

Real or not real, that is the question

(ICLR 2020 spotlight)

The Chinese University of Hong Kong

# Observation

- What makes those pictures look real or not?
  - Inharmonious facial structure and components
  - Unnatural background
  - Abnormal style combination
  - Texture distortion
  - ...



(a)



(b)

# Observation

- Discriminator predict the picture is real or not
  - Although there are several perspectives for determining, the output from the discriminator is **a single scalar**
  - The single scalar could be viewed as an abstract or summarization of multiple measures

# Idea

- RealnessGAN
  - The output of discriminator is a distribution
  - Standard GAN can be viewed as a special case of RealnessGAN

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

$$= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(\underline{D(\mathbf{x}) - 0})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(\underline{1 - D(\mathbf{x})})] \quad (2)$$

- Interpretation: difference between one dimensional vector  $D(\mathbf{x})$  and one dimensional vector  $[0] / [1]$

# Idea

- RealnessGAN

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

$$= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log(\underline{D(\mathbf{x}) - 0})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(\underline{1 - D(\mathbf{x})})] \quad (2)$$

- Interpretation: difference between one dimensional vector  $D(\mathbf{x})$  and one dimensional vector  $[0]$  /  $[1]$
- What if a multi-dimensional case?

$$\max_G \min_D V(G, D) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\mathcal{D}_{\text{KL}}(\mathcal{A}_1 \| D(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim p_g} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(\mathbf{x}))]. \quad (3)$$

- Difference between two multi-dimensional vectors (distrib.)

# Theory

Future works:

Replacing KL divergence by earth mover distance

Applying on improved architectures (progressiveGAN, styleGAN)

- Ground-truth distribution
  - A1: real, A0: fake (we call these anchor distributions)

$$\max_G \min_D V(G, D) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\mathcal{D}_{\text{KL}}(\mathcal{A}_1 \| D(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim p_g} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(\mathbf{x}))]. \quad (3)$$

- Optimality  $D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$  for a fixed  $G$

**Theorem 1.** When  $G$  is fixed, for any *outcome*  $u$  and input sample  $\mathbf{x}$ , the optimal discriminator  $D$  satisfies

$$D_G^*(\mathbf{x}, u) = \frac{\mathcal{A}_1(u)p_{\text{data}}(\mathbf{x}) + \mathcal{A}_0(u)p_g(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}. \quad (4)$$

# Theory

*Proof.* Given a fixed  $G$ , the objective of  $D$  is:

$$\min_D V(G, D) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\mathcal{D}_{\text{KL}}(\mathcal{A}_1 \| D(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim p_g} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(\mathbf{x}))], \quad (17)$$

$$= \int_{\mathbf{x}} \left( p_{\text{data}}(\mathbf{x}) \int_u \mathcal{A}_1(u) \log \frac{\mathcal{A}_1(u)}{D(\mathbf{x}, u)} du + p_g(\mathbf{x}) \int_u \mathcal{A}_0(u) \log \frac{\mathcal{A}_0(u)}{D(\mathbf{x}, u)} du \right) dx, \quad (18)$$

$$\begin{aligned} &= - \int_{\mathbf{x}} (p_{\text{data}}(\mathbf{x}) h(\mathcal{A}_1) + p_g(\mathbf{x}) h(\mathcal{A}_0)) dx \\ &\quad - \int_{\mathbf{x}} \int_u (p_{\text{data}}(\mathbf{x}) \mathcal{A}_1(u) + p_g(\mathbf{x}) \mathcal{A}_0(u)) \log D(\mathbf{x}, u) du dx, \end{aligned} \quad (19)$$

# Theory

$$\begin{aligned}
 &= - \int_{\mathbf{x}} (p_{\text{data}}(\mathbf{x})h(\mathcal{A}_1) + p_g(\mathbf{x})h(\mathcal{A}_0)) d\mathbf{x} \\
 &\quad - \int_{\mathbf{x}} \int_u (p_{\text{data}}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u)) \log D(\mathbf{x}, u) du d\mathbf{x},
 \end{aligned} \tag{19}$$

where  $h(\mathcal{A}_1)$  and  $h(\mathcal{A}_0)$  are their entropies, and the first term in equation 19 is irrelevant to  $D$ , marked as  $C_1$ . The objective thus is equivalent to:

$$\min_D V(G, D) = - \int_{\mathbf{x}} \int_u (p_{\text{data}}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u)) \log D(\mathbf{x}, u) du d\mathbf{x} + C_1, \tag{20}$$

$$\begin{aligned}
 &= - \int_{\mathbf{x}} (p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})) \int_u \frac{p_{\text{data}}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u)}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})} \log D(\mathbf{x}, u) du d\mathbf{x} + C_1, \\
 &\tag{21}
 \end{aligned}$$

where  $p_{\mathbf{x}}(u) = \frac{p_{\text{data}}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u)}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$  is a distribution defined on  $\Omega_u$ . Consequently, let  $C_2 = p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})$ , we have

$$\min_D V(G, D) = C_1 + \int_{\mathbf{x}} C_2 \left( - \int_u p_{\mathbf{x}}(u) \log D(\mathbf{x}, u) du + h(p_{\mathbf{x}}) - h(p_{\mathbf{x}}) \right) d\mathbf{x}, \tag{22}$$

$$\begin{aligned}
 &= C_1 + \int_{\mathbf{x}} C_2 \mathcal{D}_{\text{KL}}(p_{\mathbf{x}} \| D(\mathbf{x})) d\mathbf{x} + \int_{\mathbf{x}} C_2 h(p_{\mathbf{x}}) d\mathbf{x}. \\
 &\tag{23}
 \end{aligned}$$

From equation 23 we can see, for any  $\mathbf{x} \in \text{Supp}(p_{\text{data}}) \cup \text{Supp}(p_g)$ , when  $\mathcal{D}_{\text{KL}}(p_{\mathbf{x}} \| D(\mathbf{x}))$  achieves its minimum,  $D$  obtains its optimal  $D^*$ . And at that time, we have  $D^*(\mathbf{x}) = p_{\mathbf{x}}$ , which concludes the proof.  $\square$



# Theory

$$D_G^*(\mathbf{x}, u) = \frac{\mathcal{A}_1(u)p_{data}(\mathbf{x}) + \mathcal{A}_0(u)p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}. \quad (4)$$

$$\min_D V(G, D) = \mathbb{E}_{\mathbf{x} \sim p_{data}} [\mathcal{D}_{KL}(\mathcal{A}_1 \| D(\mathbf{x}))] + \mathbb{E}_{\mathbf{x} \sim p_g} [\mathcal{D}_{KL}(\mathcal{A}_0 \| D(\mathbf{x}))], \quad (17)$$

$$= \int_{\mathbf{x}} \left( p_{data}(\mathbf{x}) \int_u \mathcal{A}_1(u) \log \frac{\mathcal{A}_1(u)}{D(\mathbf{x}, u)} du + p_g(\mathbf{x}) \int_u \mathcal{A}_0(u) \log \frac{\mathcal{A}_0(u)}{D(\mathbf{x}, u)} du \right) d\mathbf{x}, \quad (18)$$

**Theorem 2.** When  $D = D_G^*$ , and there exists an *outcome*  $u \in \Omega$  such that  $\mathcal{A}_1(u) \neq \mathcal{A}_0(u)$ , the maximum of  $V(G, D_G^*)$  is achieved if and only if  $p_g = p_{data}$ .

*Proof.* When  $p_g = p_{data}$ ,  $D_G^*(\mathbf{x}, u) = \frac{\mathcal{A}_1(u) + \mathcal{A}_0(u)}{2}$ , we have:

$$V^*(G, D_G^*) = \int_u \mathcal{A}_1(u) \log \frac{2\mathcal{A}_1(u)}{\mathcal{A}_1(u) + \mathcal{A}_0(u)} + \mathcal{A}_0(u) \log \frac{2\mathcal{A}_0(u)}{\mathcal{A}_1(u) + \mathcal{A}_0(u)} du. \quad (7)$$

Subtracting  $V^*(G, D_G^*)$  from  $V(G, D_G^*)$  gives:

$$\begin{aligned} V'(G, D_G^*) &= V(G, D_G^*) - V^*(G, D_G^*) \\ &= \int_{\mathbf{x}} \int_u (p_{data}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u)) \log \frac{(p_{data}(\mathbf{x}) + p_g(\mathbf{x}))(\mathcal{A}_1(u) + \mathcal{A}_0(u))}{2(p_{data}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u))} dud\mathbf{x}, \end{aligned} \quad (8)$$

$$= -2 \int_{\mathbf{x}} \int_u \frac{p_{data}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u)}{2} \log \frac{\frac{p_{data}(\mathbf{x})\mathcal{A}_1(u) + p_g(\mathbf{x})\mathcal{A}_0(u)}{2}}{\frac{(p_{data}(\mathbf{x}) + p_g(\mathbf{x}))(\mathcal{A}_1(u) + \mathcal{A}_0(u))}{4}} dud\mathbf{x}, \quad (9)$$

$$= -2\mathcal{D}_{KL}\left(\frac{p_{data}\mathcal{A}_1 + p_g\mathcal{A}_0}{2} \parallel \frac{(p_{data} + p_g)(\mathcal{A}_1 + \mathcal{A}_0)}{4}\right). \quad (10)$$

# Theory

$$= -2\mathcal{D}_{\text{KL}}\left(\frac{p_{\text{data}}\mathcal{A}_1 + p_g\mathcal{A}_0}{2} \parallel \frac{(p_{\text{data}} + p_g)(\mathcal{A}_1 + \mathcal{A}_0)}{4}\right). \quad (10)$$

Since  $V^*(G, D_G^*)$  is a constant with respect to  $G$ , maximizing  $V(G, D_G^*)$  is equivalent to maximizing  $V'(G, D_G^*)$ . The optimal  $V'(G, D_G^*)$  is achieved if and only if the KL divergence reaches its minimum, where:

$$\frac{p_{\text{data}}\mathcal{A}_1 + p_g\mathcal{A}_0}{2} = \frac{(p_{\text{data}} + p_g)(\mathcal{A}_1 + \mathcal{A}_0)}{4}, \quad (11)$$

$$(p_{\text{data}} - p_g)(\mathcal{A}_1 - \mathcal{A}_0) = 0, \quad (12)$$

for any valid  $\mathbf{x}$  and  $u$ . Hence, as long as there exists a valid  $u$  that  $\mathcal{A}_1(u) \neq \mathcal{A}_0(u)$ , we have  $p_{\text{data}} = p_g$  for any valid  $\mathbf{x}$ .  $\square$

# Discussion

- Number of outcomes
  - Increment of the number of outcomes makes G become rigorous and G needs more effort to learn
  - It is related to the ratio of the number of updates between G and D

# Discussion

- Objective of G

As shown in the theoretical analysis, the ideal objective for  $G$  is maximizing the KL divergence between  $D(\mathbf{x})$  of generated samples and  $\mathcal{A}_0$ :

$$(G_{\text{objective1}}) \quad \min_G -\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(G(\mathbf{z})))]. \quad (14)$$

- Alternative objectives (regularization)
  - Since D is not always optimal (A1: real, A0: fake)

$$(G_{\text{objective2}}) \quad \min_G \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}, \mathbf{z} \sim p_{\mathbf{z}}} [\mathcal{D}_{\text{KL}}(D(\mathbf{x}) \| D(G(\mathbf{z})))] - \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(G(\mathbf{z})))], \quad (15)$$

$$(G_{\text{objective3}}) \quad \min_G \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\mathcal{D}_{\text{KL}}(\mathcal{A}_1 \| D(G(\mathbf{z})))] - \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(G(\mathbf{z})))]. \quad (16)$$

# Experiment

- Synthetic dataset
  - Real distribution: the mixture of nine normal distributions
  - Generator and discriminator are consists of four FC layers
  - Input latent from 32 dimensional normal distribution

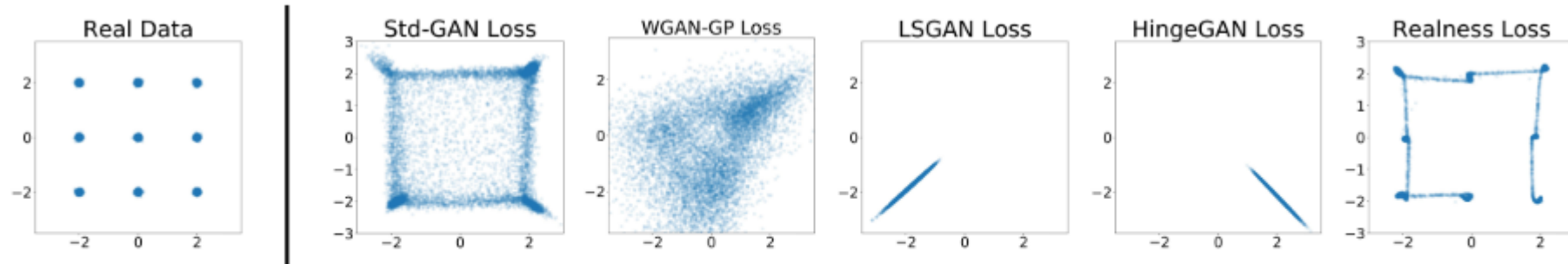
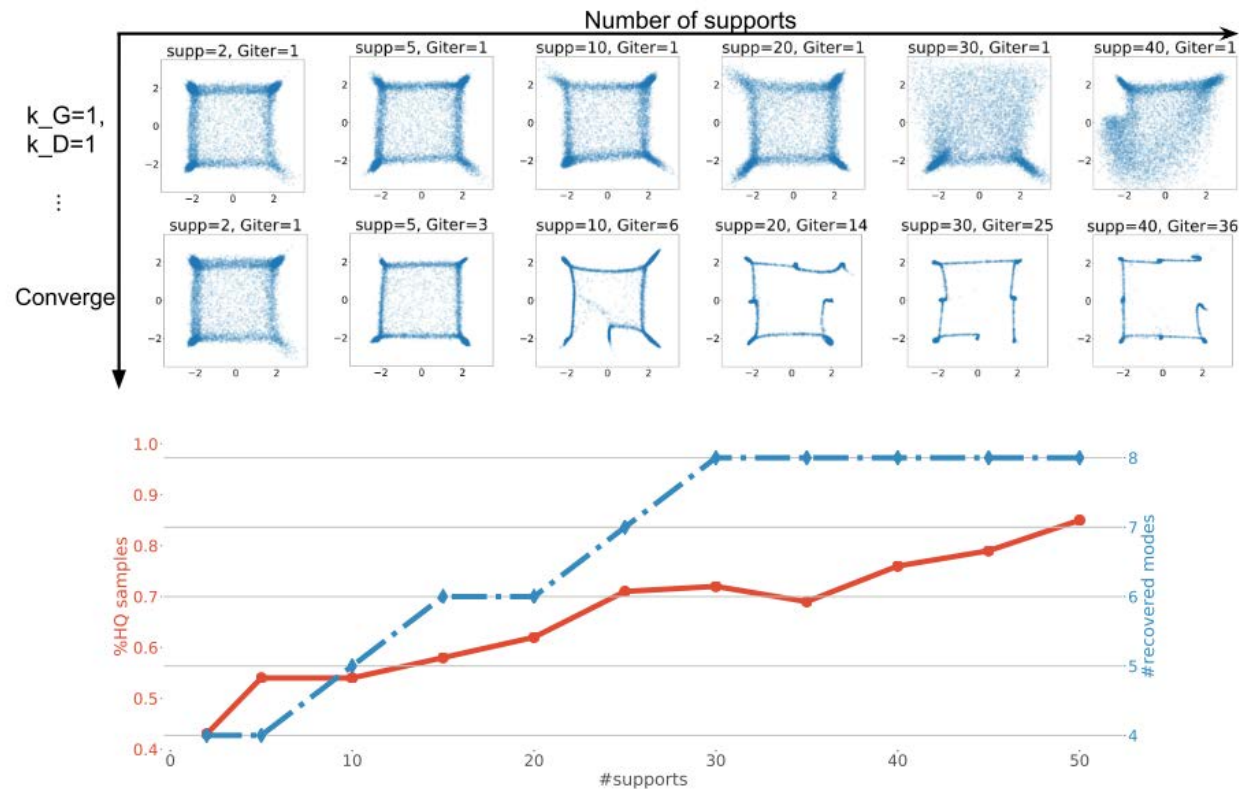


Figure 2: Left: real data sampled from the mixture of 9 Gaussian distributions. Right: samples generated by *Std-GAN*, *WGAN-GP*, *LSGAN*, *HingeGAN* and *RealnessGAN*.

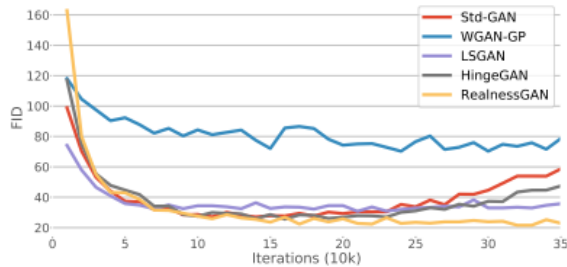
# Experiments

- Synthetic dataset
  - Effect of adjusting the number of outcomes (supports)

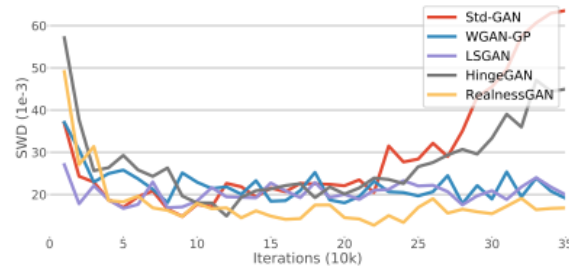


# Experiments

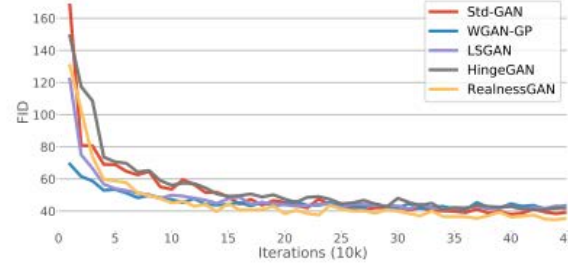
- Real-world datasets
  - CelebA, CIFAR-10 and FFHQ (32 by 32, 256 by 256 and 1,024 by 1,024)
  - Use DCGAN



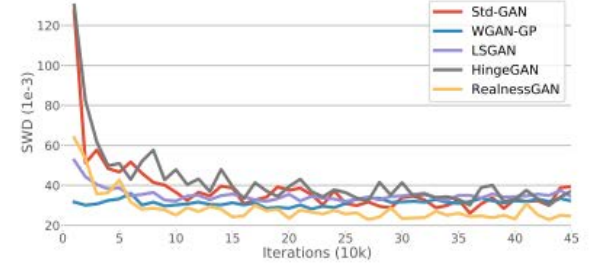
(a) FID on CelebA



(b) SWD on CelebA



(c) FID on CIFAR10



(d) SWD on CIFAR10

# Experiments

- Objective of G

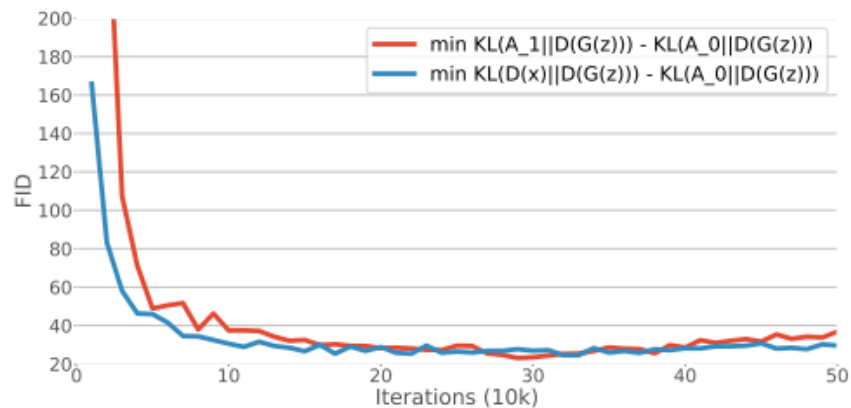
$$(G_{\text{objective1}}) \quad \min_G -\mathbb{E}_{z \sim p_z} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(G(z)))]. \quad (14)$$

$$(G_{\text{objective2}}) \quad \min_G \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}, z \sim p_z} [\mathcal{D}_{\text{KL}}(D(\mathbf{x}) \| D(G(z)))] - \mathbb{E}_{z \sim p_z} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(G(z)))], \quad (15)$$

$$(G_{\text{objective3}}) \quad \min_G \mathbb{E}_{z \sim p_z} [\mathcal{D}_{\text{KL}}(\mathcal{A}_1 \| D(G(z)))] - \mathbb{E}_{z \sim p_z} [\mathcal{D}_{\text{KL}}(\mathcal{A}_0 \| D(G(z)))]. \quad (16)$$

Table 3: In this table we compare different objectives of  $G$  on CIFAR10.

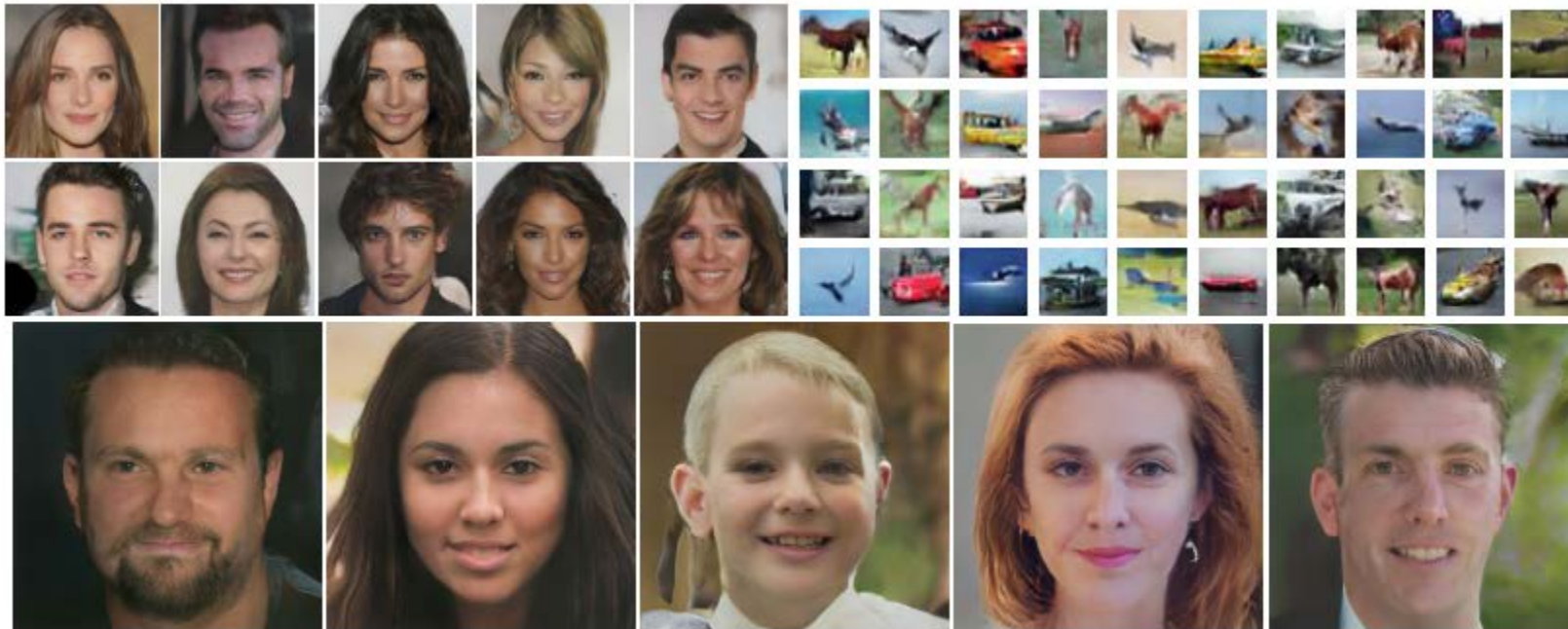
Objective	FID
Ideal Objective (equation 14)	36.73
Objective 2 (equation 15)	34.59
Objective 3 (equation 16)	36.21





# Experiments

- Real-world datasets
  - 1024 by 1024 images from DCGAN with RealnessGAN



# Experiments

We further compute FID as the quantitative result. Specifically, RealnessGAN yields a FID score of \*17.18\*. For reference, we also re-implement StyleGAN and train it using the same setting, resulting in a FID score of \*16.12\*.

- Real-world datasets
  - 1024 by 1024 images from DCGAN with RealnessGAN

Table 1: Minimum (min), maximum (max), mean and standard deviation (SD) of FID and SWD on CelebA and CIFAR10, calculated at 20k, 30k, ... iterations. The best indicators in baseline methods are underlined.

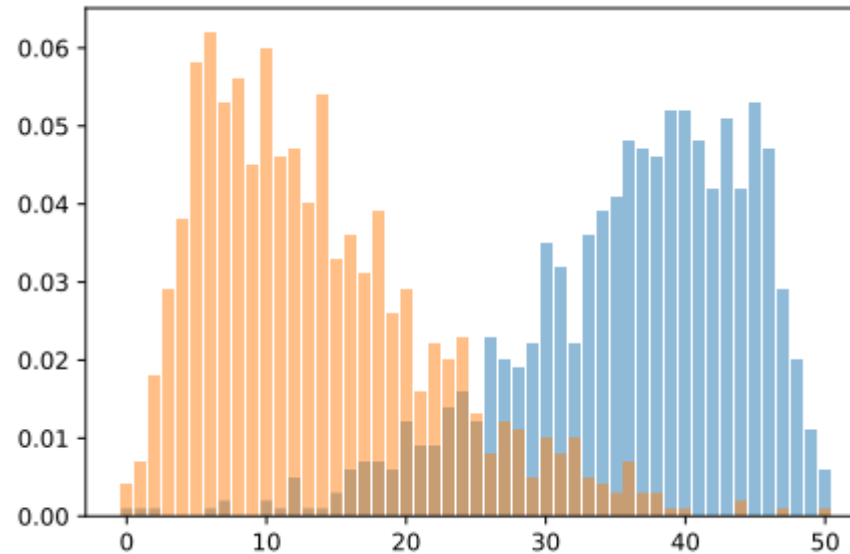
	Method	FID ↓				SWD ( $\times 10^3$ ) ↓			
		Min	Max	Mean	SD	Min	Max	Mean	SD
CelebA	Std-GAN	27.02	70.43	34.85	9.40	<u>14.81</u>	68.06	30.58	15.39
	WGAN-GP	70.28	104.60	81.15	8.27	<u>17.85</u>	30.56	22.09	2.93
	LSGAN	30.76	<u>57.97</u>	34.99	<u>5.15</u>	16.72	<u>23.99</u>	<u>20.39</u>	<u>2.25</u>
	HingeGAN	<u>25.57</u>	<u>75.03</u>	<u>33.89</u>	10.61	14.91	54.30	28.86	10.34
	RealnessGAN	<b>23.51</b>	81.3	<b>30.82</b>	7.61	<b>12.72</b>	31.39	<b>17.11</b>	3.59
CIFAR10	Std-GAN	<u>38.56</u>	88.68	47.46	15.96	28.76	57.71	37.55	7.02
	WGAN-GP	<u>41.86</u>	79.25	<u>46.96</u>	<u>5.57</u>	<u>28.17</u>	<u>36.04</u>	30.98	<u>1.78</u>
	LSGAN	42.01	<u>75.06</u>	48.41	7.72	<u>31.99</u>	40.46	<u>34.75</u>	2.34
	HingeGAN	42.40	117.49	57.30	20.69	32.18	61.74	41.85	7.31
	RealnessGAN	<b>34.59</b>	102.98	<b>42.30</b>	11.84	<b>22.80</b>	53.38	<b>26.98</b>	5.47

# Experiments



# Issues

- Selecting anchor distribution
  - In this paper, the author chose two skewed normal distribution



# Issues

- Multi-dimensional output VS multiple discriminator
  - Is it same as an ensemble of discriminators? (Ensemble GAN)
  - No!
- Conceptually, RealnessGAN and EnsembledGAN are orthogonal
- RealnessGAN could serve as one of the discriminators of EnsembleGAN
- Technically, EnsembleGAN uses multiple discriminators that could have different architectures and weights
- RealnessGAN uses a single discriminator
- EnsembleGAN using DCGAN fails on FFHQ

# Issues

- Role of each outcomes
  - 리버탈에 서술했지만, 아무 말 대단치
  - Future work

# Etc

- Code: will be uploaded
- Content in magenta
- Review score: 8, 3→6, 3→6