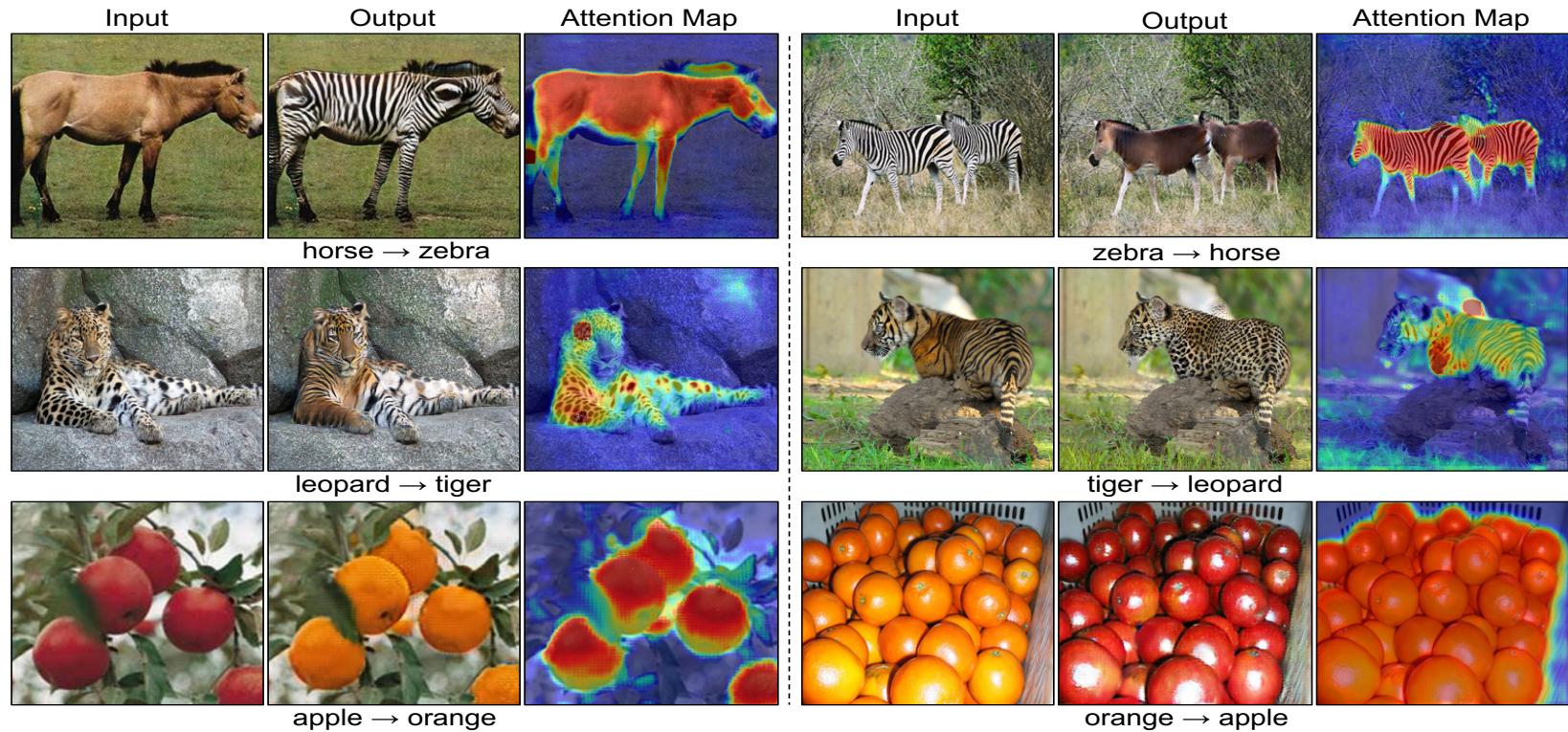
$$\begin{aligned} & \mathcal{L}(T_X, T_Y, D_X, D_Y, A_X, A_Y) \\ &= \mathcal{L}_{GAN}(\mathcal{G}, D_Y, X, Y) + \mathcal{L}_{GAN}(\mathcal{F}, D_X, X, Y) \\ &+ \lambda_{cyc} \mathcal{L}_{cyc}(T_X, T_Y, D_X, D_Y, A_X, A_Y) + \lambda_{A_{cyc}} \mathcal{L}_{A_{cyc}}(A_X, A_Y) \\ &+ \lambda_{sparse} \mathcal{L}_{sparse}(A_X, A_Y) \end{aligned}$$

The networks can be solved from the following min-max game:

$$\arg \min_{T_X, T_Y, A_X, A_Y} \max_{D_X, D_Y} \mathcal{L}(T_X, T_Y, D_X, D_Y, A_X, A_Y).$$

Experiment

➤ Qualitative results



Human Study

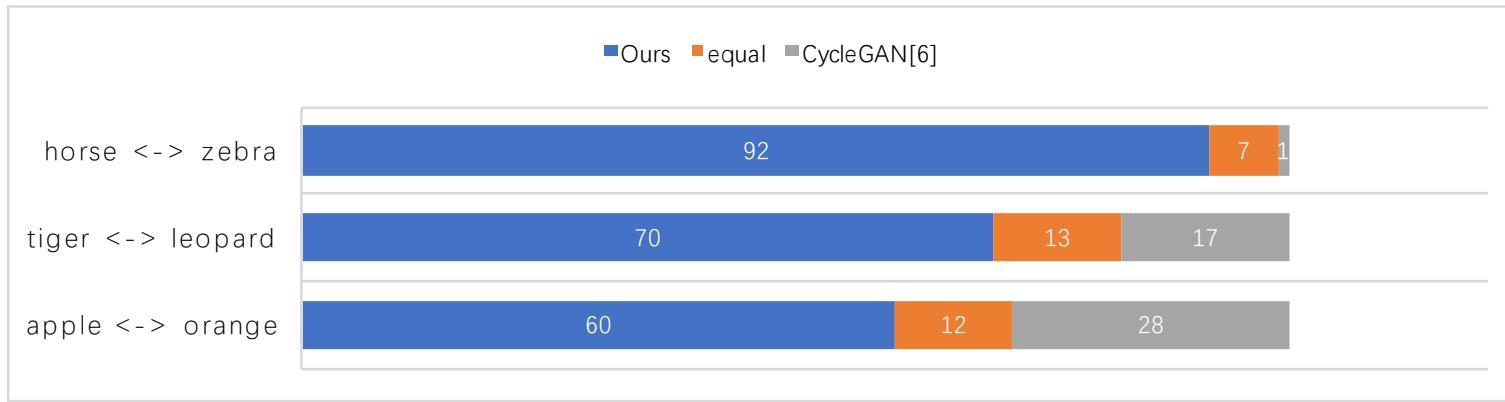


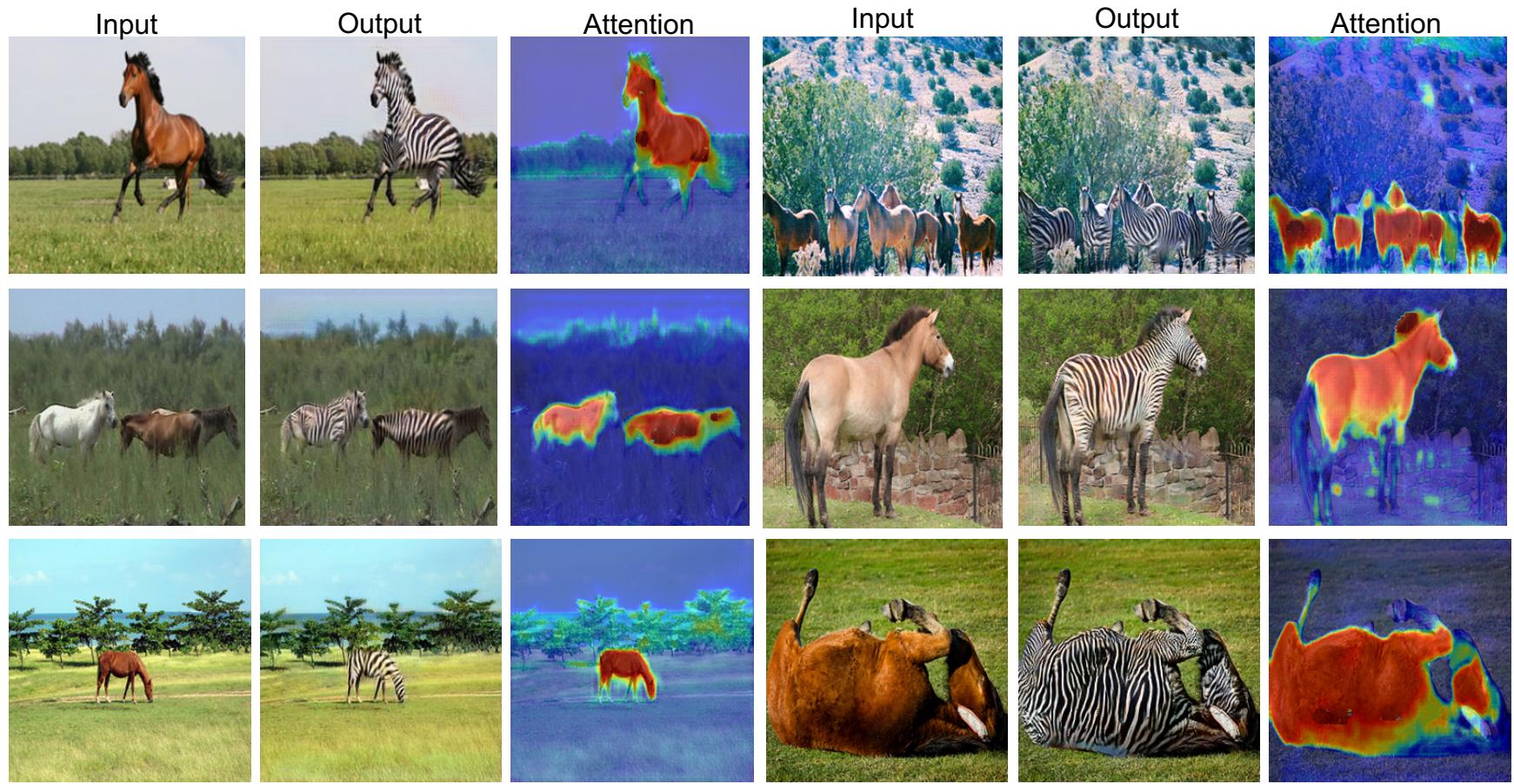
Figure. The stacked bar chart of participants preferences for our methods compared to CycleGAN [21]. The blue bar indicates the number of images that more participants prefer our results. The gray bar indicates the number of images that more participants prefer CycleGAN's results. The orange bar indicates the number of images where two methods get a equal number of votes from 10 participants.

Experiments

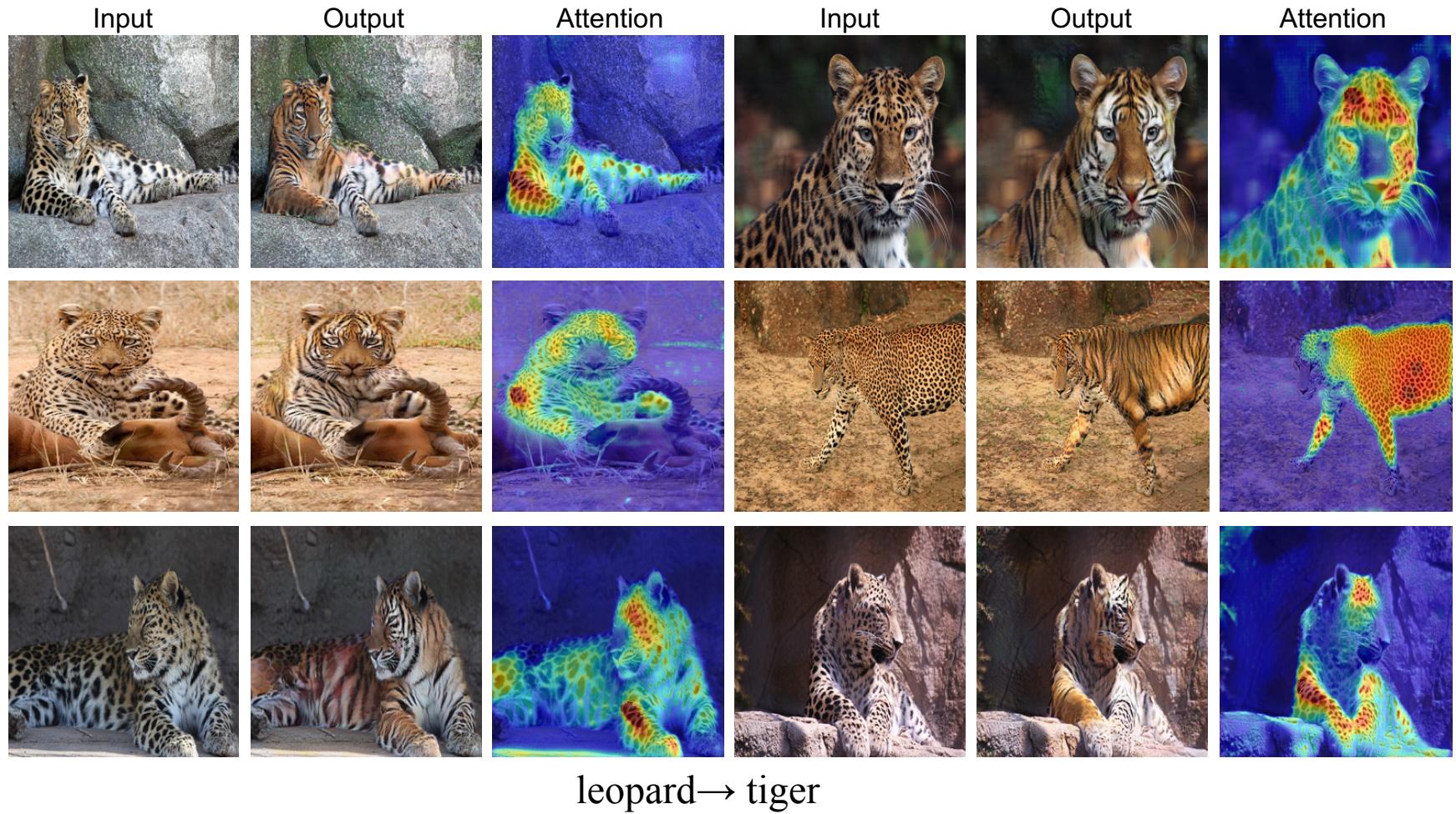
➤ Quantitative comparison:

	Task	CycleGAN	DistanceGAN	Ours (Unsupervised)	Ours (Supervised)
PSNR	horse → zebra	18.1875	11.1896	22.2629	24.589
	zebra → horse	18.1021	10.1153	21.5360	23.9330
SSIM	horse → zebra	0.6725	0.2630	0.9003	0.9482
	zebra → horse	0.7155	0.3627	0.8988	0.9534

More Results



More Results



More Results

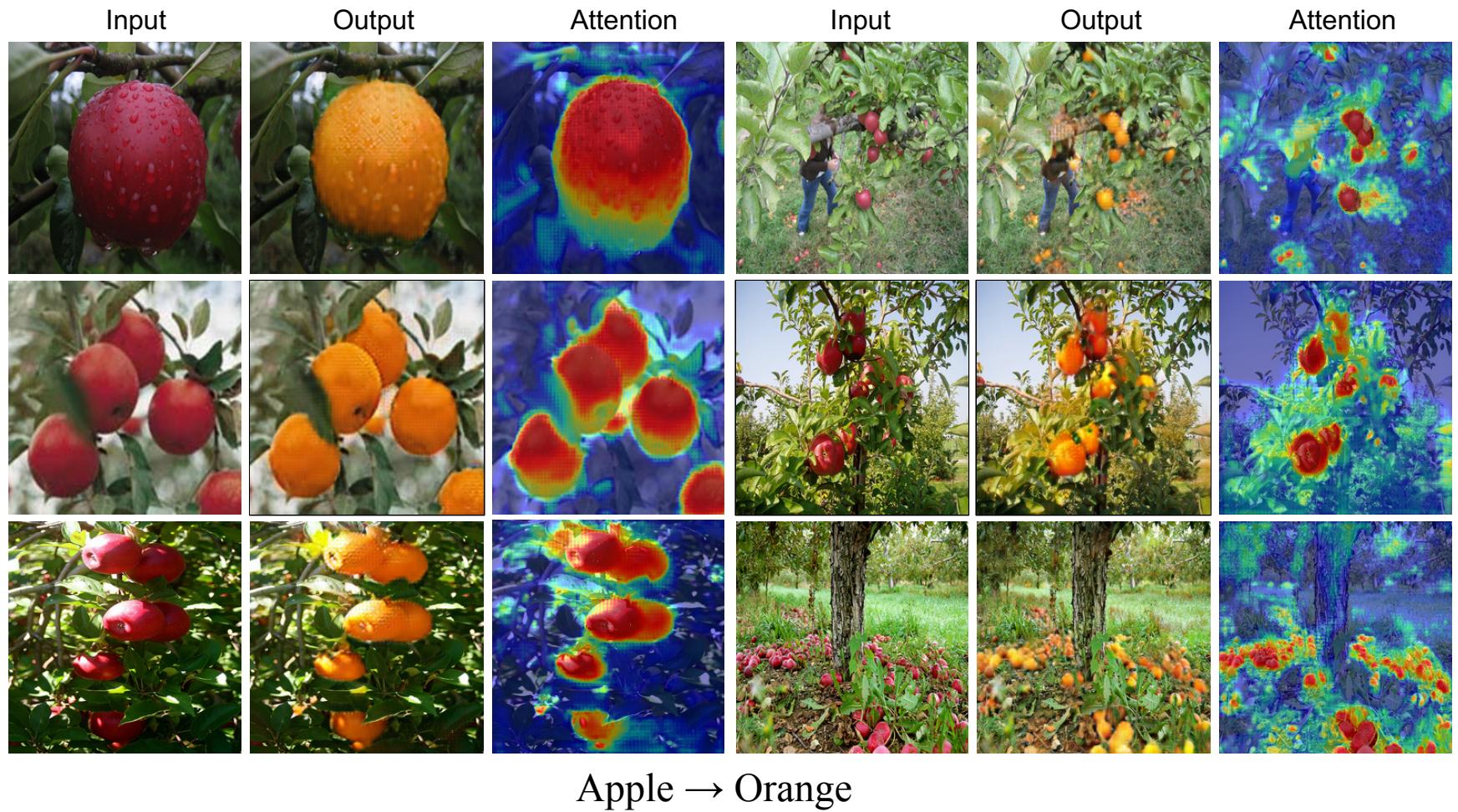
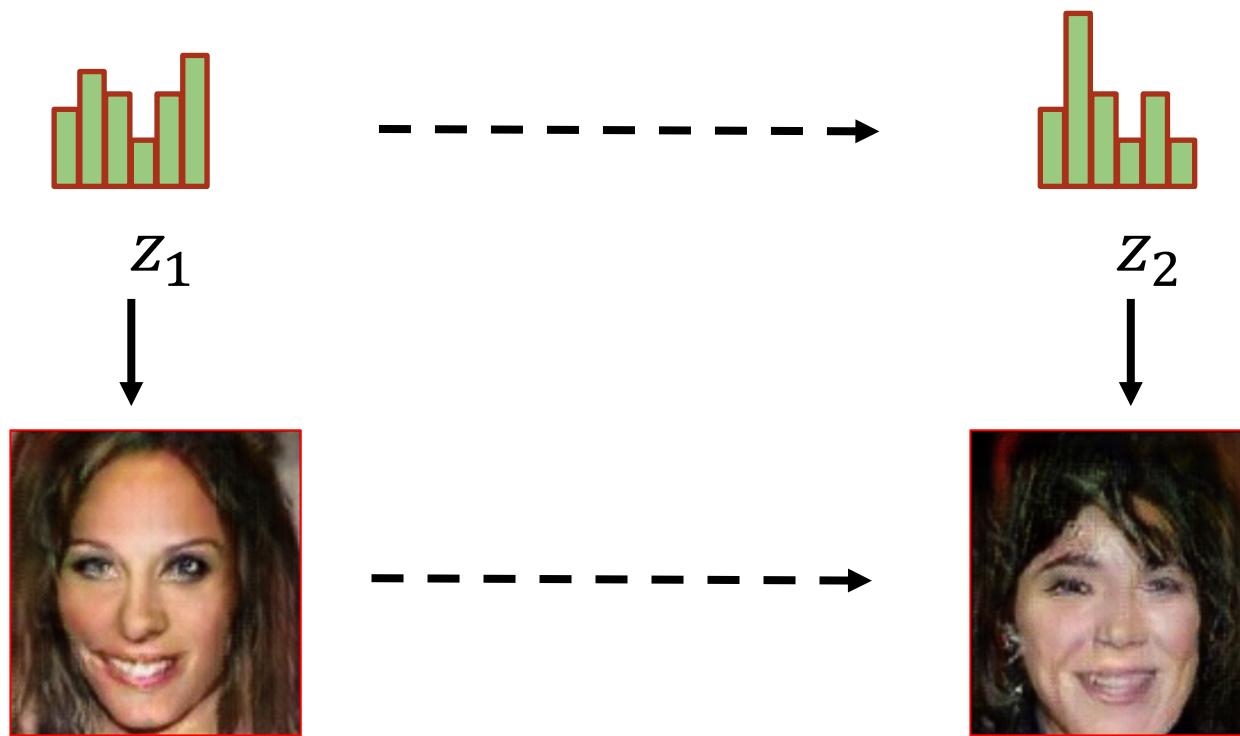


Image Interpolation



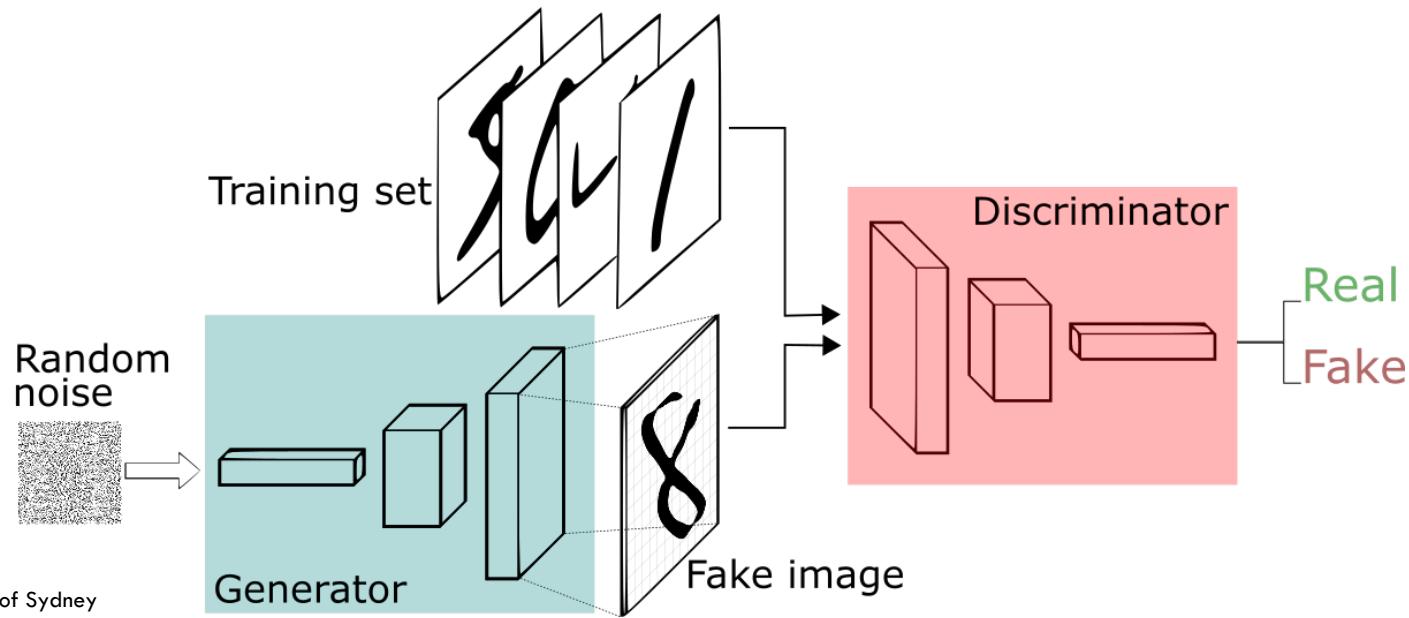
Noise vector inputs for GANs can be taken as low-dimensional representations of images. As widely accepted in representation learning, the closeness of two data points is supposed to be preserved before and after transformation.

Smooth Deep Image Generator from Noises [AAAI 2019]

A **smooth generator** is then expected to have the following pixel-wise property:

$$|G_i(z + \delta) - G_i(z)| < \epsilon$$

where G is the generator, δ stands for a small perturbation over input noise z , and $\epsilon > 0$ is a small constant number.



Smooth Generator

$$|G_i(\mathbf{z} + \boldsymbol{\delta}) - G_i(\mathbf{z})| < \epsilon$$

How to specify the small perturbation $\boldsymbol{\delta}$?

Theorem 1. Fix $\epsilon > 0$. Given $\mathbf{z} \in \mathbb{R}^d$ as the noise input of generator G_i , if the perturbation $\boldsymbol{\delta} \in \mathbb{R}^d$ satisfies

$$\|\boldsymbol{\delta}\|_q < \frac{\epsilon}{\max_{\hat{\mathbf{z}} \in B_p(\mathbf{z}, R)} \|\nabla_{\hat{\mathbf{z}}} G_i(\hat{\mathbf{z}})\|_p}, \quad (3)$$

with $\frac{1}{p} + \frac{1}{q} = 1$, we have $|G_i(\mathbf{z} + \boldsymbol{\delta}) - G_i(\mathbf{z})| < \epsilon$.

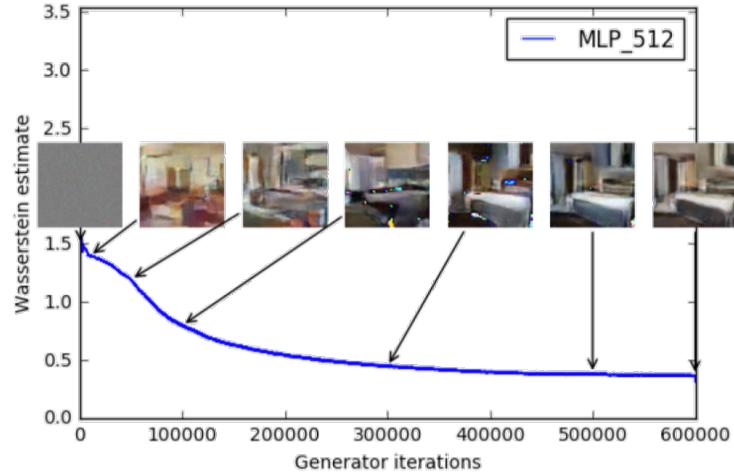
By minimizing the denominator in the above equation, the model is expected to tolerate larger perturbation $\boldsymbol{\delta}$ under fixed difference ϵ on the i -th pixel.

Smooth GANs

Besides the smoothness of generated images, how about their quality?

In W-GAN, loss function could reflect the image quality:

$$\min_D \mathbb{E}_{z \sim \mathbb{P}_z} [D(G(z))] - \mathbb{E}_{x \sim \mathbb{P}_d} [D(x)].$$



To decrease the quality gap between images generated from closed noise inputs, we need to ensure that $D(G(z))$ would not drop significantly when the input variates

$$|D[G(z + \delta)] - D[G(z)]| < \epsilon.$$

Smooth GANs

Theorem 2. Fix $\epsilon > 0$. Consider generator G and discriminator D in GANs. Given a noise input $\mathbf{z} \in \mathbb{R}^d$, the generated image is $\hat{\mathbf{x}} = G(\hat{\mathbf{z}})$. If the perturbation $\boldsymbol{\delta} \in \mathbb{R}^d$ satisfies

$$\|\boldsymbol{\delta}\|_q < \frac{\epsilon}{\max_{\hat{\mathbf{z}} \in B_p(\mathbf{z}, R)} \|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_p \|\nabla_{\hat{\mathbf{z}}} G(\hat{\mathbf{z}})\|_p}, \quad (17)$$

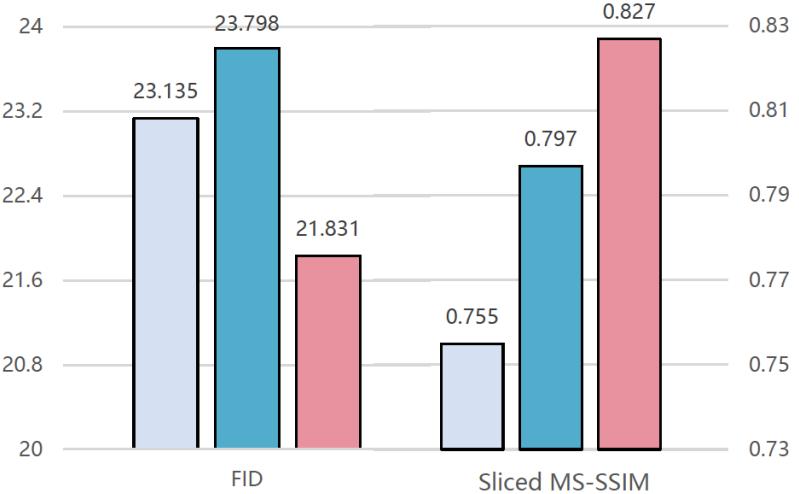
with $\frac{1}{p} + \frac{1}{q} = 1$, we have $|D[G(\mathbf{z} + \boldsymbol{\delta})] - D[G(\mathbf{z})]| < \epsilon$.

- Connection with WGAN-GP;
- Encourage the norm of the gradient of generator to stay at a lower level.

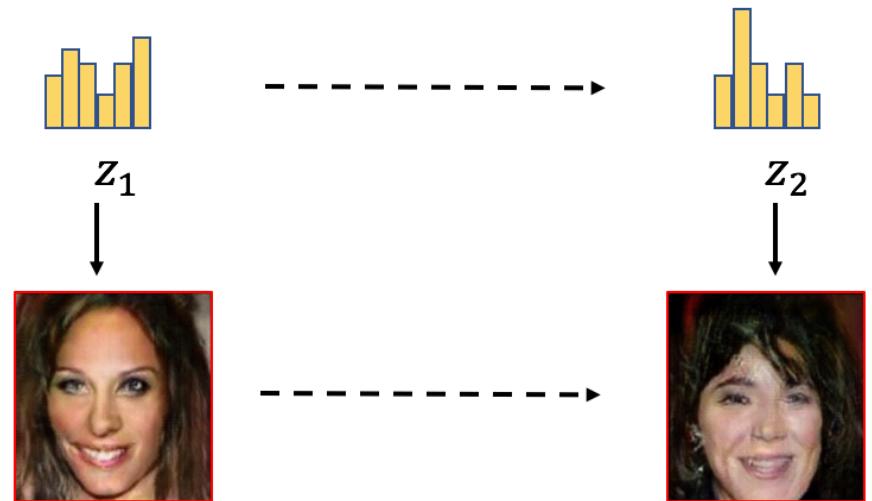
Experiments



□ DCGAN □ WGAN-GP ■ Smooth GAN



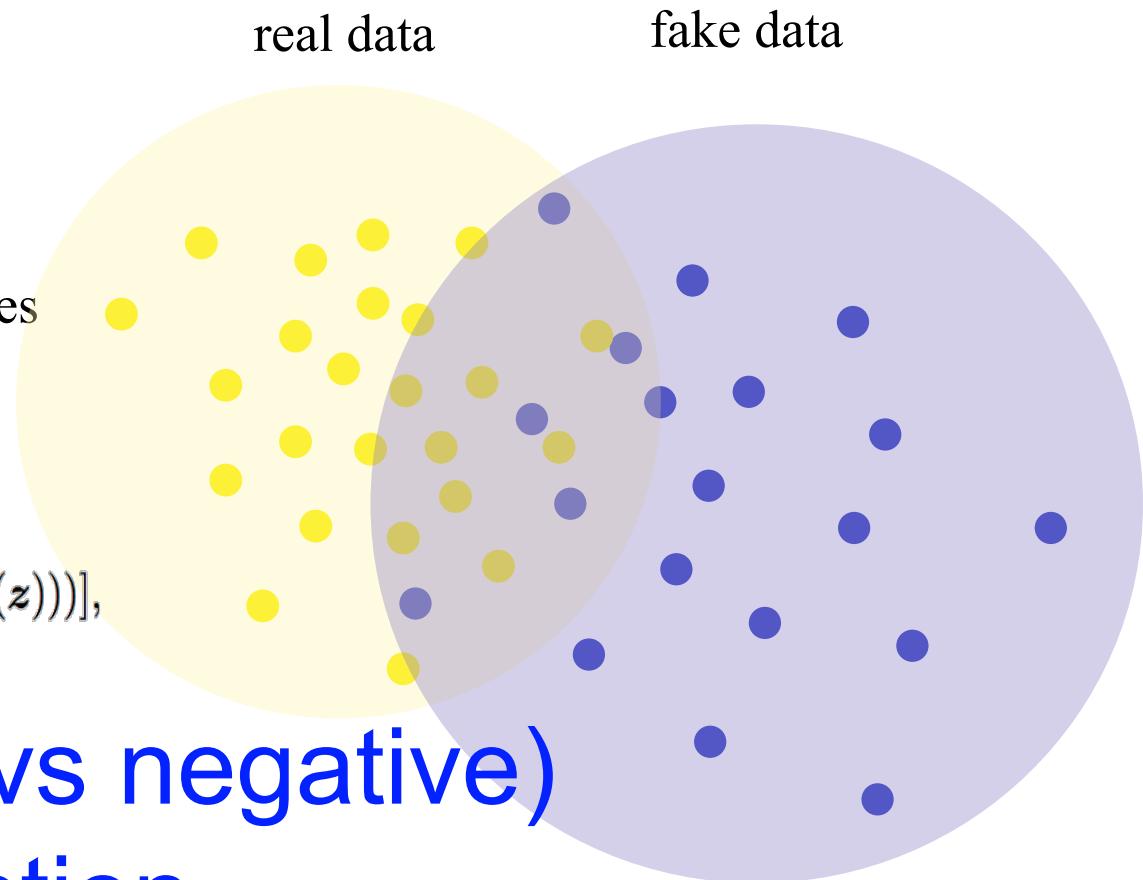
FID and sliced MS-SSIM of interpolated images on the CIFAR-10 dataset.



On Positive-Unlabeled Classification in GAN [CVPR 2020]

Traditional GANs

- Separating real and fake data
- Ignoring good generated samples



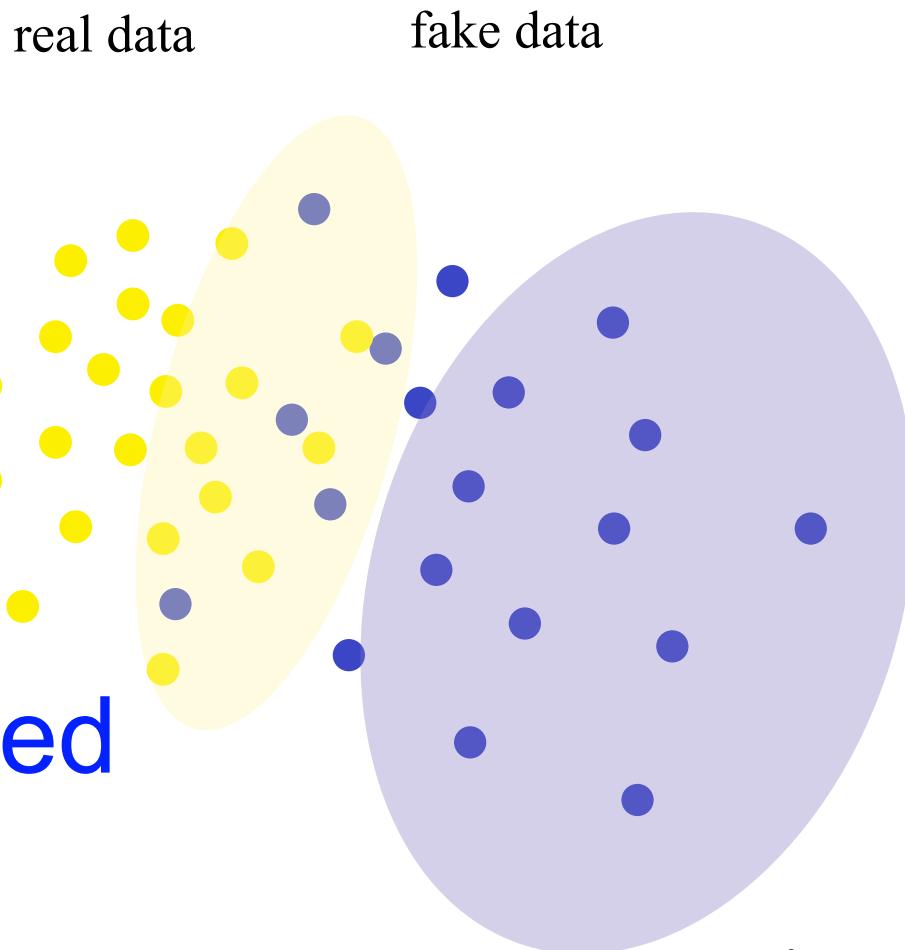
$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Binary (positive vs negative)
classification

Positive-unlabeled GAN (PUGAN)

- Separating good samples from generated samples referring to real data
- Focusing on low-quality samples

$$\begin{aligned} \max_G \min_D V(D, G) = & \pi \mathbb{E}_{p_{data}} [f_1(D(\mathbf{x}))] \\ & + \max \{0, \mathbb{E}_{p_z} [f_2(D(G(\mathbf{z})))]\} \\ & - \pi \mathbb{E}_{p_{data}} [f_2(D(\mathbf{x}))]. \end{aligned}$$



Positive and unlabeled classification

Positive-unlabeled GAN (PUGAN)

Loss	SGAN	PUSGAN	HingeGAN	PUHingeGAN	LSGAN	PULSGAN	WGAN-GP	PUWGAN-GP
MNIST	18.65	16.53	21.48	16.91	13.47	13.96	17.42	14.77
Fashion	25.72	24.33	28.39	25.31	32.05	26.72	26.50	23.12
CIFAR-10	43.39	31.02	43.85	36.37	27.64	22.32	36.86	31.85
CelebA-64	48.44	43.89	46.13	43.94	51.67	47.80	36.09	35.72
CAT-64	46.13	16.90	29.52	25.72	57.22	26.79	21.86	17.35

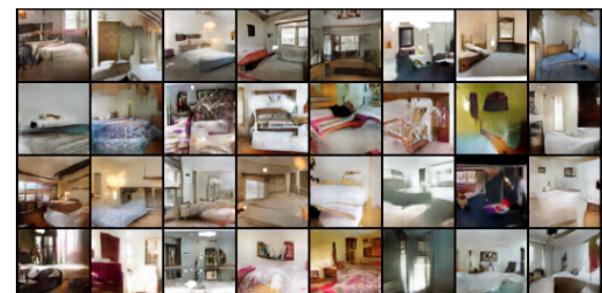
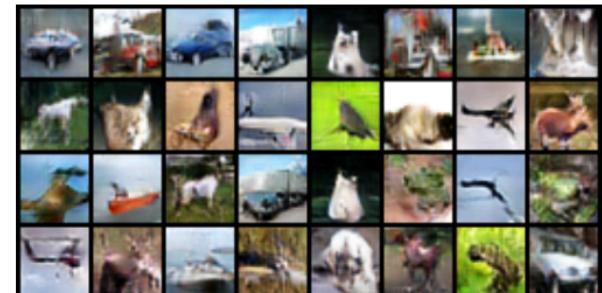
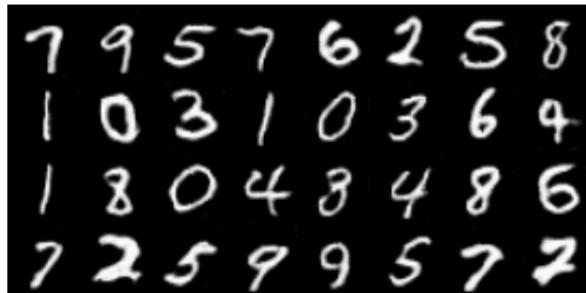
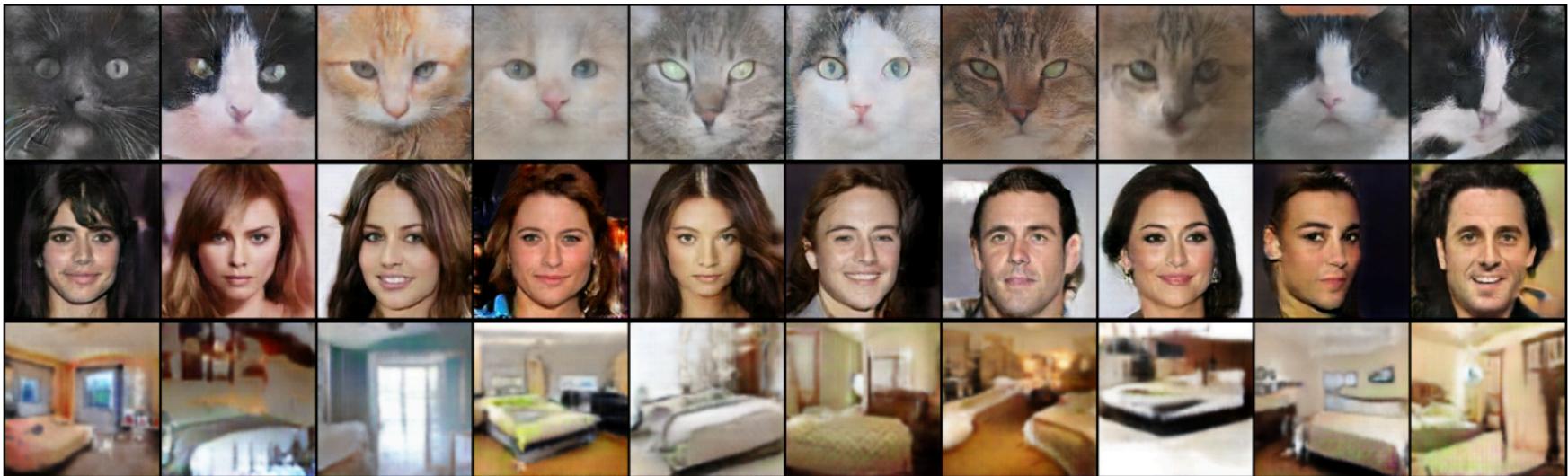


Figure 1: Generated samples obtained by the proposed method on image datasets.



(a) Generation samples on 128 resolution datasets.



(b) Generation samples on the CAT-256 dataset.

Figure 2: High-resolution generated samples.

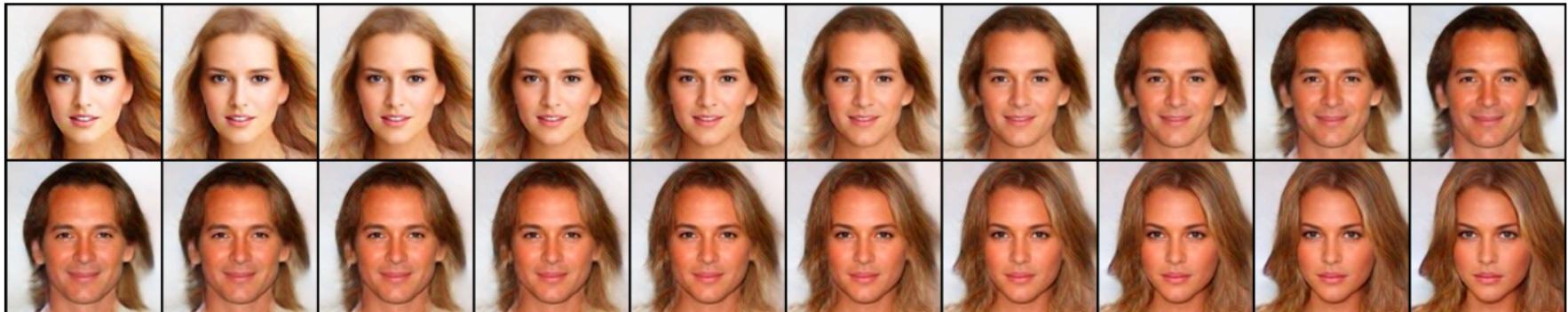


Figure 3: High-resolution generated samples.