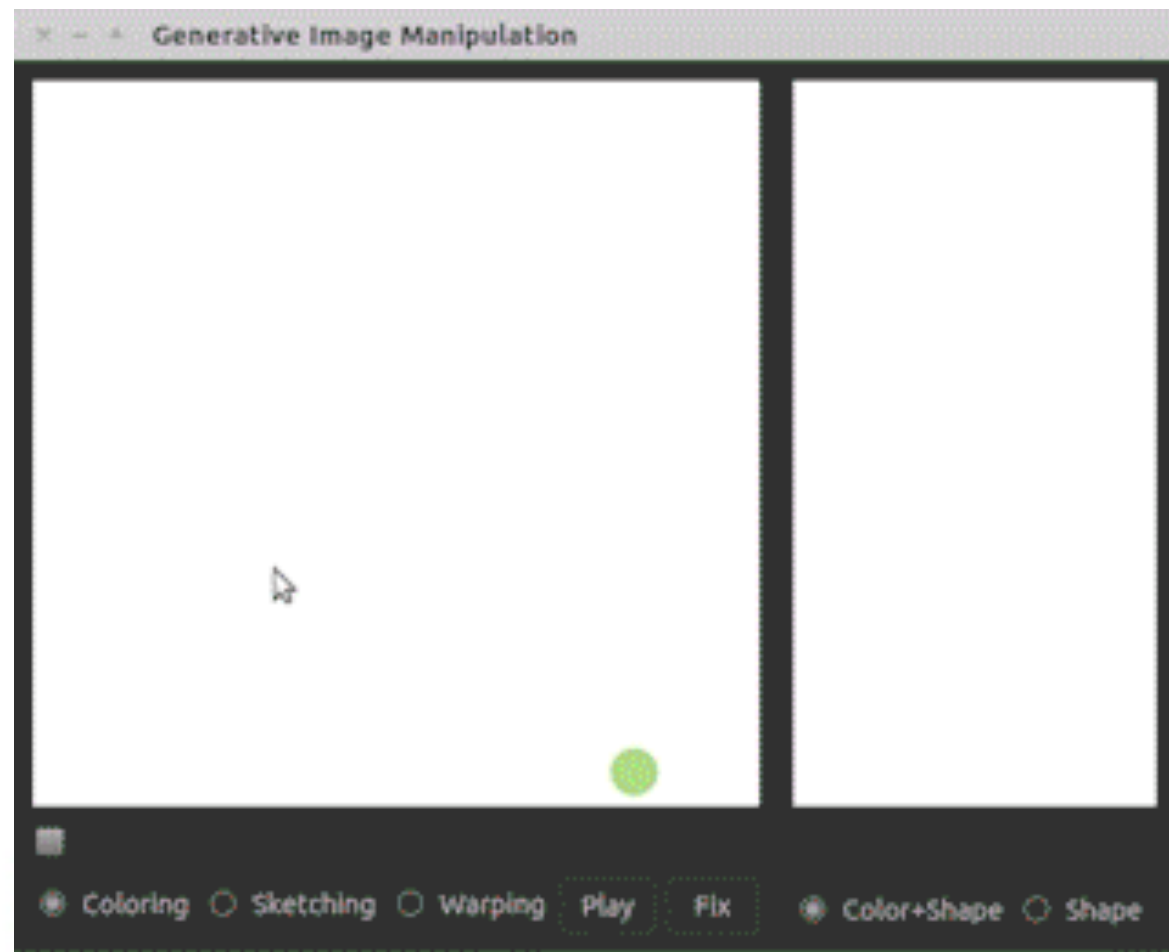


MedGAN ID-CGAN CoGAN LR-GAN CGAN IcGAN  
b-GAN LS-GAN AffGAN DiscoGAN AdaGAN  
LSGAN CatGAN LAPGAN iGAN IAN  
InfoGAN AMGAN MPM-GAN MIX+GAN  
McGAN **Head First Generative Adversarial Networks**  
C-RNN-GAN **From Theoretic View** DR-GAN  
MGAN GoGAN BS-GAN  
C-VAE-GAN FF-GAN  
3D-GAN CCGAN AC-GAN DCGAN  
GAWWN DualGAN BiGAN  
Bayesian GAN The Hong Kong Polytechnic University GP-GAN  
EBGAN Context-RNN-GAN AnoGAN DTN  
ALI MARTA-GAN ArtGAN f-GAN MAGAN MAD-GAN  
BEGAN MaIGAN AL-CGAN

獨釣寒江雪  
孤舟蓑笠翁  
萬徑人踪滅  
千山鳥飛絕

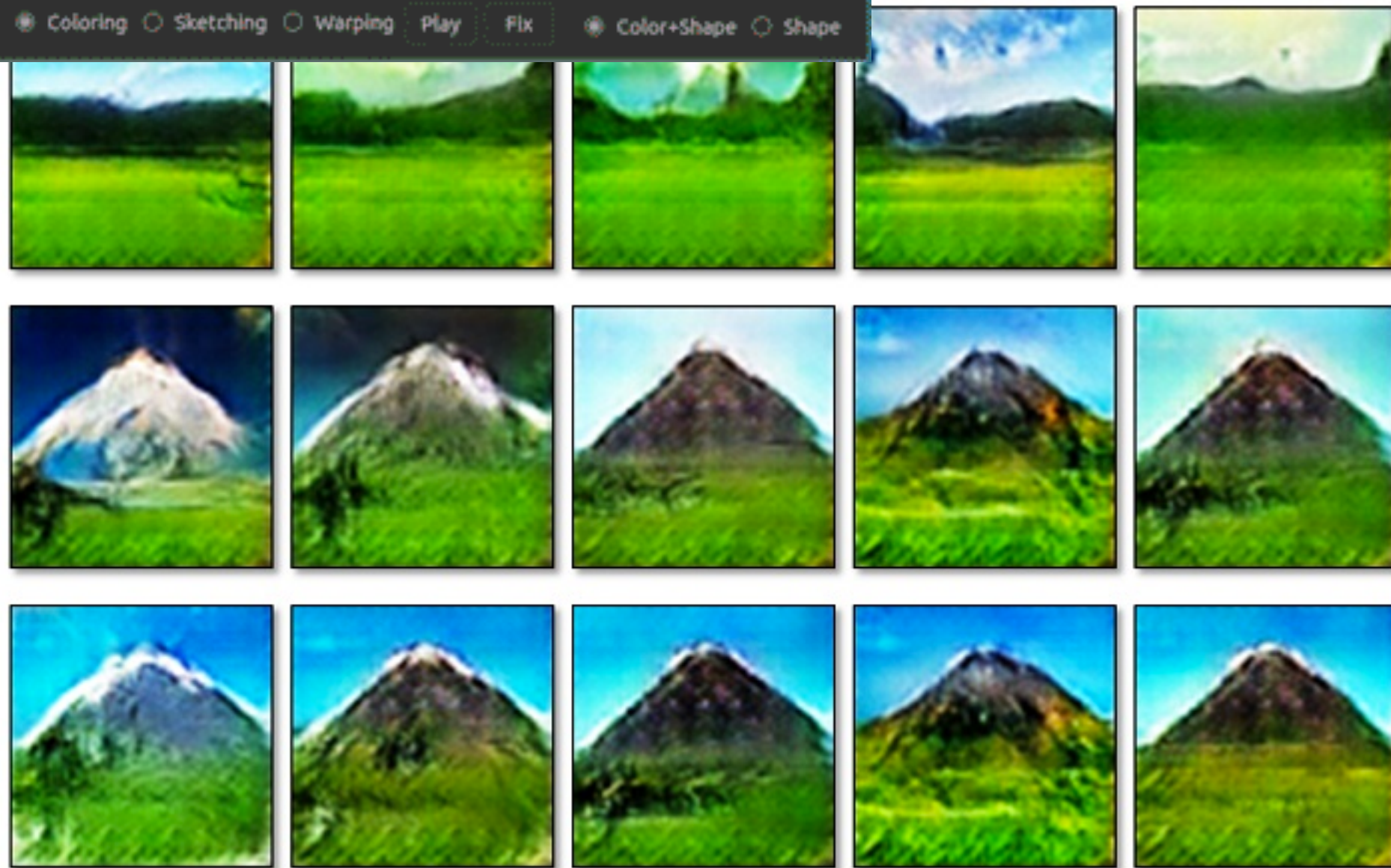
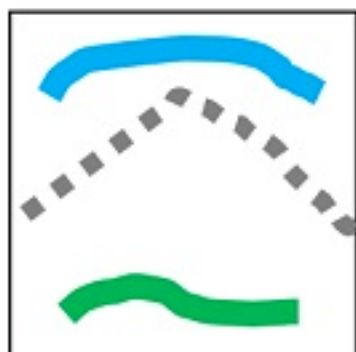
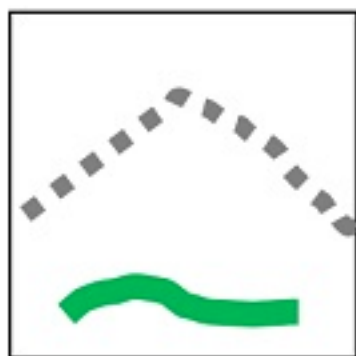
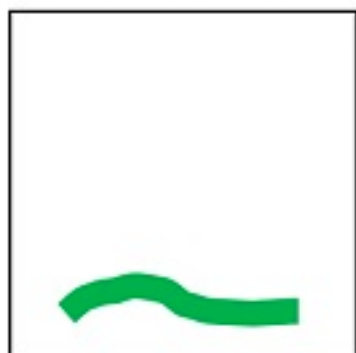
久	久	腸	腸	柞	柞	裘	裘	肋	肋
亮	亮	樞	樞	叉	叉	哲	哲	苗	苗
閑	閑	洧	洧	材	材	簿	簿	裝	裝
序	序	魁	魁	坤	坤	瞬	瞬	晉	晉
骷	骷	熟	熟	又	又	吹	吹	瀾	瀾
矜	矜	綢	綢	弘	弘	取	取	嘆	嘆
齧	齧	涉	涉	否	否	椎	椎	訟	訟
餓	餓	蚕	蚕	湖	湖	陰	陰	涂	涂
汨	汨	丞	丞	铎	铎	瑗	瑗	拉	拉
緯	緯	蒯	蒯	啼	啼	燠	燠	虾	虾
祿	祿	糸	糸	閑	閑	厠	厠	荏	荏
特	特	鉏	鉏	祝	祝	敦	敦	鏃	鏃
犖	犖	亏	亏	话	话	ッ	ッ	瑾	瑾
吗	吗	鏹	鏹	繡	繡	髡	髡	徒	徒
乳	乳	磁	磁	球	球	脍	脍	隙	隙
蛩	蛩	節	節	林	林	情	情	翎	翎
在	在	泥	泥	蚊	蚊	努	努	睨	睨
ぢ	ぢ	鳶	鳶	鵠	鵠	筆	筆	瑋	瑋





(<https://github.com/junyanz/iGAN>)

User edits





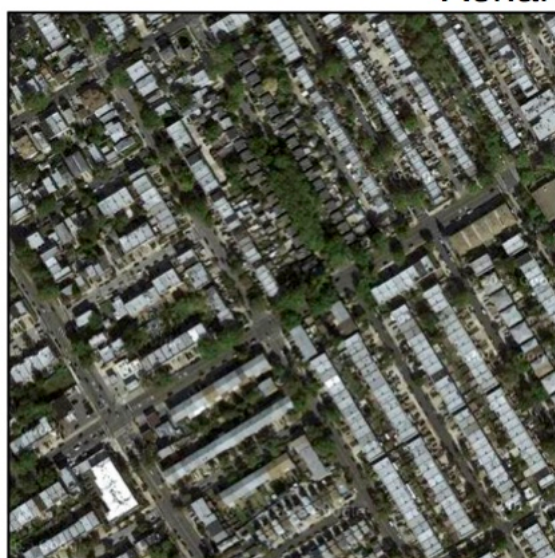


input



output

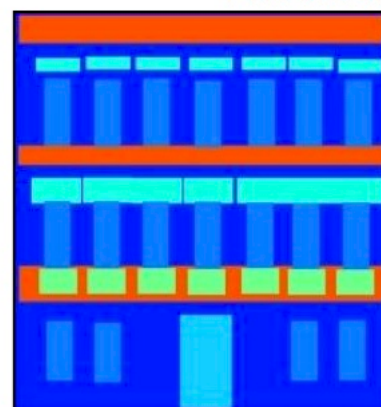
Aerial to Map



input



output

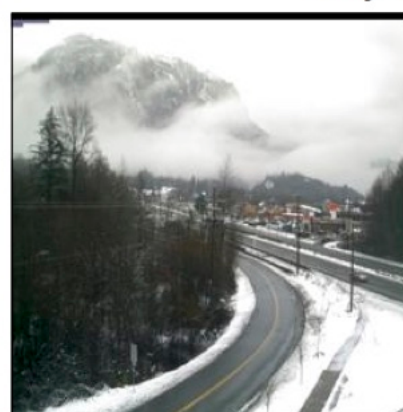


input



output

Day to Night



input



output



input

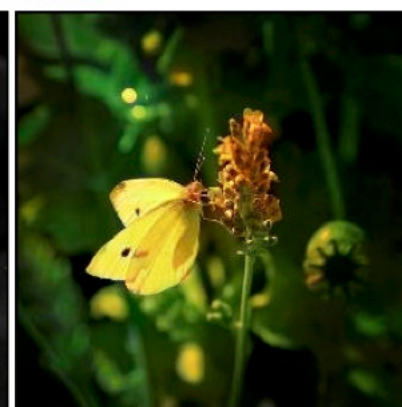


output

Edges to Photo



input



output









(<https://junyanz.github.io/CycleGAN/>)

Monet ↔ Photos



Monet → photo



photo → Monet

Zebras ↔ Horses



zebra → horse



horse → zebra

Summer ↔ Winter



summer → winter



winter → summer



Photograph



Monet



Van Gogh



Cezanne



Ukiyo-e





(genekogan@Twitter)

# Content

- Generative Adversarial Networks
  - Basics and Attractiveness
  - Difficulties
- Solution 1: Partial and Fine-grained Guidance
- Solution 2: Encoder-incorporated
- Solution 3: Wasserstein Distance

# Content

- Generative Adversarial Networks
  - Basics and Attractiveness
  - Difficulties
- Solution 1: Partial and Fine-grained Guidance
- Solution 2: Encoder-incorporated
- Solution 3: Wasserstein Distance

# Generative Adversarial Networks

0.6551  
(RNG)



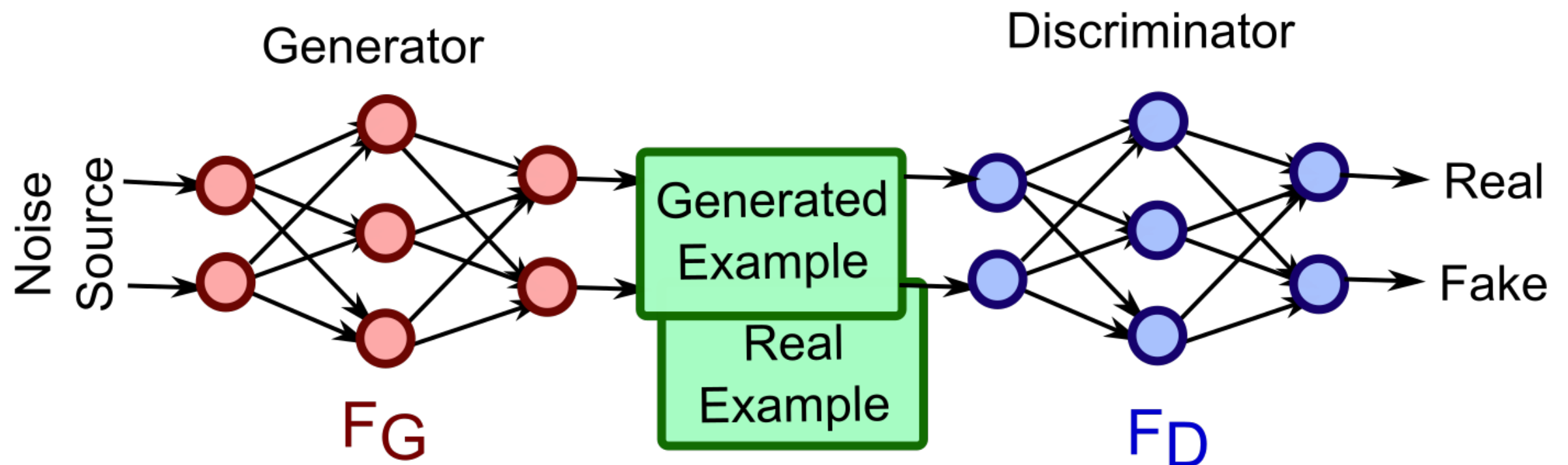


# Generative Adversarial Networks

- A **counterfeiter-police game** between two components: a generator ***G*** and a discriminator ***D***
- ***G***: counterfeiter, trying to fool police with fake currency
- ***D***: policy, trying to detect the counterfeit currency
- Competition drives both to improve, until counterfeits are *indistinguishable* from genuine currency

# Generative Adversarial Networks

- A **min-max game** between two components: a generator  $G$  and a discriminator  $D$



# Generative Adversarial Networks

- A **min-max game** between two components: a generator ***G*** and a discriminator ***D***

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

D predicting that  
real data is genuine

D predicting that  
G's generated data is fake

# Generative Adversarial Networks

- A **min-max game** between two components: a generator ***G*** and a discriminator ***D***

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

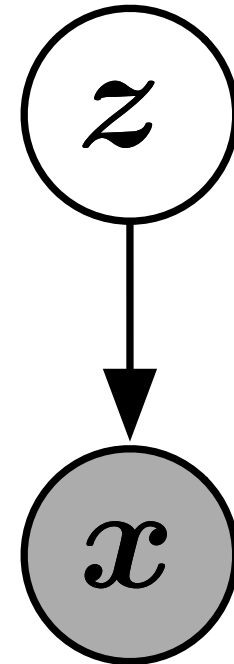
*D* predicting that  
real data is genuine

*D* predicting that  
*G*'s generated data is fake

- ***D***'s goal: maximize  $V(D, G)$   
***G***'s goal: minimize  $\max V(D, G)$



# Attractiveness



- Generator Networks  $x = G(z; \theta^{(G)})$
- It is only required that,  $G$  is differentiable.
- So, having training data  $x \sim p_{\text{data}}(x)$   
what we want is a model that can draw samples  $x \sim p_{\text{model}}(x)$ , where  $p_{\text{model}} \approx p_{\text{data}}$
- **Don't write a formula for  $p_{\text{data}}(x)$** , just learn to draw sample directly.

“There’s no free lunch.”

–From Economics



Original



Generated



Original



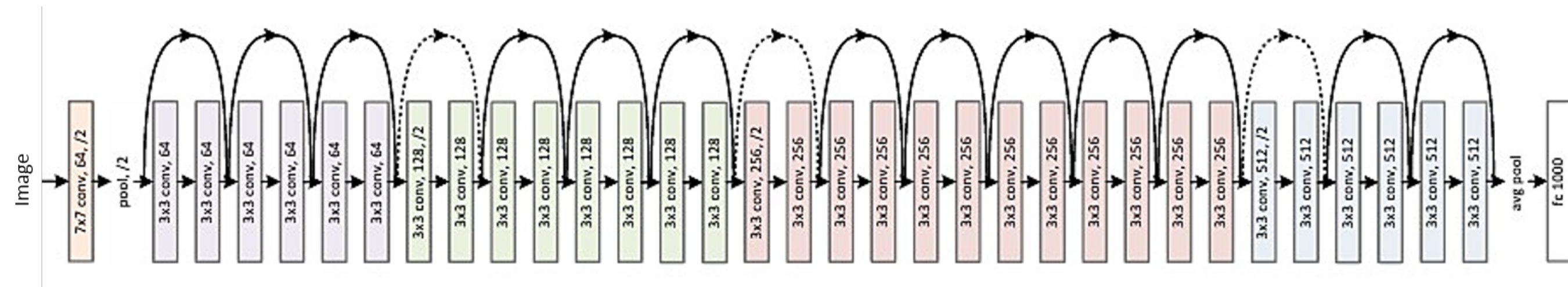
Generated





# Difficulty 1

- The gradient issues existed in deep neural networks
- The deeper, the more difficult



# Objectives for GAN

- The objective of ***D***:

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

- The objective of ***G***:

- the original:  $\mathbb{E}_{z \sim p(z)} [\log(1 - D(g_\theta(z)))]$

- the alternative:  $\mathbb{E}_{z \sim p(z)} [-\log D(g_\theta(z))]$

- *Why alternative?*

# Difficulty 2

- using the original form of the objective of ***G***

$$\mathbb{E}_{z \sim p(z)} [\log(1 - D(g_{\theta}(z)))]$$

will result in gradient vanishing issue of ***D*** for ***G*** because *intuitively*, at the very early phase of training, ***D*** is very easy to be confident in detecting ***G***, so ***D*** will output almost always 0

# Difficulty 2

- using the original form of the objective of ***G***

$$\mathbb{E}_{z \sim p(z)} [\log(1 - D(g_\theta(z)))]$$

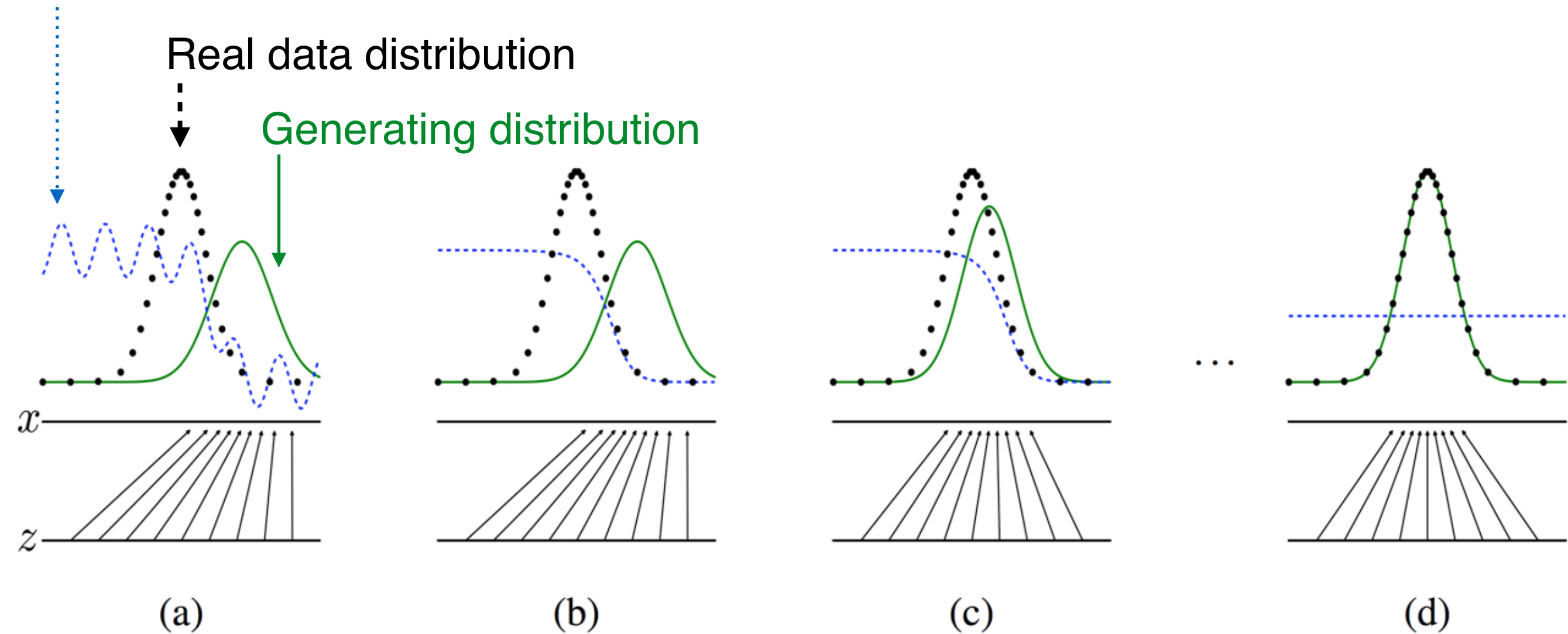
will result in gradient vanishing issue of ***D*** for ***G*** because *theoretically*, when ***D*** is *optimal*, minimizing the loss is equal to minimizing the *JS divergence* (Arjovsky & Bottou, 2017)

# Difficulty 2

Discriminative distribution

Real data distribution

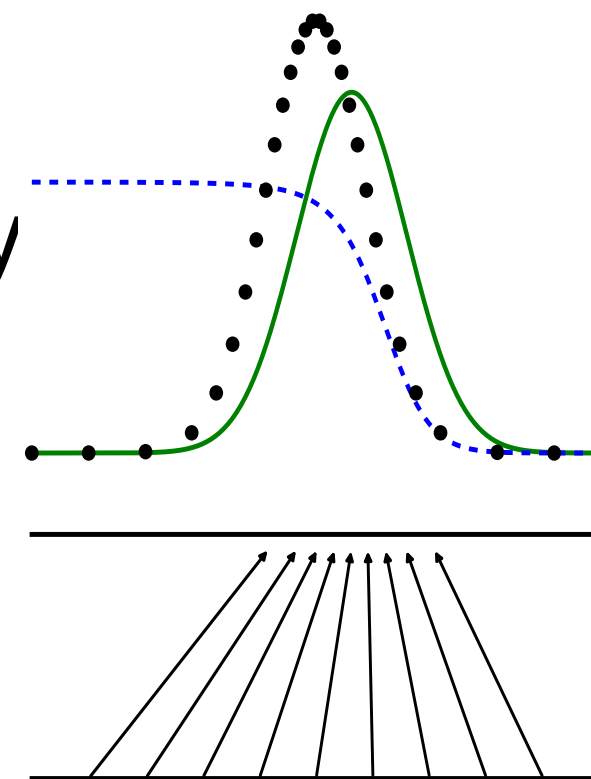
Generating distribution



# Difficulty 2

- The optimal  $D$  for any  $P_r$  and  $P_g$  is always

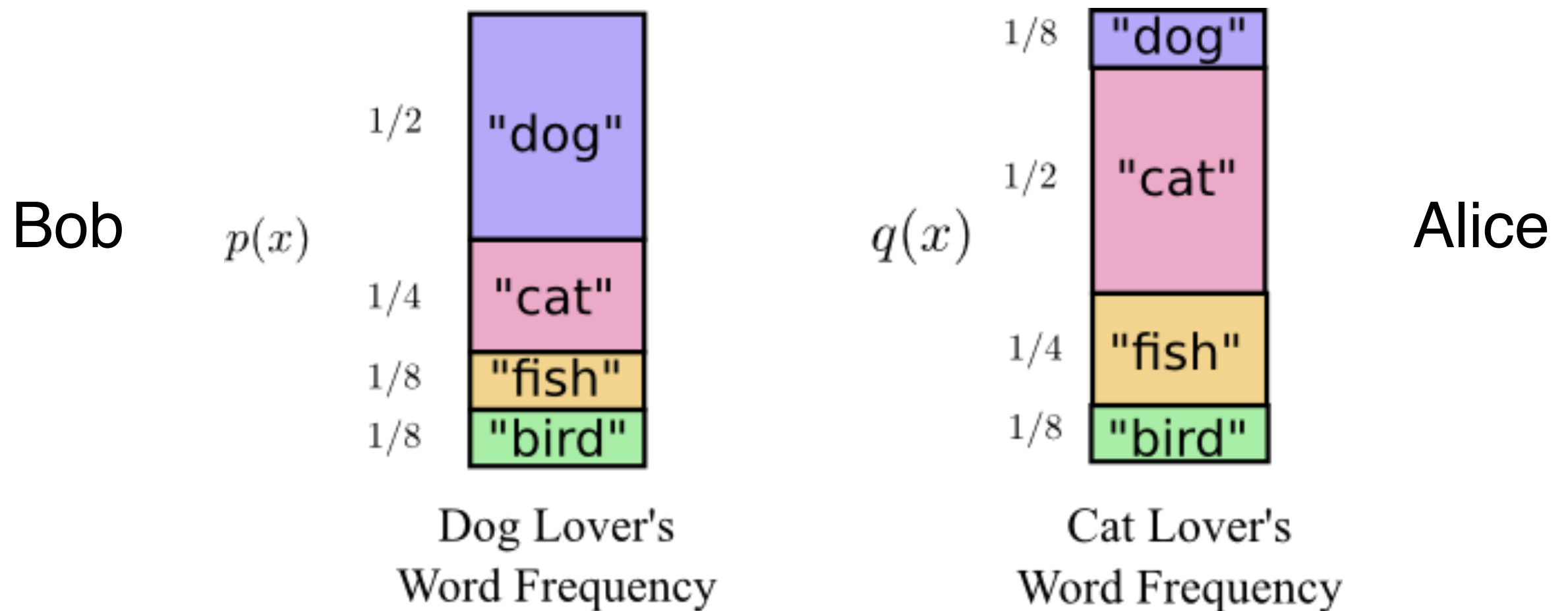
$$D^*(x) = \frac{P_r(x)}{P_r(x) + P_g(x)}$$



and that  $L(D^*, g_\theta) = 2JSD(\mathbb{P}_r || \mathbb{P}_g) - 2 \log 2$

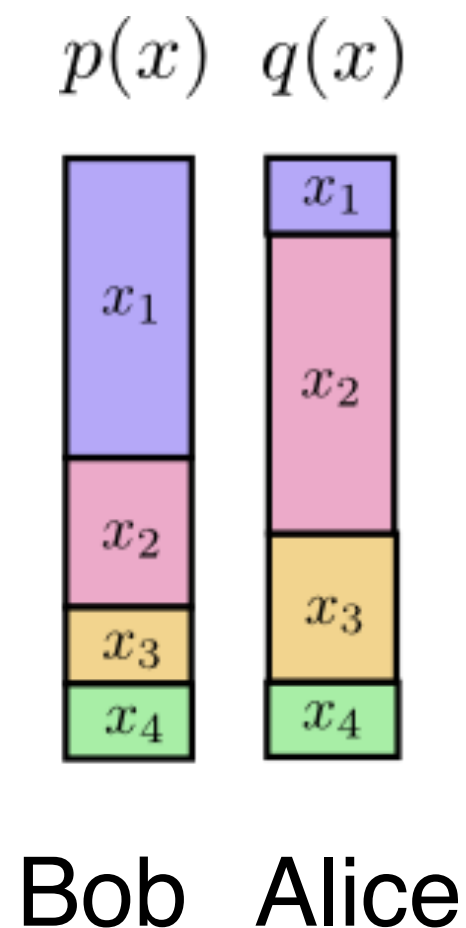
so, when  $\mathbf{D}$  is *optimal*, minimizing the loss is equal to minimizing the *JS divergence* (Arjovsky & Bottou, 2017)

# Recall KL and JS Divergence





# Recall KL and JS Divergence

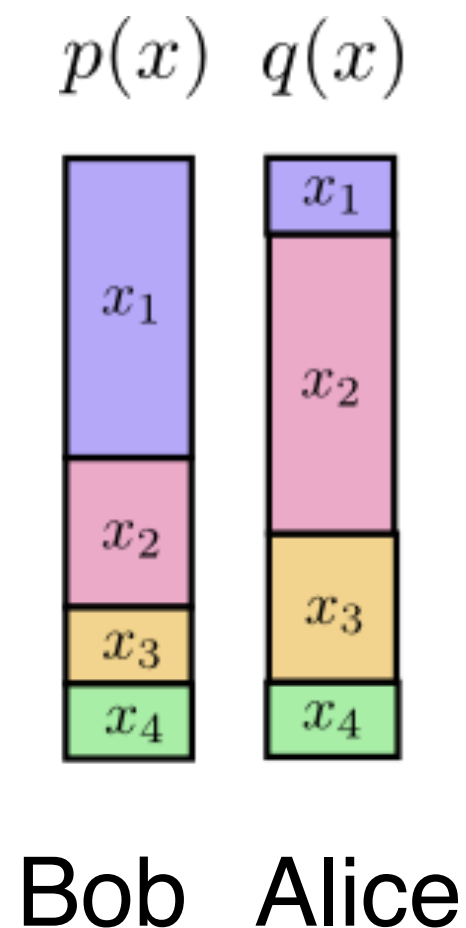


Cross-Entropy:  $H_p(q)$

Average Length  
of message from  $q(x)$   
using code for  $p(x)$ .

$$H_p(q) = \sum_x q(x) \log_2 \left( \frac{1}{p(x)} \right)$$

# Recall KL and JS Divergence

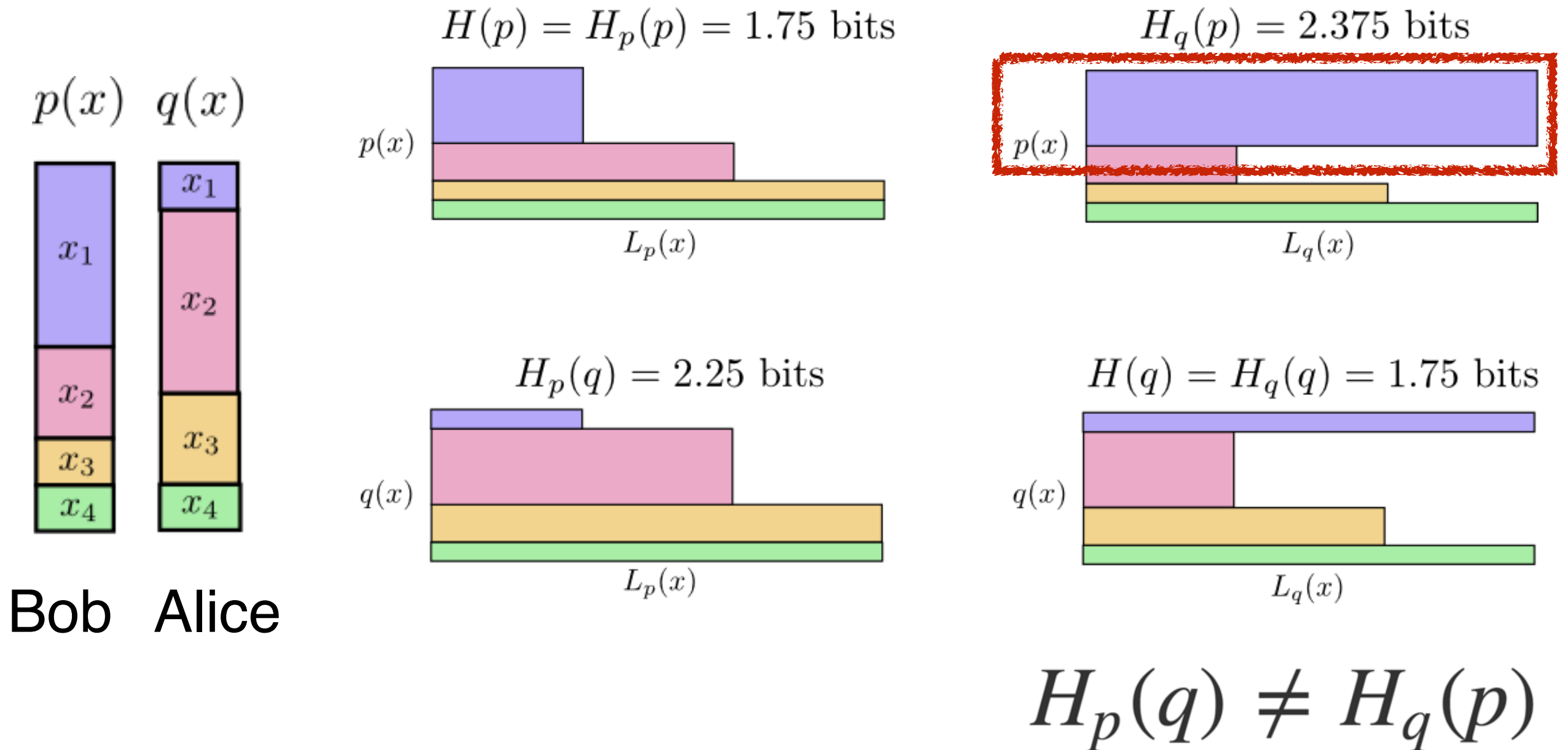


Cross-Entropy:  $H_p(q)$

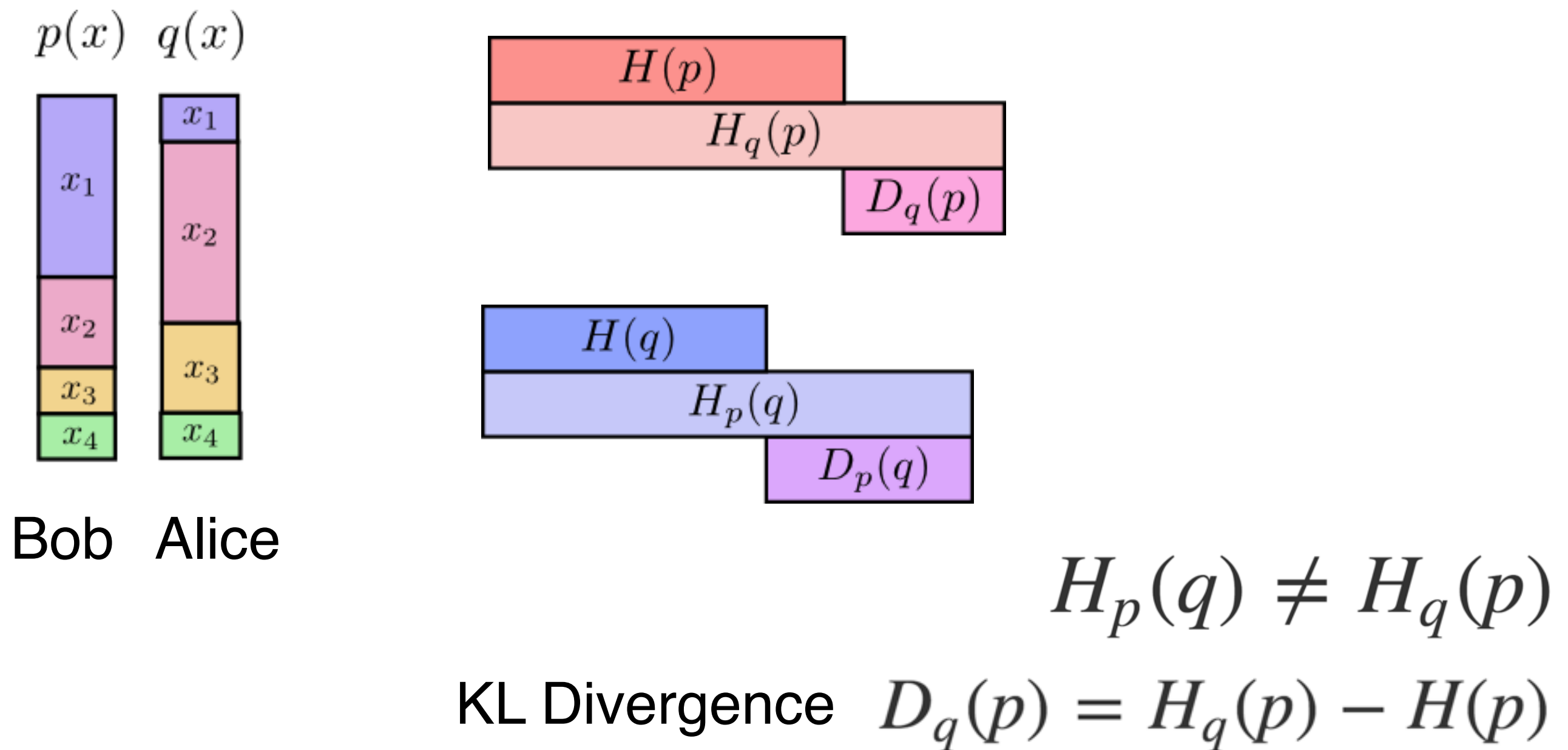
Average Length  
of message from  $q(x)$   
using code for  $p(x)$ .

$$H_p(q) = \sum_x q(x) \log_2 \left( \frac{1}{p(x)} \right)$$

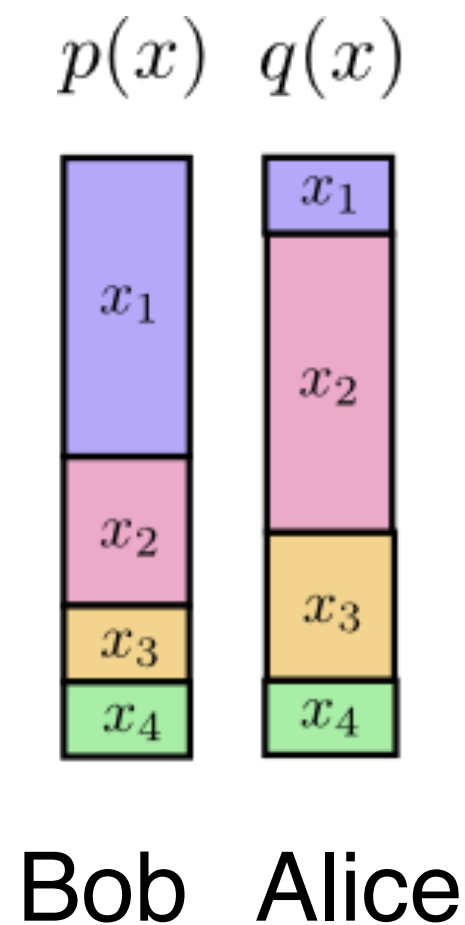
# Recall KL and JS Divergence



# Recall KL and JS Divergence



# Recall KL and JS Divergence



JS Divergence

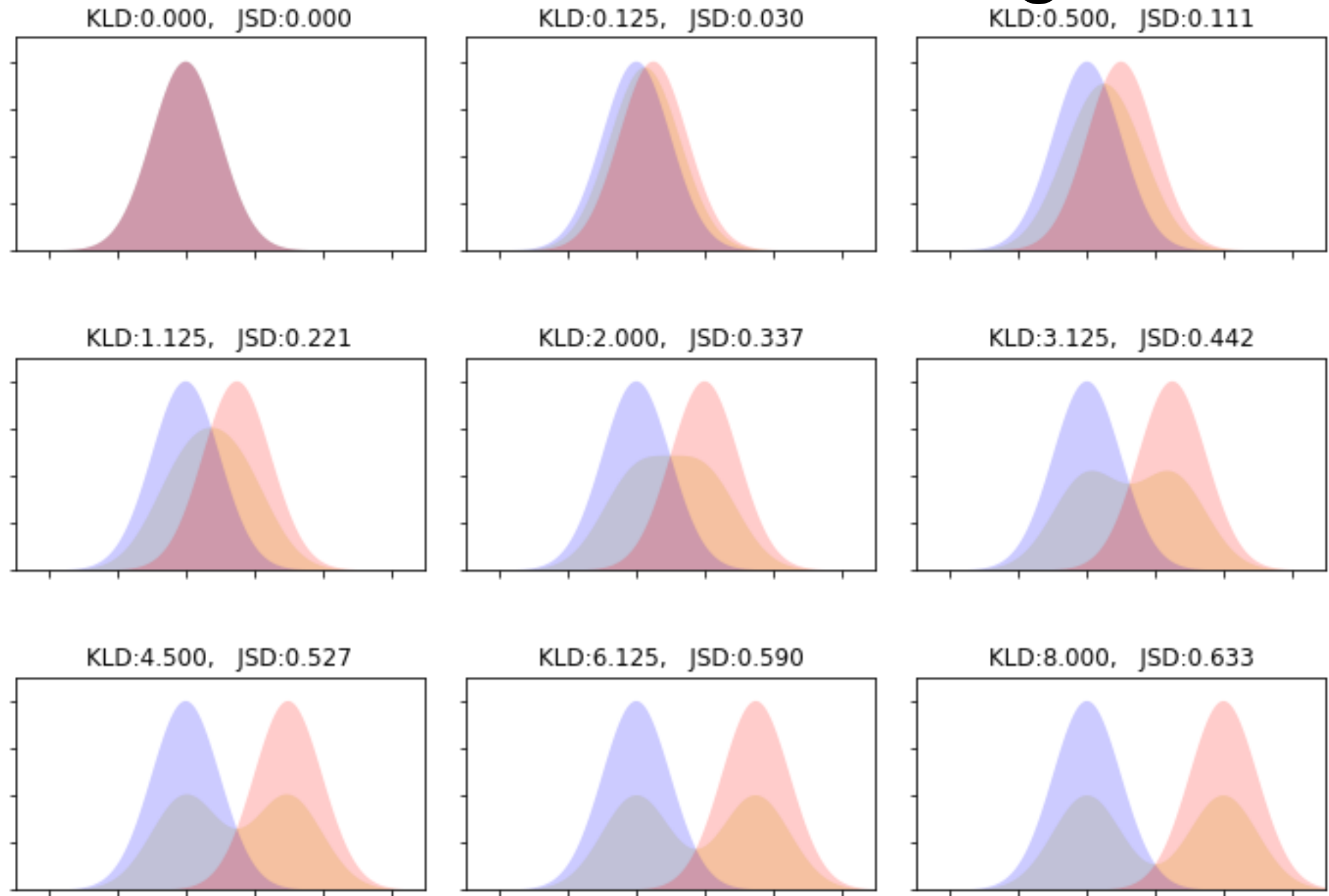
$$D_{JS}(p|q) = D_{JS}(q|p) = \frac{1}{2}D_{KL}(p|r) + \frac{1}{2}D_{KL}(q|r)$$
$$r = \frac{1}{2}(p + q)$$

Be symmetric!

$$H_p(q) \neq H_q(p)$$

KL Divergence  $D_q(p) = H_q(p) - H(p)$

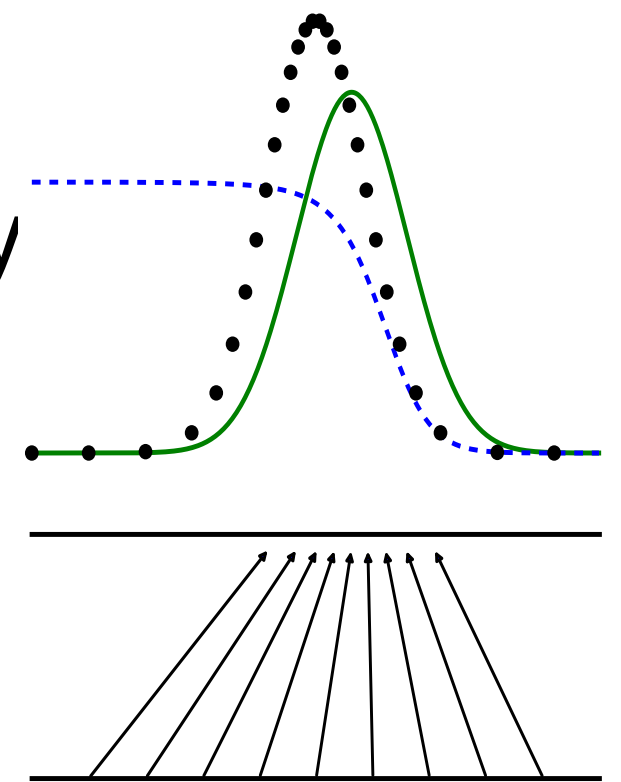
# Recall KL and JS Divergence



# Difficulty 2

- The optimal  $D$  for any  $P_r$  and  $P_g$  is always

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



and that  $L(D^*, g_\theta) = 2JSD(\mathbb{P}_r \parallel \mathbb{P}_g) - 2 \log 2$

so, when  $D$  is *optimal*, minimizing the loss is equal to minimizing the *JS divergence* (Goodfellow et al., 2014)



# Difficulty 2

- when:

$$L(D^*, g_\theta) = 2JSD(\mathbb{P}_r \parallel \mathbb{P}_g) - 2 \log 2.$$

- The JS divergence for the two distributions  $P_r$  and  $P_g$  is (almost) always  $\log 2$  because  $P_r$  and  $P_g$  hardly can overlap (Arjovsky & Bottou, 2017, Theorem 2.1~2.3)
- This results in vanishing gradient in theory!

# The alternative objective

- The alternative objective of ***G***:

$$\mathbb{E}_{z \sim p(z)} [-\log D(g_\theta(z))]$$

- Instead of minimizing, let ***G*** maximize the log-probability of the discriminator being mistaken
- It is heuristically motivated that generator can still learn even when discriminator successfully rejects all generator samples, but not theoretically guaranteed

# Difficulty 3

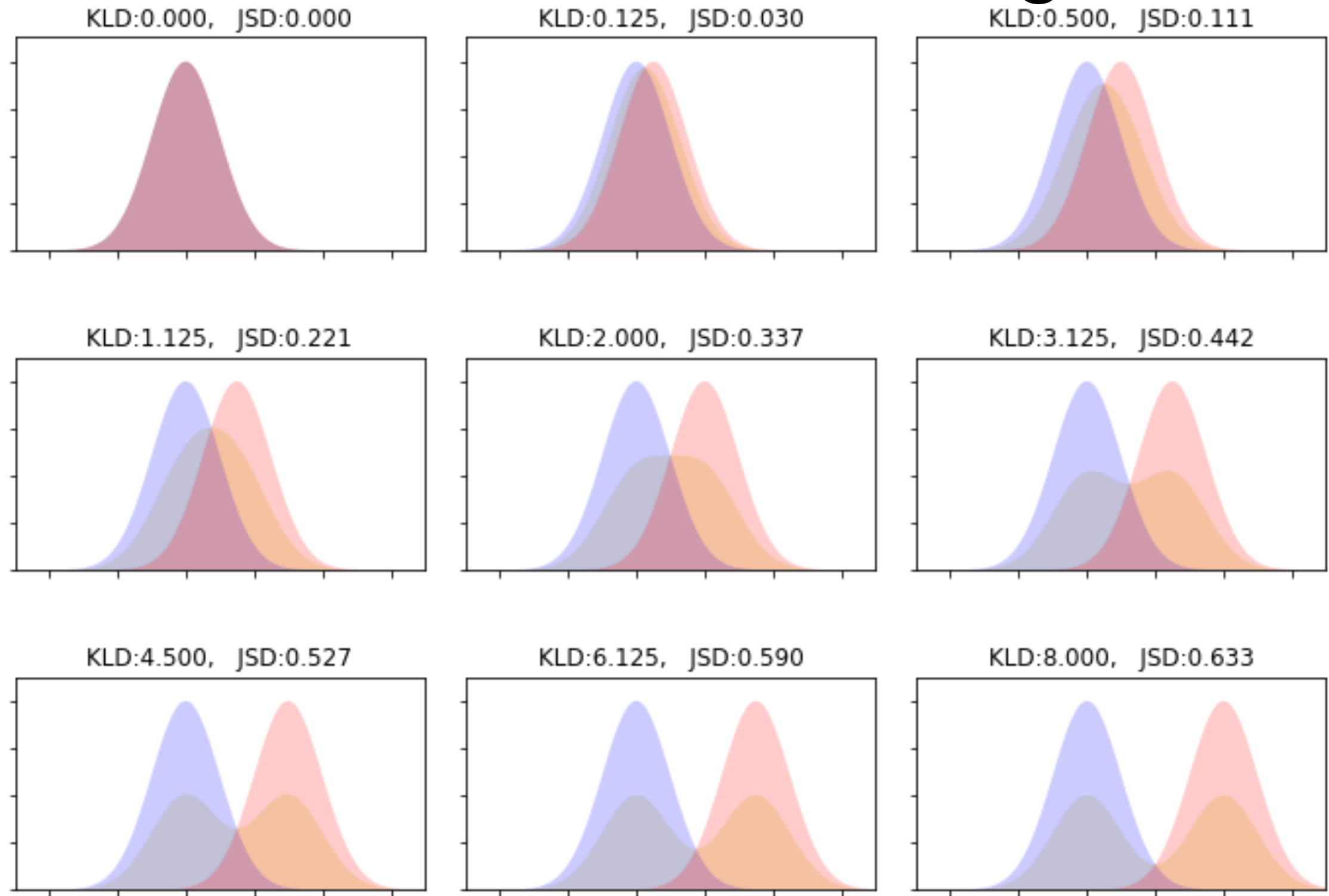
- using the alternative form of the objective of ***G***

$$\mathbb{E}_{z \sim p(z)} [-\log D(g_\theta(z))]$$

will result in gradient unstable issue and mode missing problem because *theoretically*, when ***D*** is *optimal*, minimizing the loss is equal to **minimizing** the *KL divergence* meanwhile **maximizing** the *JS divergence* (Arjovsky & Bottou, 2017, Theorem 2.5):

$$KL(\mathbb{P}_{g_\theta} || \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_\theta} || \mathbb{P}_r)]$$

# Recall KL and JS Divergence



# Difficulty 3

- **minimizing** the *KL divergence* meanwhile **maximizing** the *JS divergence* is crazy:

$$KL(\mathbb{P}_{g_\theta} || \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_\theta} || \mathbb{P}_r)]$$

- which results in gradient unstable issue

# Difficulty 3

- **minimizing** the *KL divergence* only **is biased**:

$$KL(\mathbb{P}_{g_\theta} || \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_\theta} || \mathbb{P}_r)]$$

- because *KL divergence* is asymmetric, and thus it is not equally treated when **G** generates an unreal sample and when **G** fails to generate real sample
- Therefore, **G** will generate too many few-mode (less diverse) but real samples , a safer strategy



# Content

- Generative Adversarial Networks
  - Basics and Attractiveness
  - Difficulties
- **Solution 1: Partial and Fine-grained Guidance**
- Solution 2: Encoder-incorporated
- Solution 3: Wasserstein Distance

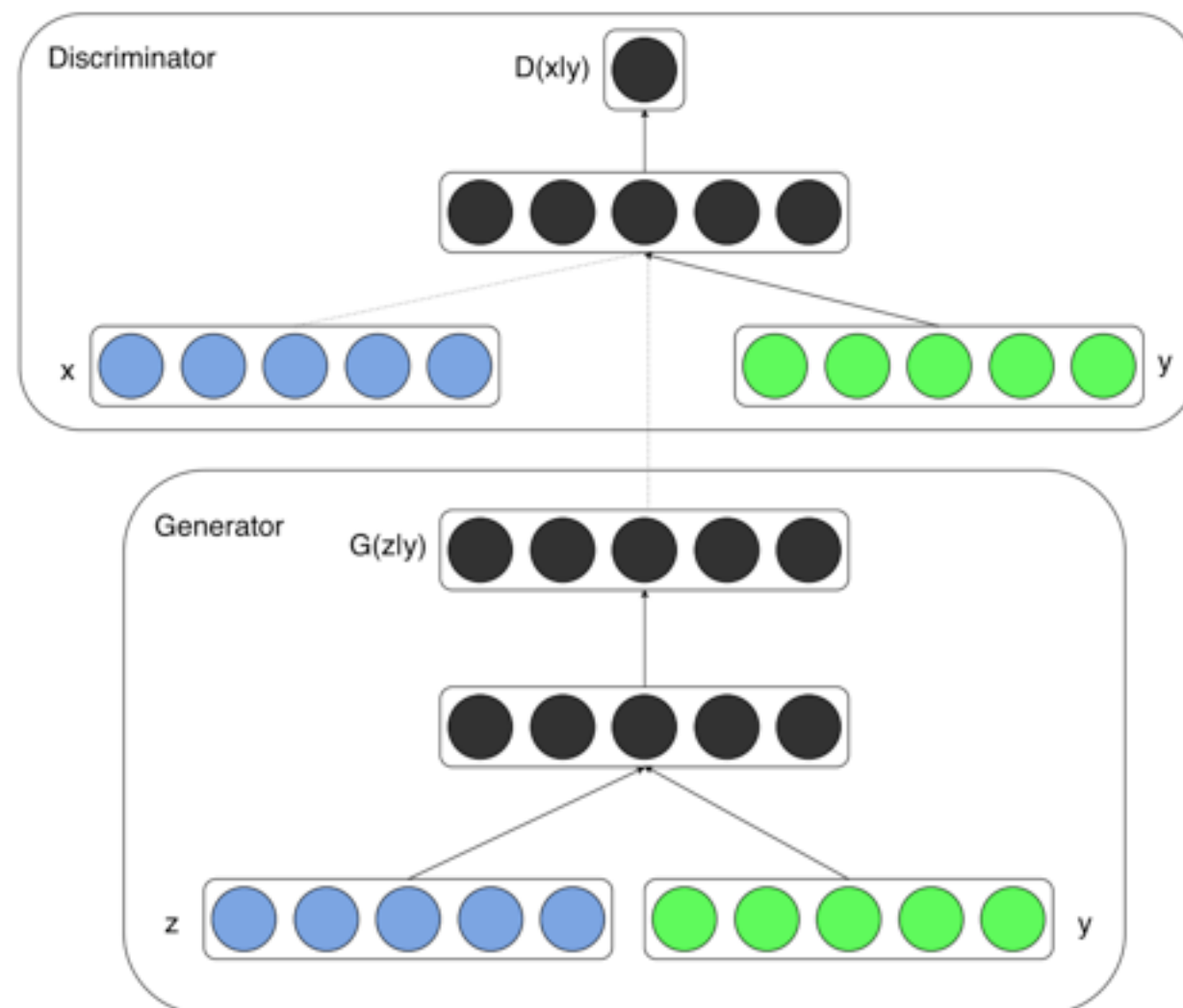
# Solution 1.1: Partial Guidance

- Conditional GANs (Mirza & Osindero, 2014)
- Improved GAN (Salimans et al., 2016)
- iGAN/GVM (Zhu et al., 2016)
- pix2pix (Isola et al., 2017)
- GP-GAN (Wu et al., 2017)

# Solution 1.1: Partial Guidance

- Conditional GANs (Mirza & Osindero, 2014)

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





# Solution 1.1: Partial Guidance

- Improved GAN (Salimans et al., 2016)

- feature matching

$$||\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \mathbf{f}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} \mathbf{f}(G(\mathbf{z}))||_2^2$$

- minibatch discrimination

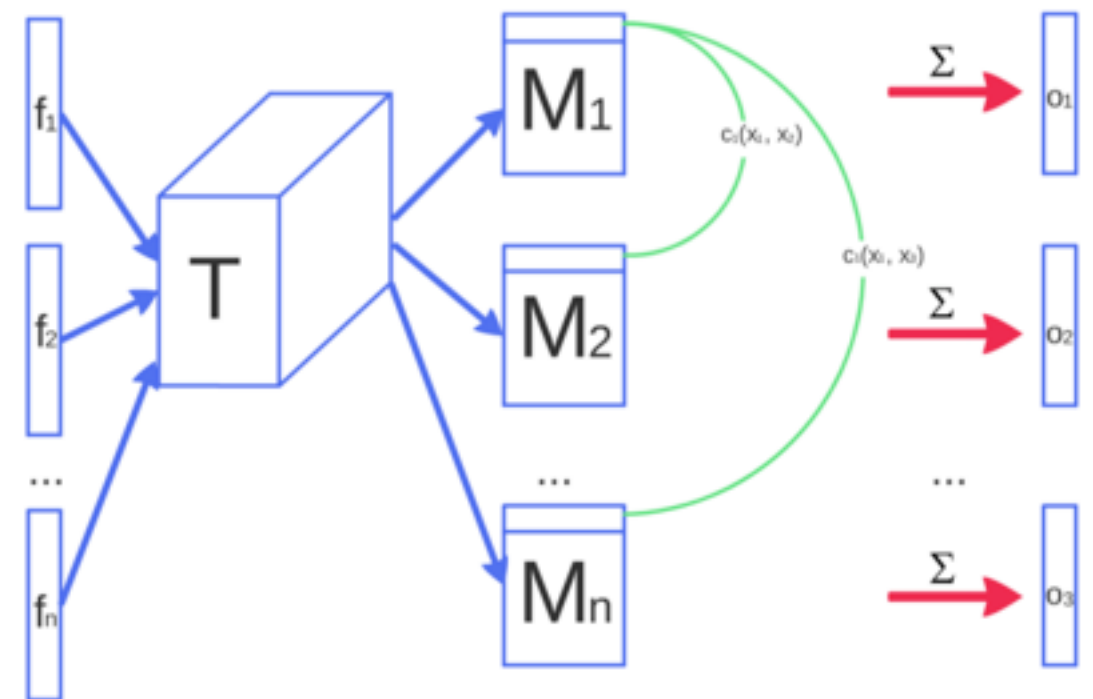


Figure 1: Figure sketches how minibatch discrimination works. Features  $\mathbf{f}(\mathbf{x}_i)$  from sample  $\mathbf{x}_i$  are multiplied through a tensor  $T$ , and cross-sample distance is computed.

# Solution 1.1: Partial Guidance

- iGAN/GVM (Zhu et al., 2016)



(a) original photo



(b) projection on manifold



(e) different degree of image manipulation



Edit Transfer



(d) smooth transition between the original and edited projection

# Solution 1.1: Partial Guidance

- iGAN/GVM (Zhu et al., 2016)



(a) User constraints  $v_g$  at different update steps



$G(z_0)$

(b) Updated images according to user edits

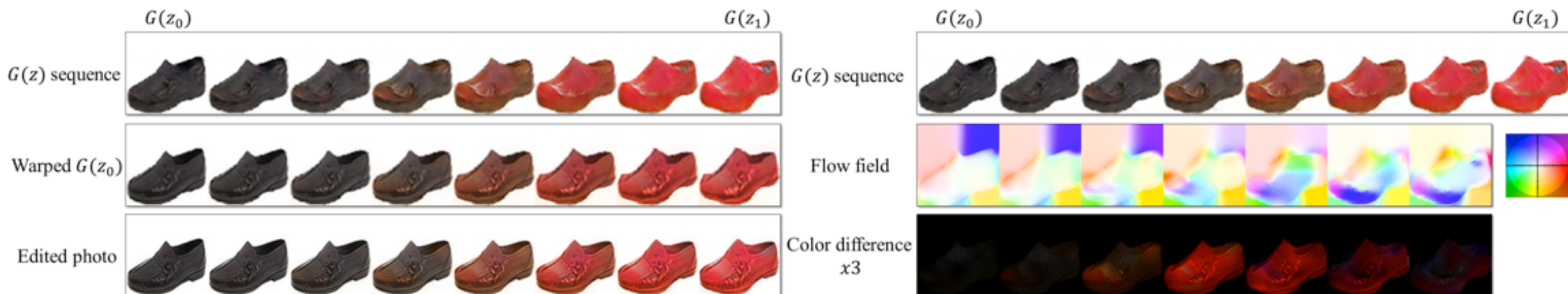
$G(z_1)$



(c) Linear interpolation between  $G(z_0)$  and  $G(z_1)$

# Solution 1.1: Partial Guidance

- iGAN/GVM (Zhu et al., 2016)



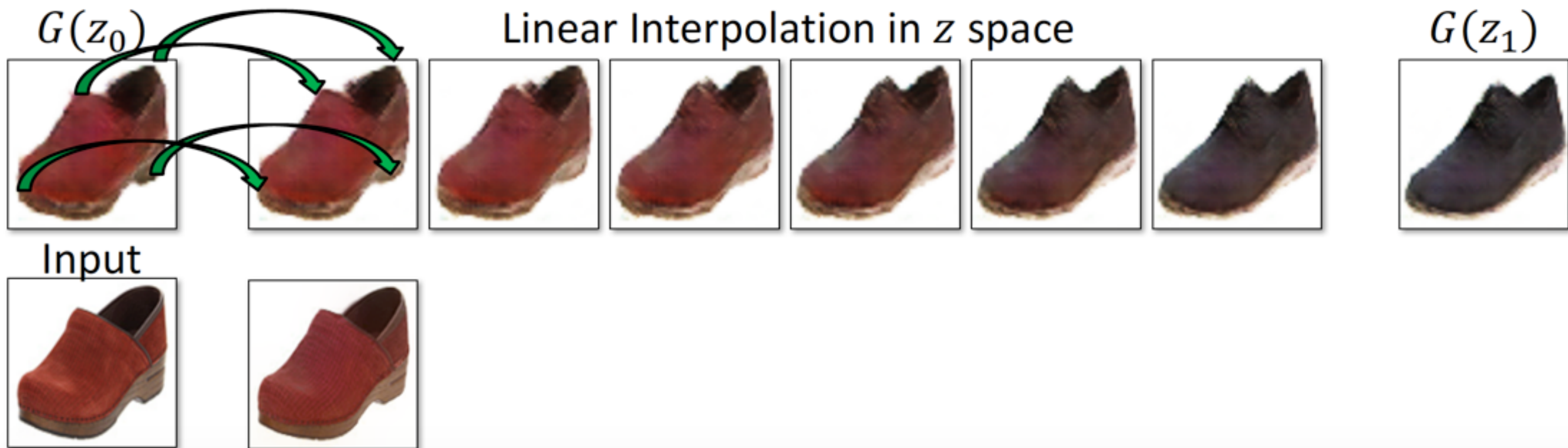


# Solution 1.1: Partial Guidance

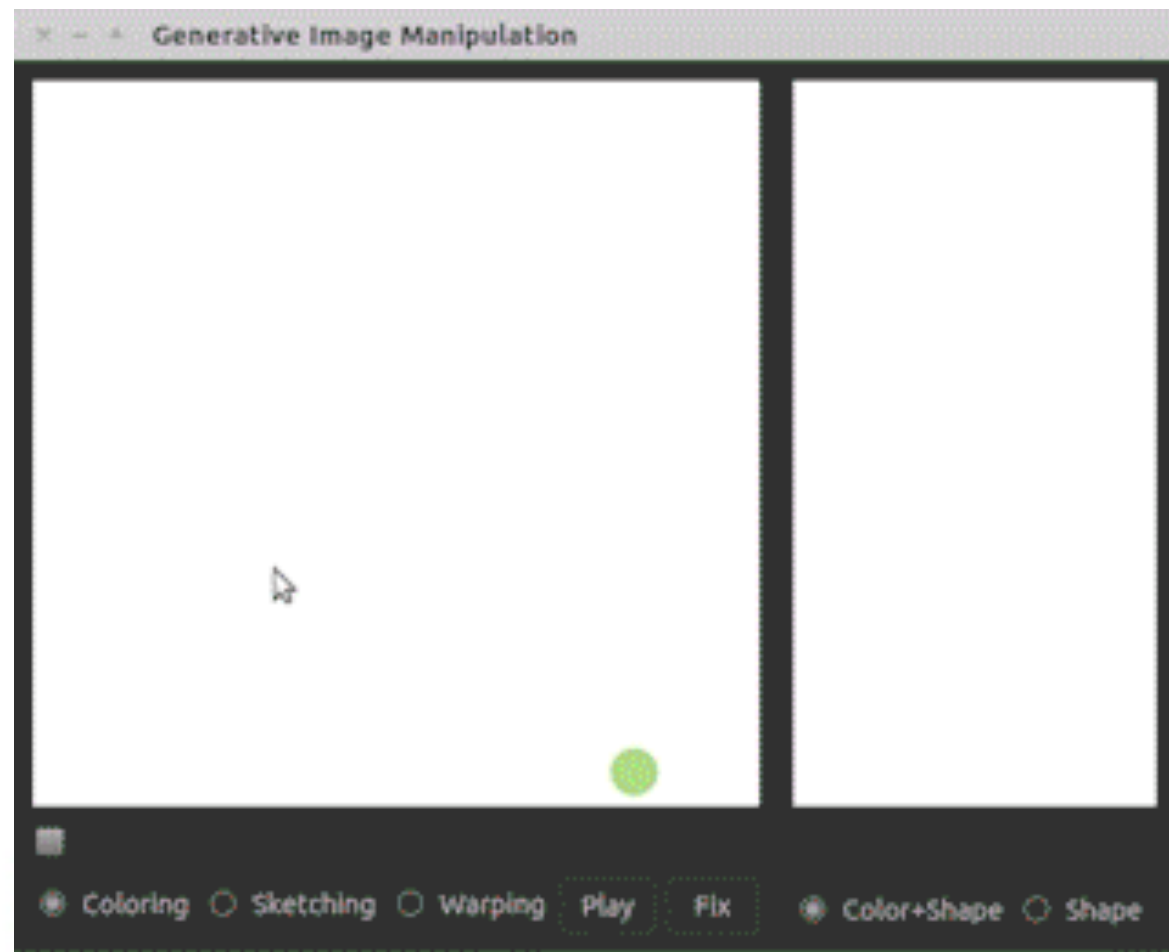
- iGAN/GVM (Zhu et al., 2016)

**Motion** ( $u, v$ ) + **Color** ( $A_{3 \times 4}$ ): estimate per-pixel geometric and color variation

$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x+u, y+v, t+1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dx dy$$

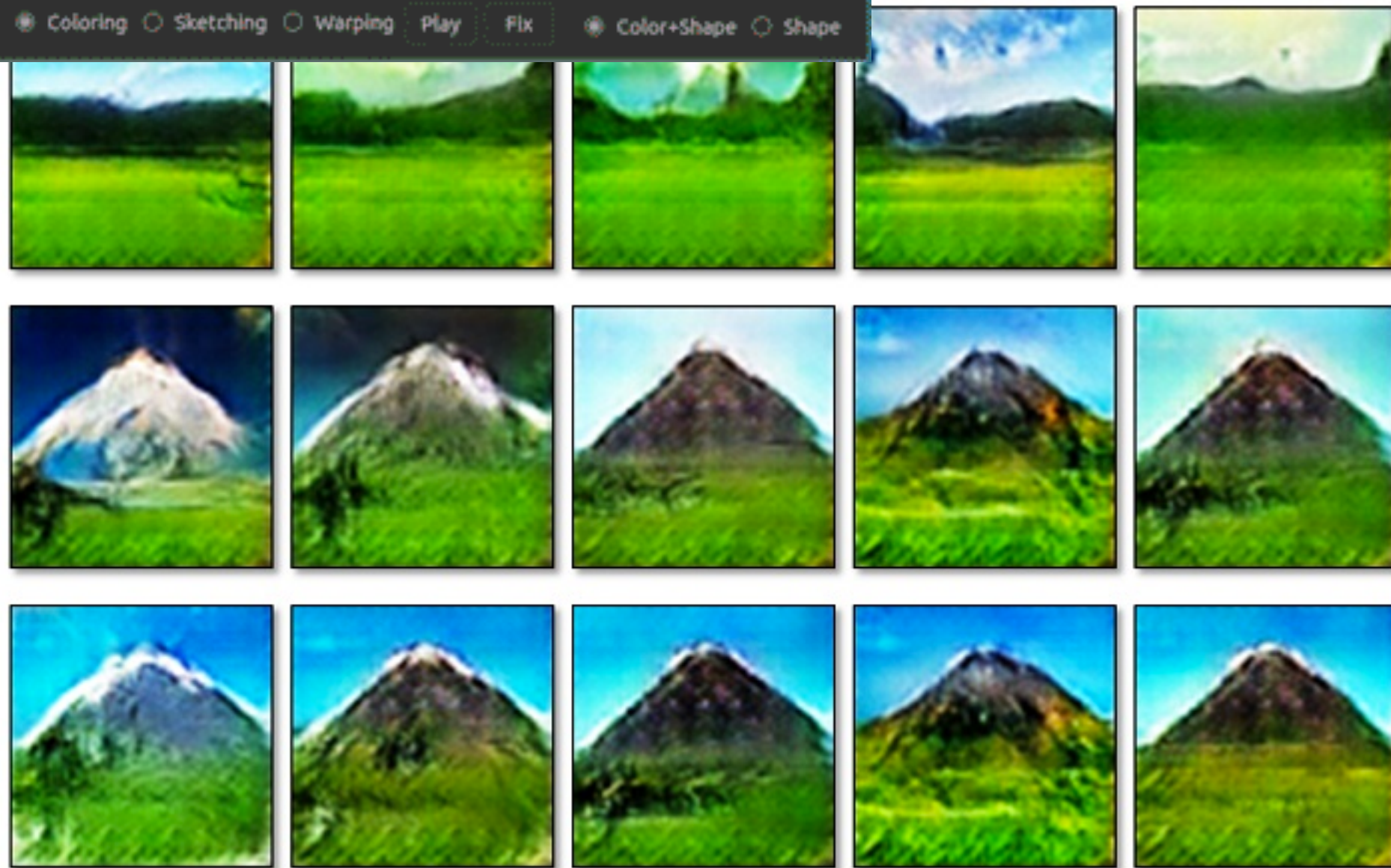
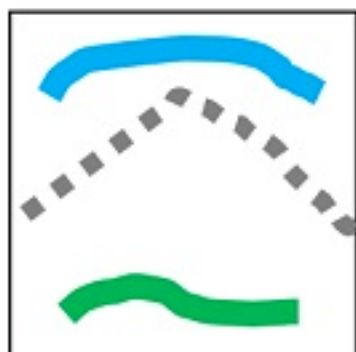
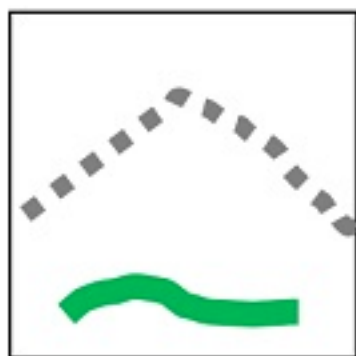
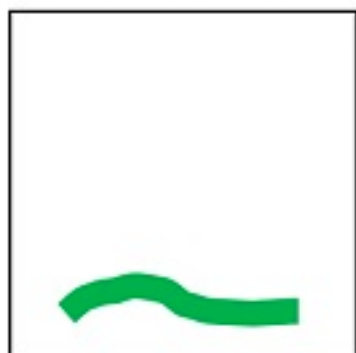






(<https://github.com/junyanz/iGAN>)

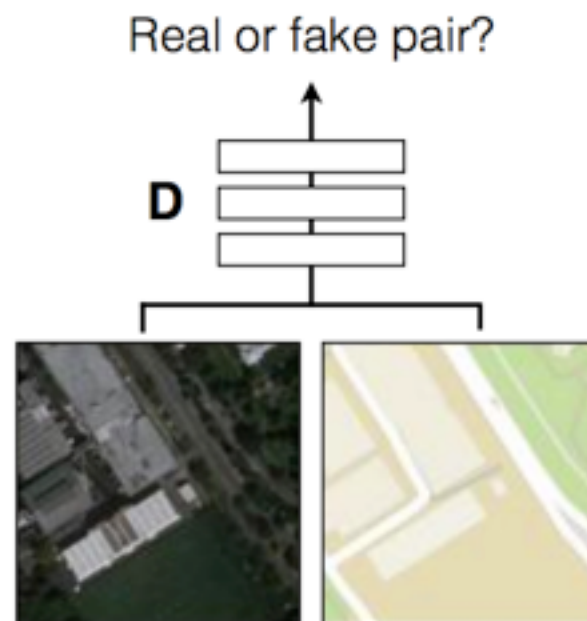
User edits



# Solution 1.1: Partial Guidance

- pix2pix (Isola et al., 2017)

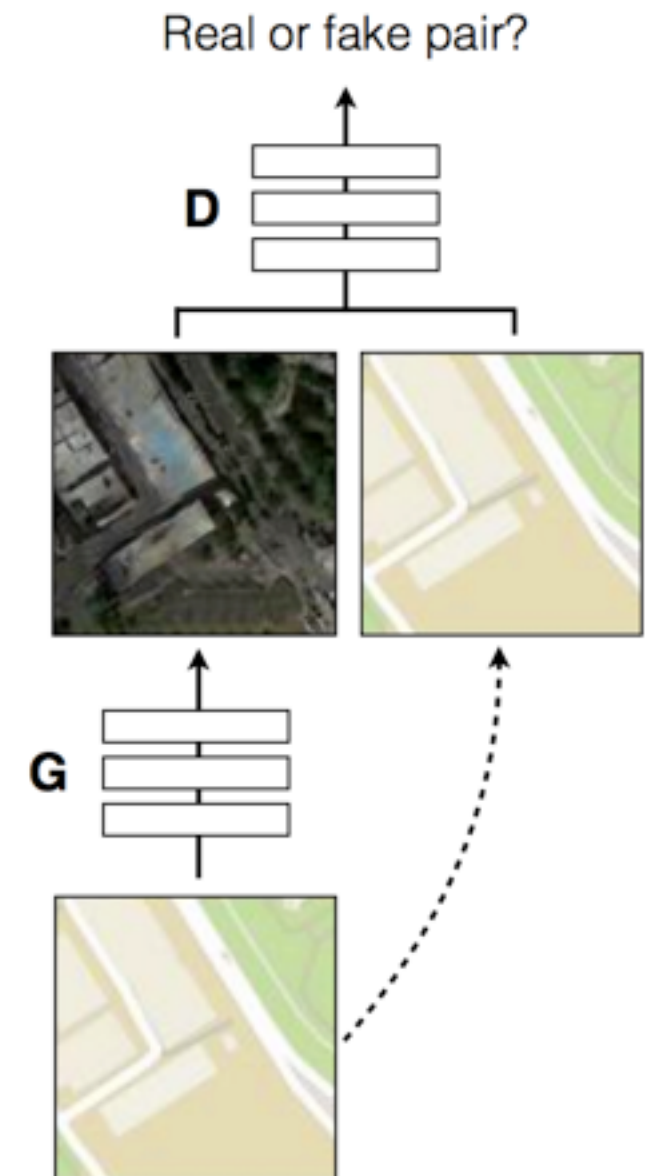
Positive examples



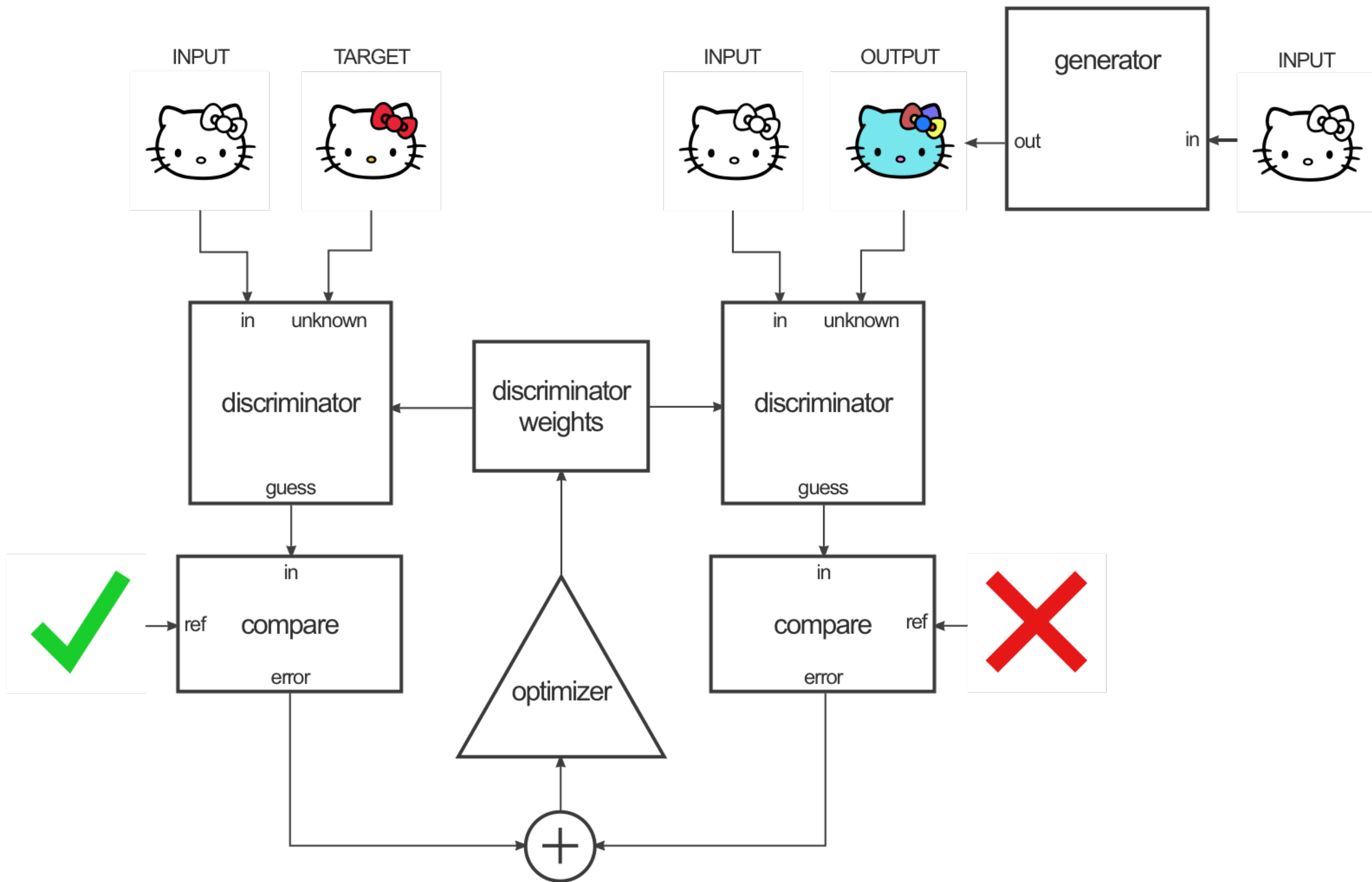
**G** tries to synthesize fake images that fool **D**

**D** tries to identify the fakes

Negative examples



(Isola et al., 2017)







Labels to Street Scene

Labels to Facade

BW to Color

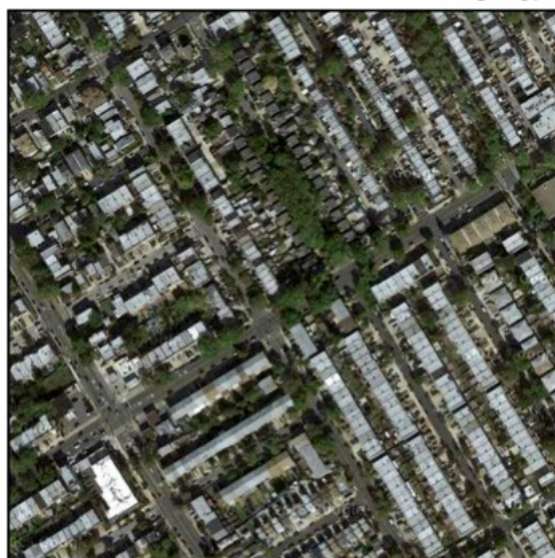


input



output

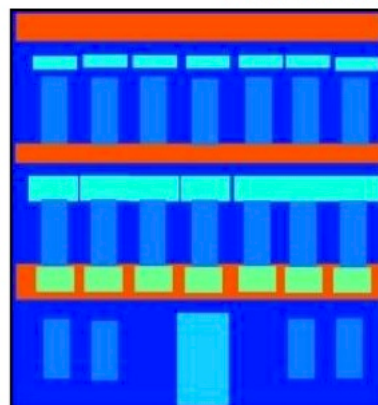
Aerial to Map



input



output

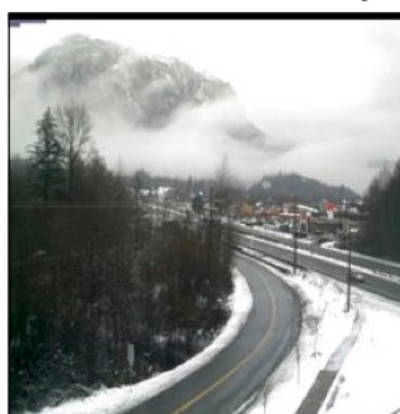


input



output

Day to Night



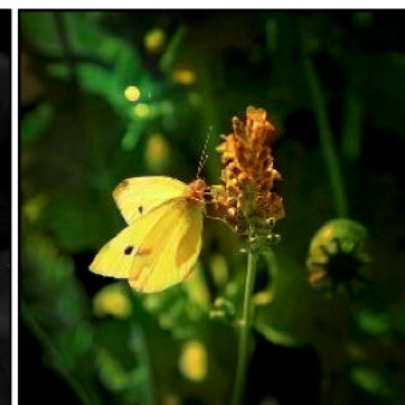
input



output



input



output

Edges to Photo



input



output

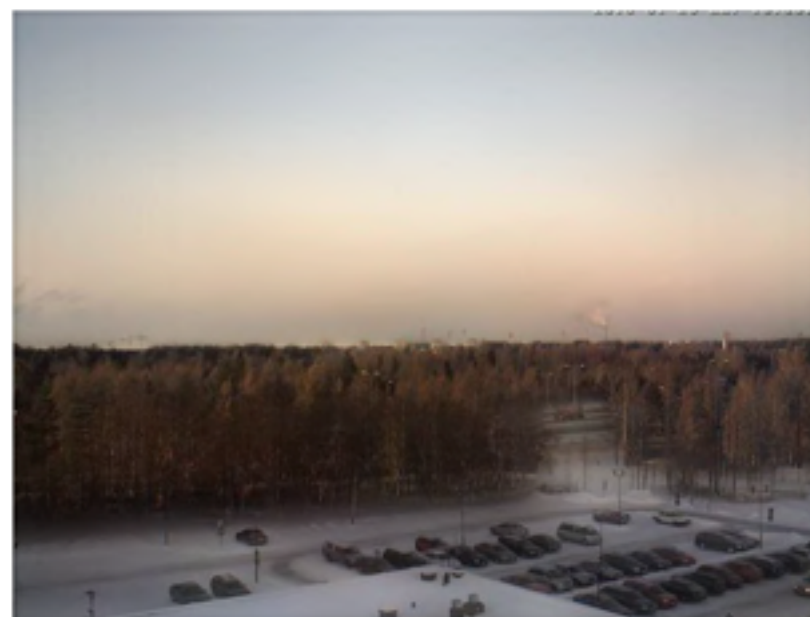


# Solution 1.1: Partial Guidance

- GP-GAN (Wu et al., 2017)



(a)

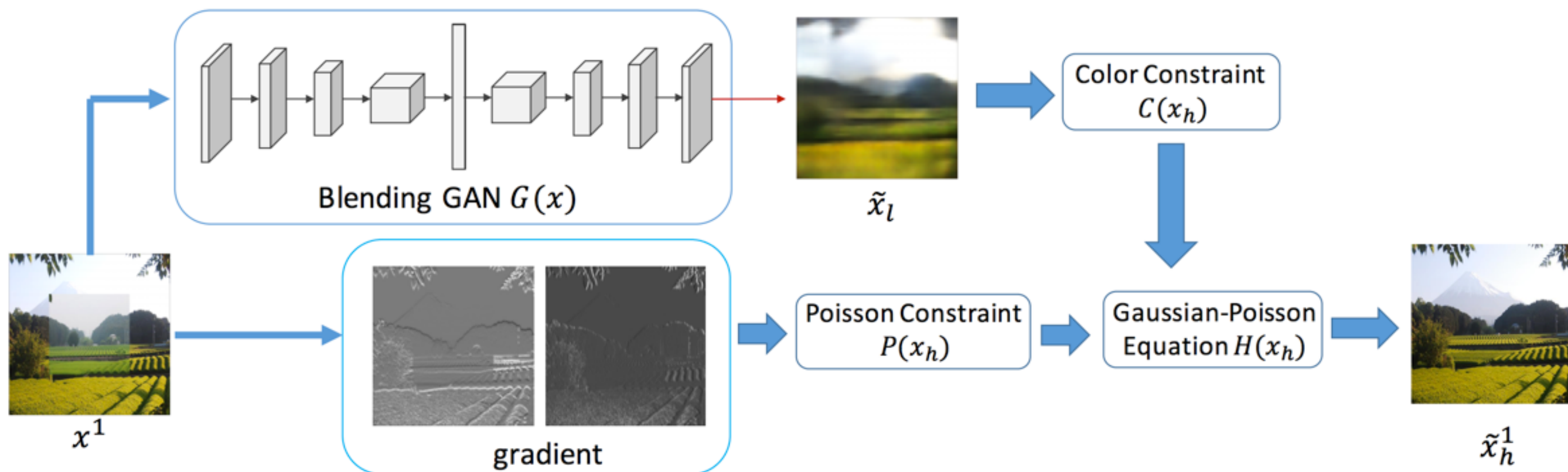


(d)



# Solution 1.1: Partial Guidance

- GP-GAN (Wu et al., 2017)



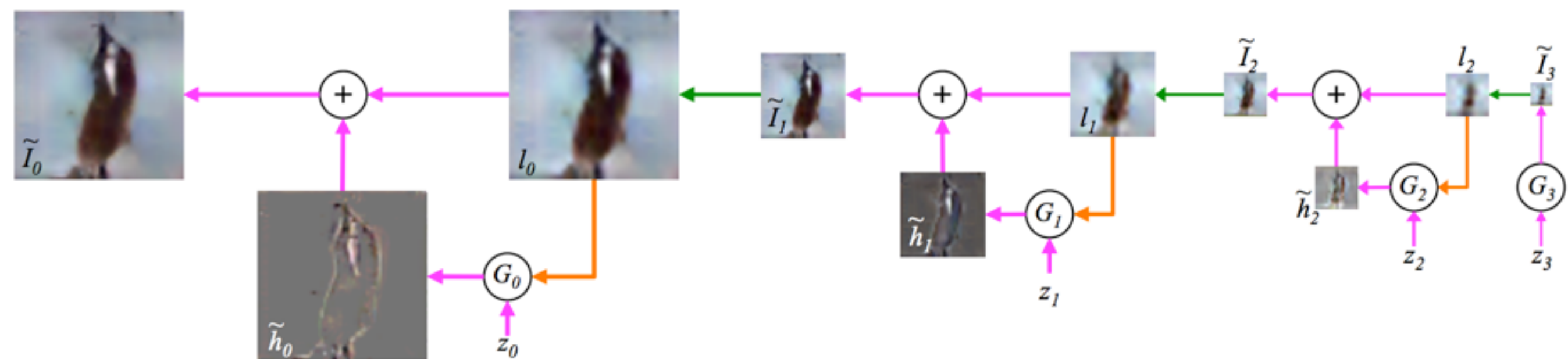
# Solution 1.2: Fine-grained Guidance

- LAPGAN (Denton et al., 2015)
- Matching-aware Discriminator (Reed et al., 2016)
- StackGAN (Zhang et al., 2016)
- PPGN (Nguyen et al., 2017)

# Solution 1.2: Fine-grained Guidance

- LAPGAN (Denton et al., 2015)

$$\min_G \max_D \mathbb{E}_{h, l \sim p_{\text{Data}}(h, l)} [\log D(h, l)] + \mathbb{E}_{z \sim p_{\text{Noise}}(z), l \sim p_l(l)} [\log(1 - D(G(z, l), l))]$$



# Solution 1.2: Fine-grained Guidance

- Matching-aware Discriminator (Reed et al., 2016)
  - implicitly separate two sources of error: unrealistic images (for any text), and realistic images of the wrong class that mismatch the conditioning

$\hat{x} \leftarrow G(z, h)$  {Forward through generator}

$s_r \leftarrow D(x, h)$  {real image, right text}

$s_w \leftarrow D(x, \hat{h})$  {real image, wrong text}

$s_f \leftarrow D(\hat{x}, h)$  {fake image, right text}

$\mathcal{L}_D \leftarrow \log(s_r) + (\log(1 - s_w) + \log(1 - s_f))/2$



## Text descriptions (content)

## Images (style)



The bird has a **yellow breast** with **grey** features and a small beak.

This is a large **white** bird with **black wings** and a **red head**.

A small bird with a **black head and wings** and features grey wings.

This bird has a **white breast**, brown and white coloring on its head and wings, and a thin pointy beak.

A small bird with **white base** and **black stripes** throughout its belly, head, and feathers.

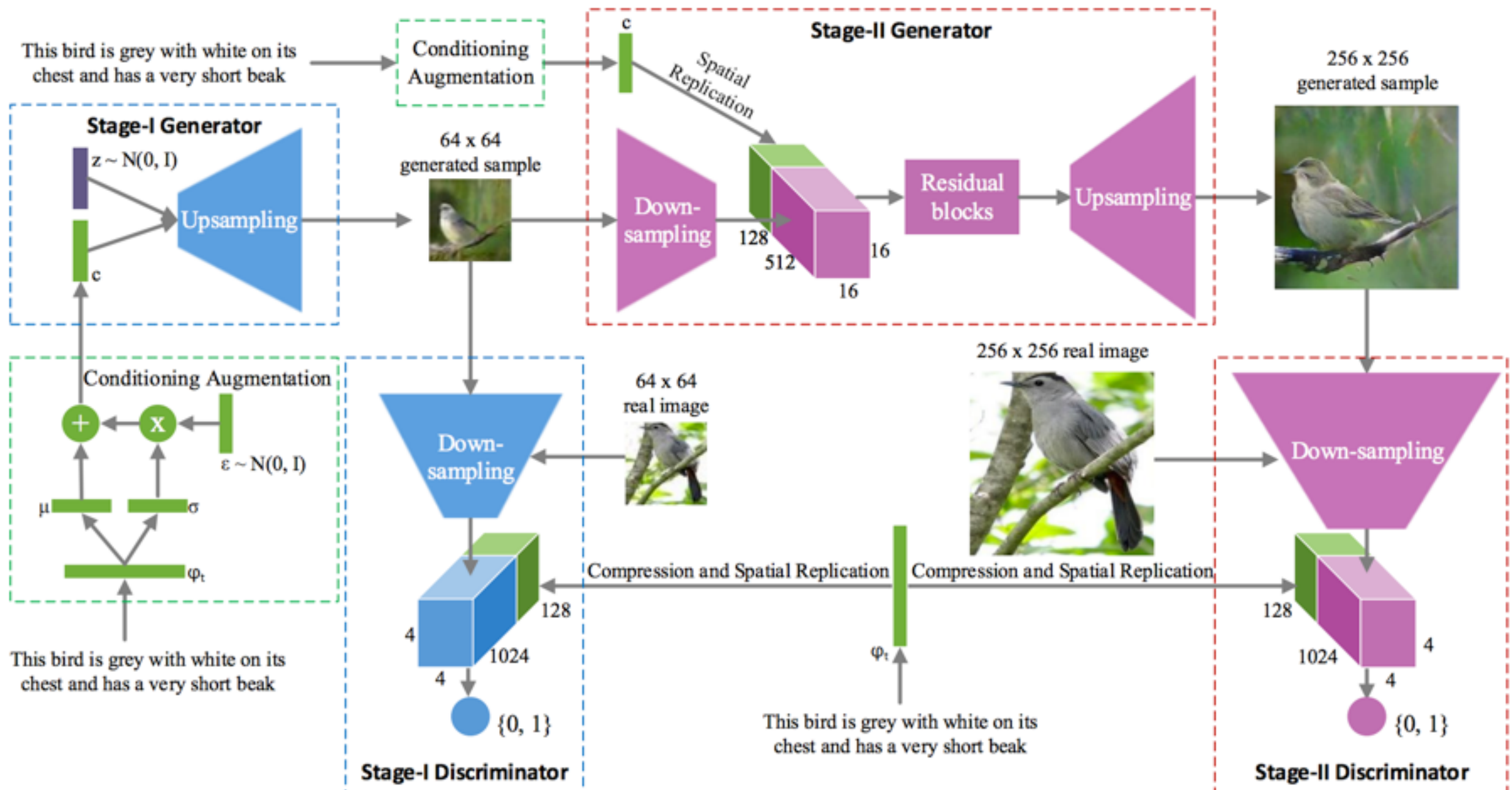
A small sized bird that has a cream belly and a short pointed bill.










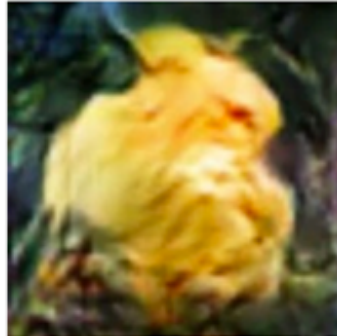



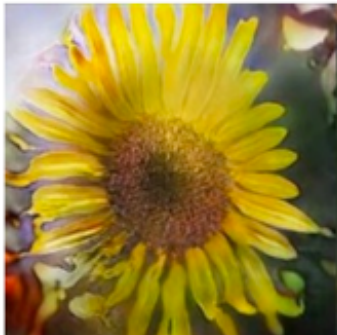


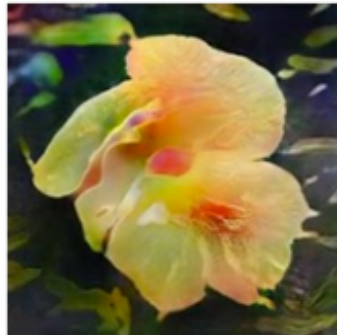

# Solution 1.2: Fine-grained Guidance

- StackGAN (Zhang et al., 2016)



# Solution 1.2: Fine-grained Guidance

- StackGAN (Zhang et al., 2016)

Text description	This flower has petals that are white and has pink shading	This flower has a lot of small purple petals in a dome-like configuration	This flower has long thin yellow petals and a lot of yellow anthers in the center	This flower is pink, white, and yellow in color, and has petals that are striped	This flower is white and yellow in color, with petals that are wavy and smooth	This flower has upturned petals which are thin and orange with rounded edges	This flower has petals that are dark pink with white edges and pink stamen
64x64 GAN-INT-CLS [22]							
256x256 StackGAN							

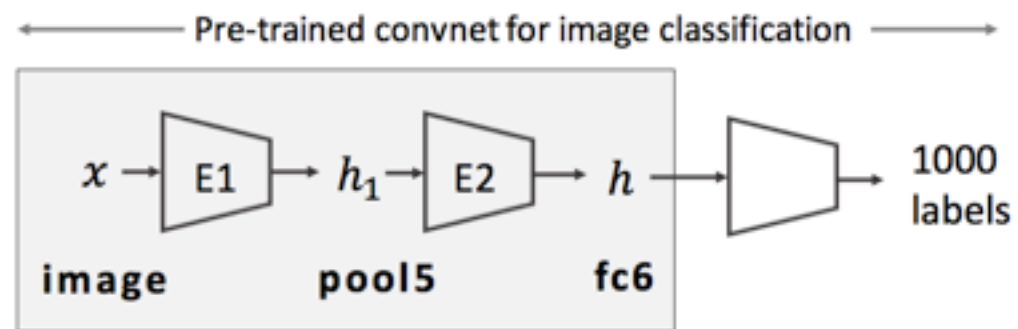


# Solution 1.2: Fine-grained Guidance

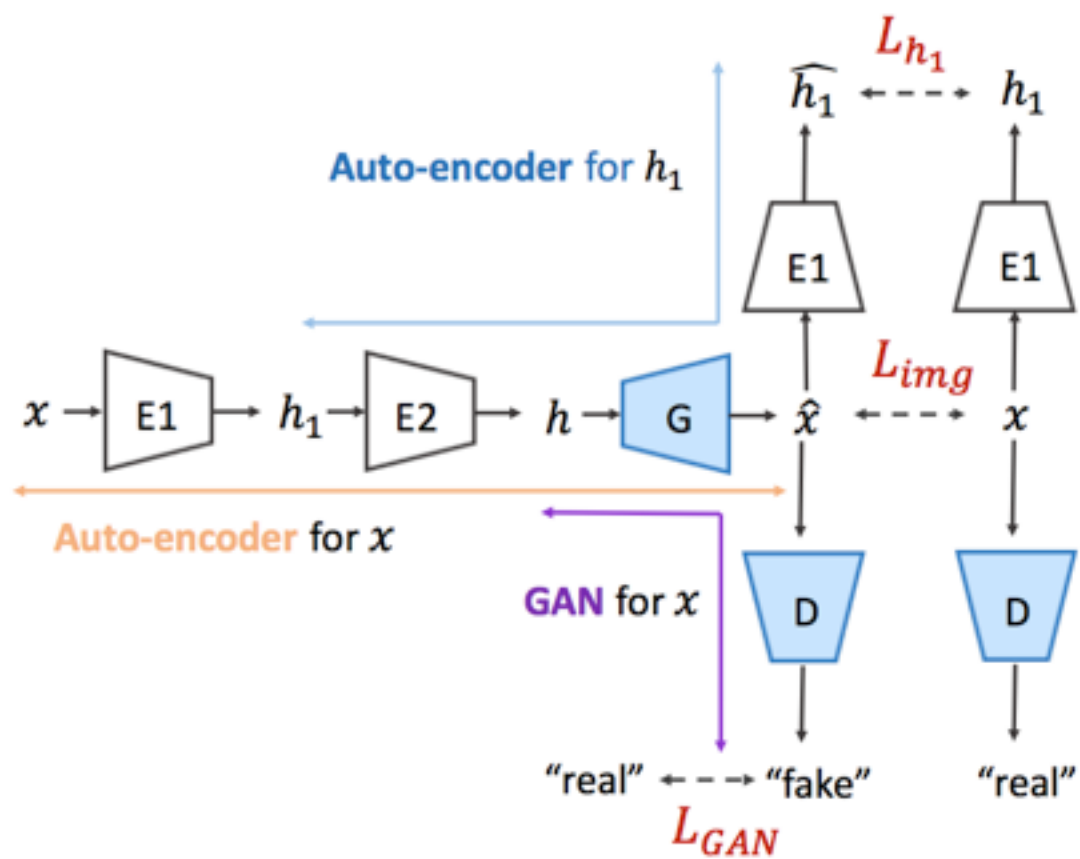
- PPGN (Nguyen et al., 2017)



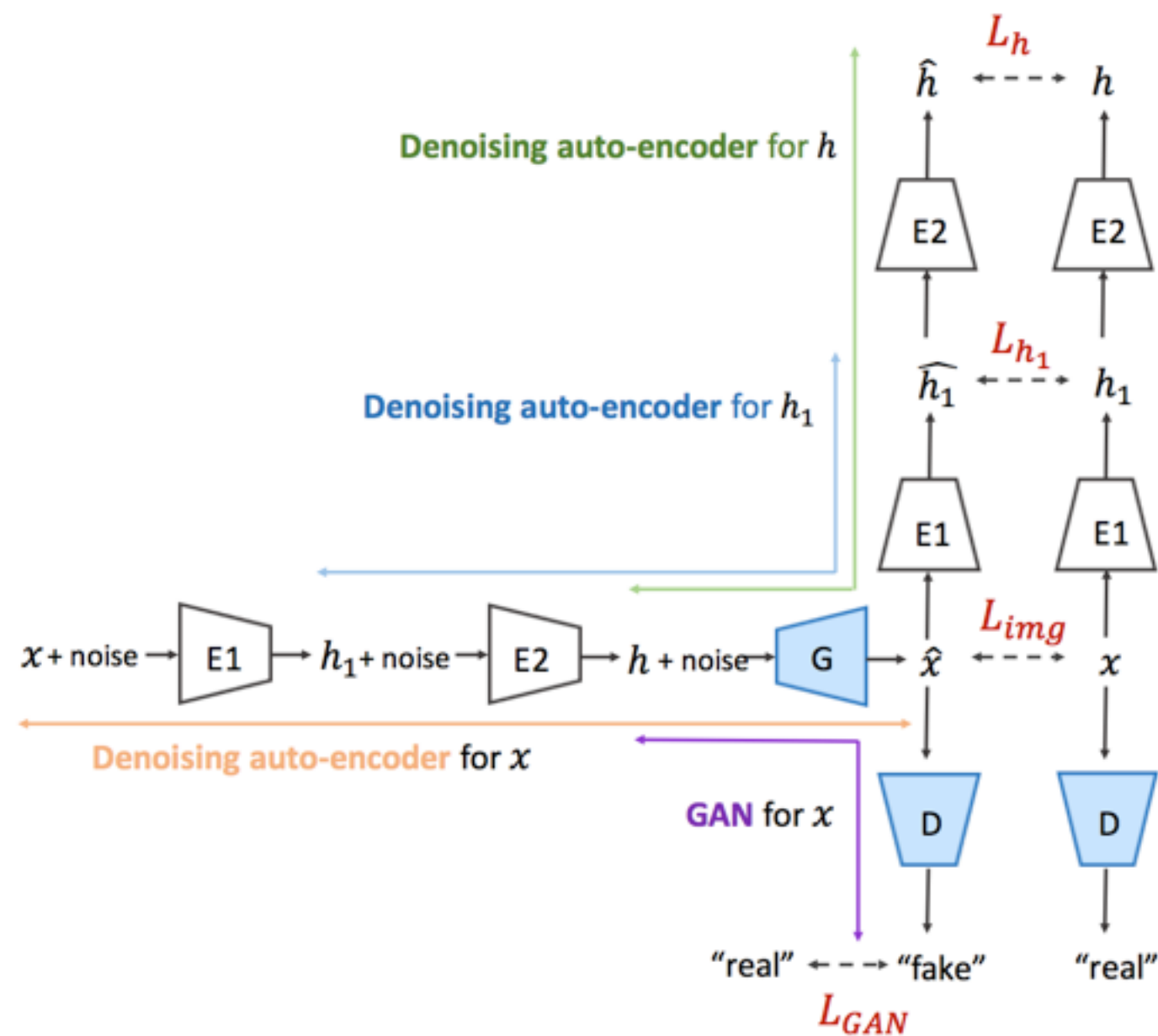




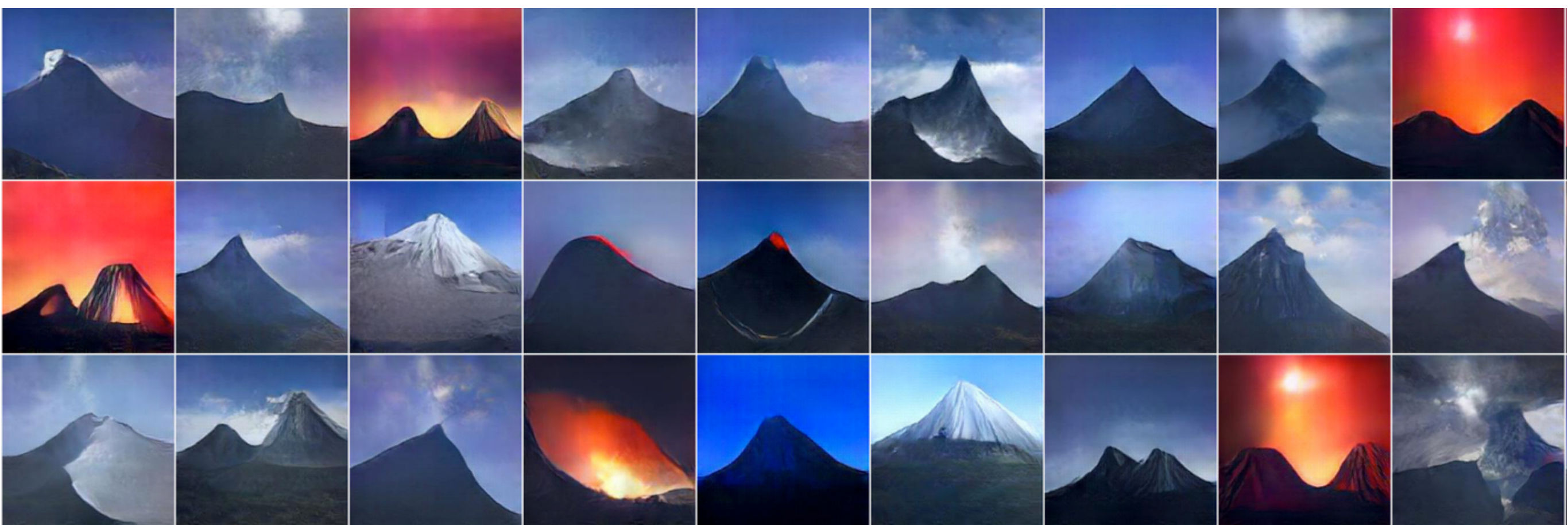
(a) Encoder network E



(b) Noiseless joint PPGN-h



(c) Joint PPGN-h



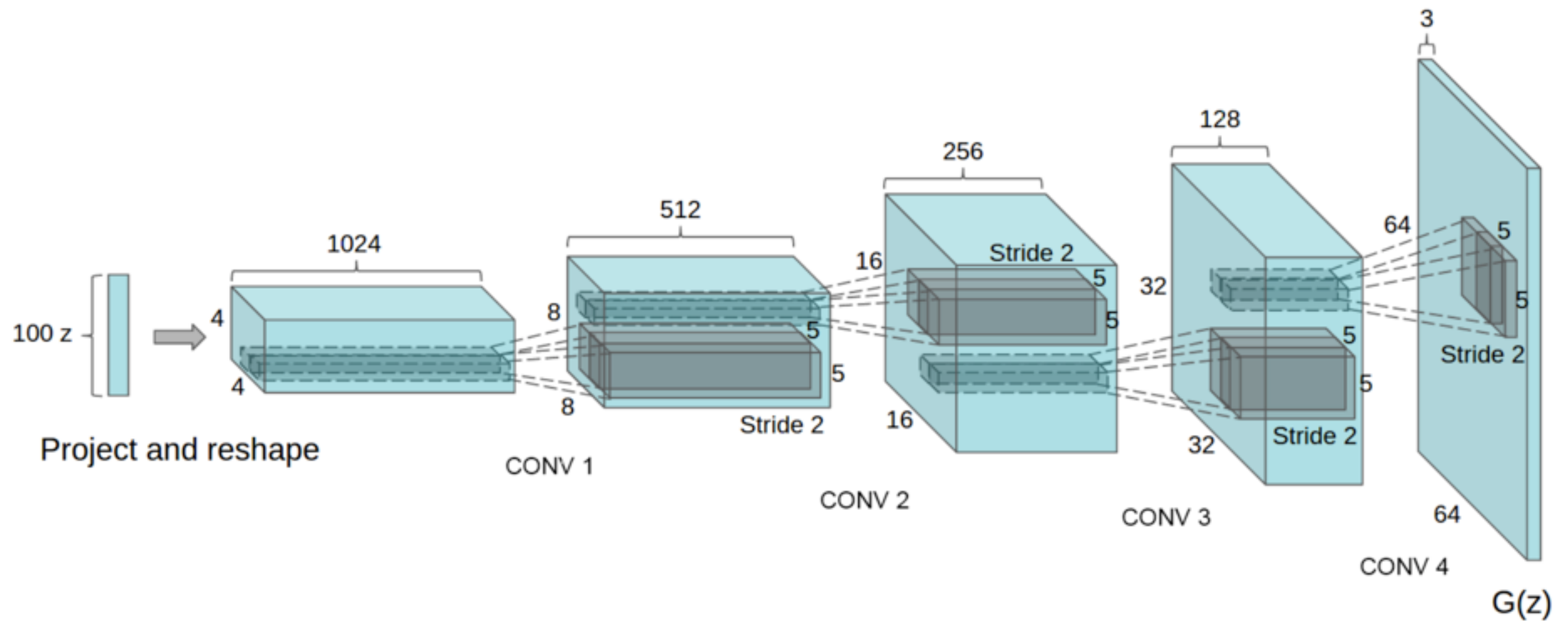


# Solution 1.3: Special Architecture

- DCGAN (Radford et al., 2016)
- pix2pix (Isola et al., 2017)
- GP-GAN (Wu et al., 2017)

# Solution 1.3: Special Architecture

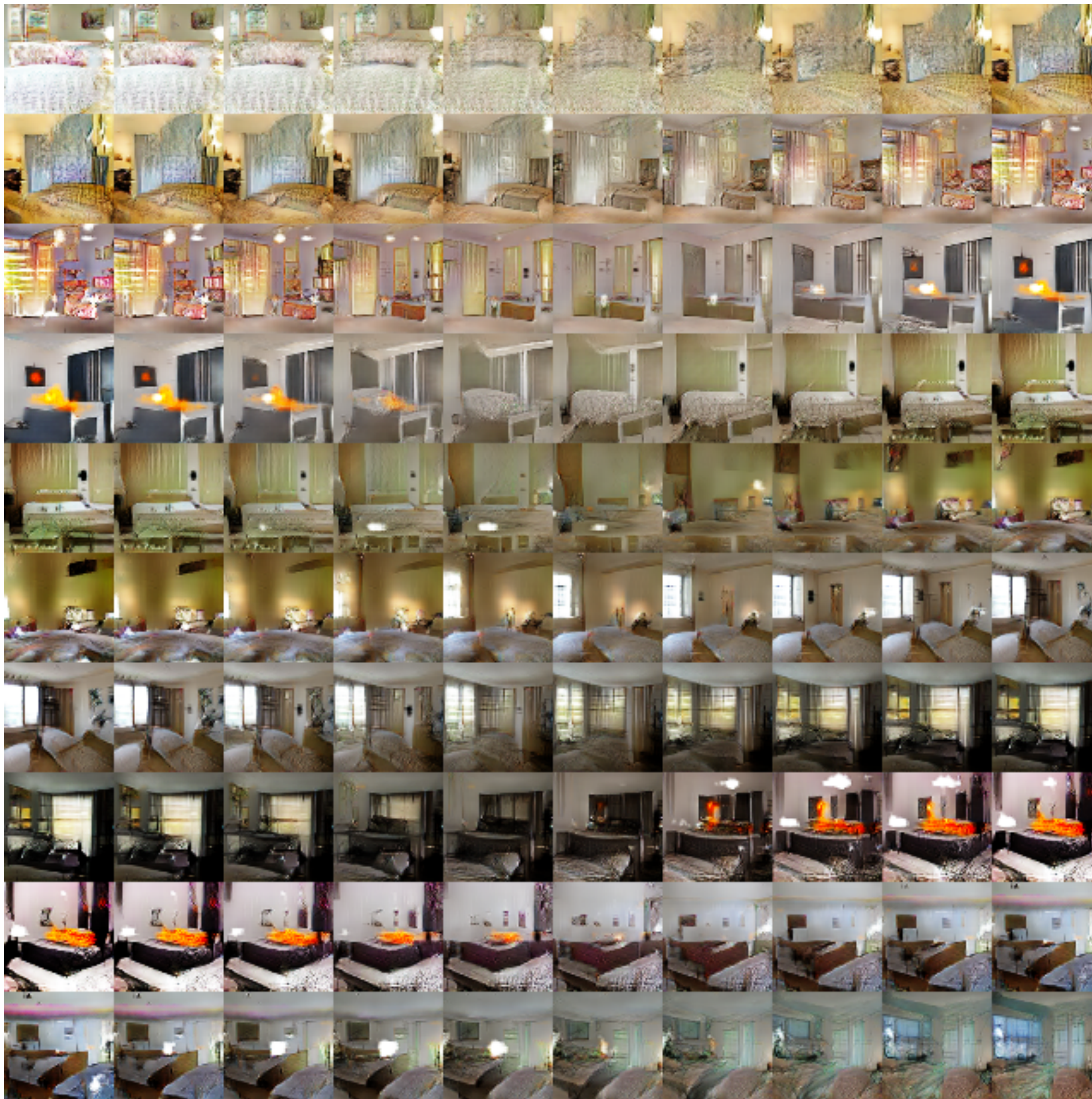
- DCGAN (Radford et al., 2016)



# Solution 1.3: Special Architecture

- DCGAN (Radford et al., 2016)
  - Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator)
  - Use batchnorm in both the generator and the discriminator
  - Remove fully connected hidden layers
  - Use ReLU activation in generator for all layers except for the output, which uses Tanh; Use LeakyReLU activation in the discriminator for all layers



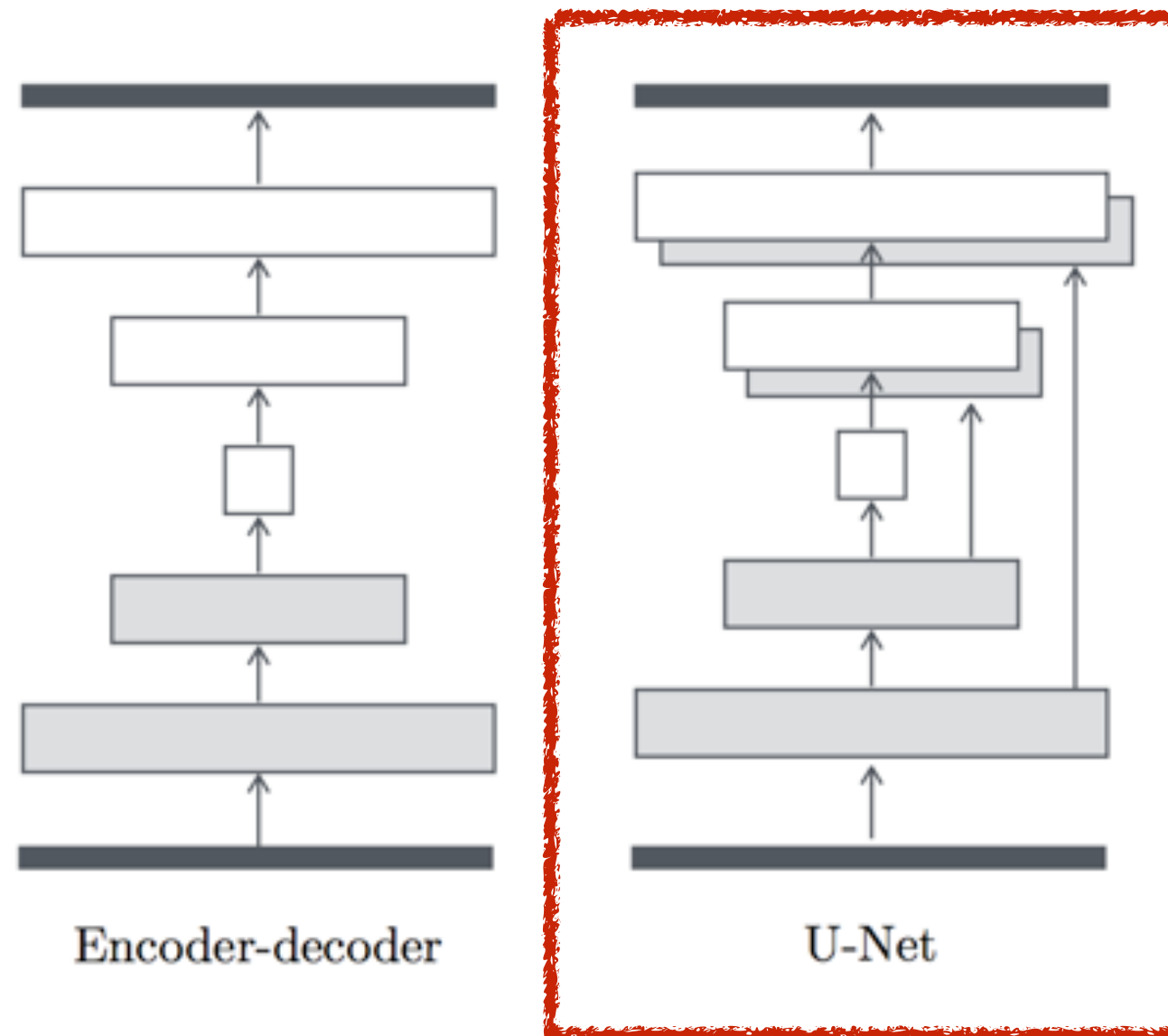


(Radford et al., 2016)



# Solution 1.3: Special Architecture

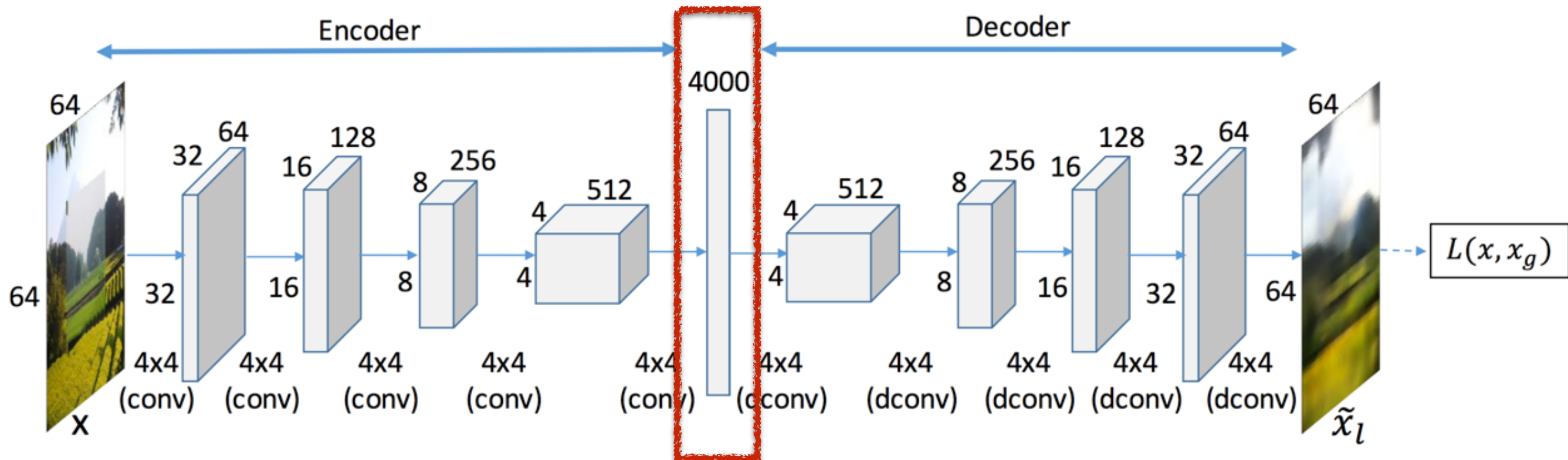
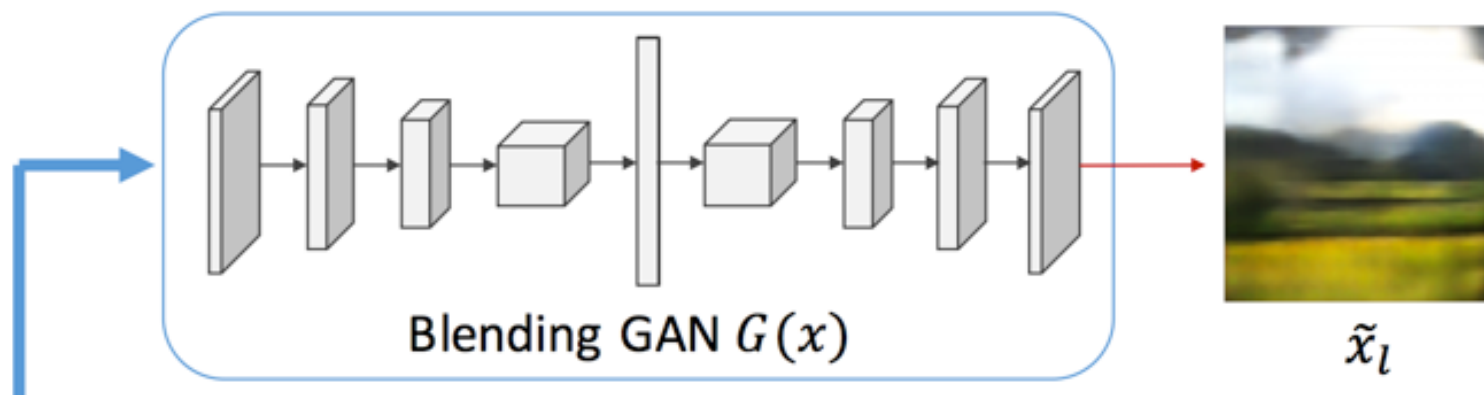
- pix2pix (Isola et al., 2017)





# Solution 1.3: Special Architecture

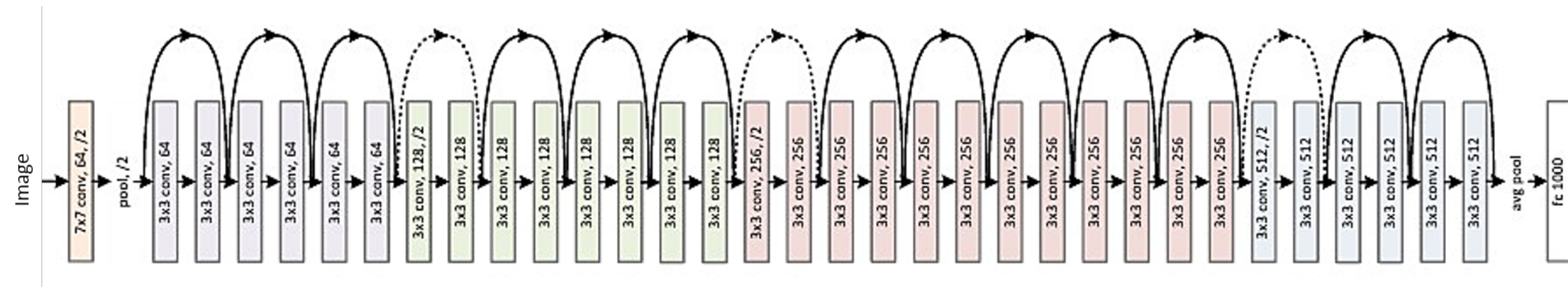
- GP-GAN (Wu et al., 2017)



# Difficulty 1

**Tackled!**

- The gradient issues existed in deep neural networks
- The deeper, the more difficult



# Content

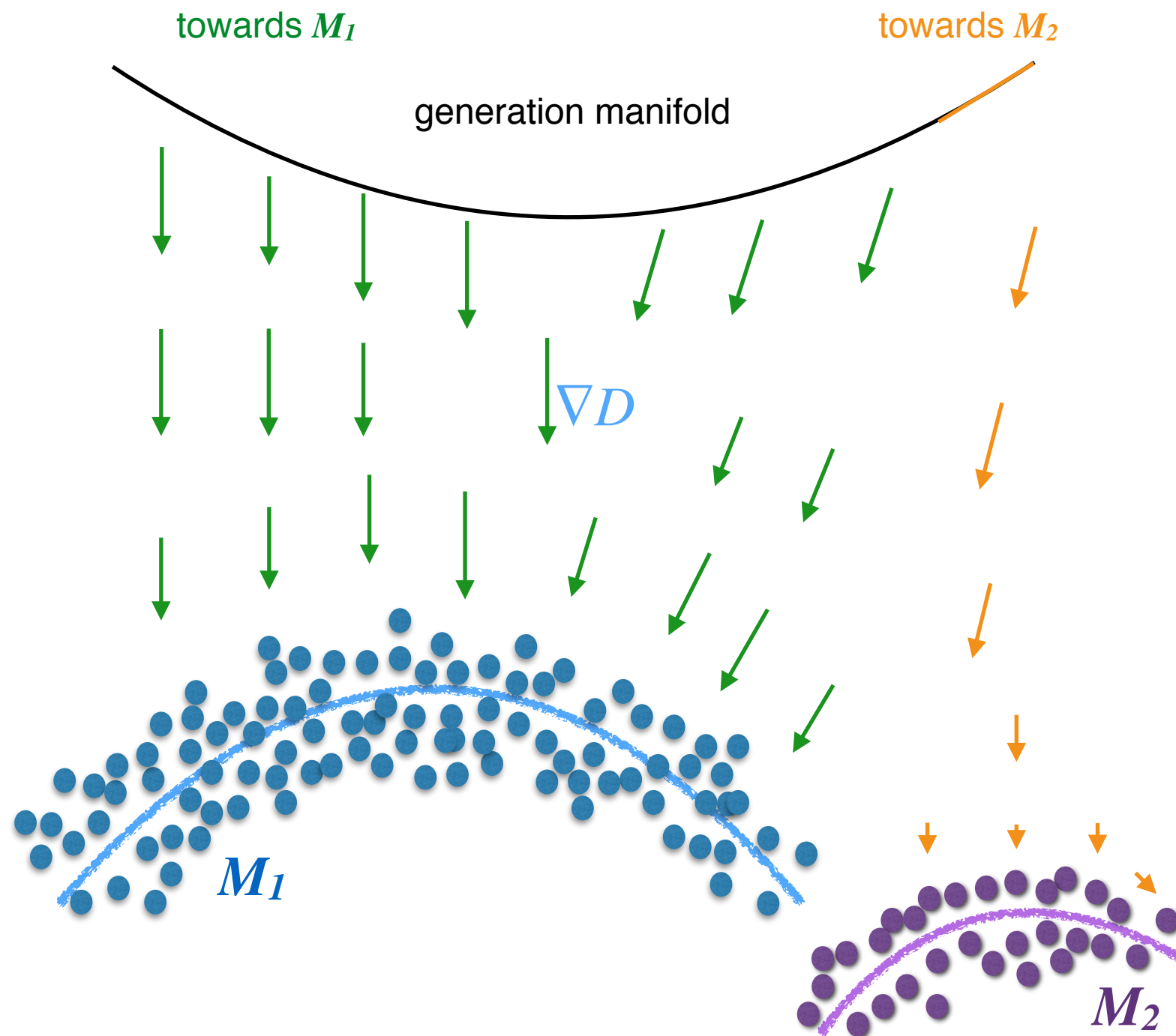
- Generative Adversarial Networks
  - Basics and Attractiveness
  - Difficulties
- Solution 1: Partial and Fine-grained Guidance
- **Solution 2: Encoder-incorporated**
- Solution 3: Wasserstein Distance

# Solution 2: Encoder-incorporated

- Mode Regularized GANs (Che et al., 2017)
- Tackling the gradient vanishing issue and mode missing problem by incorporating an additional encoder  $E$  to:
  - (1) “enforce”  $P_r$  and  $P_g$  overlap
  - (2) “build a bridge” between *fake data* and *real data*



# Mode Missing Problem



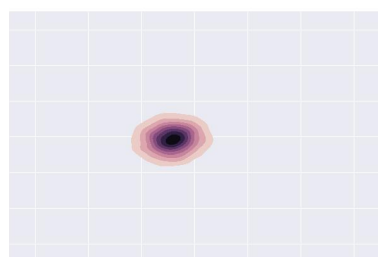
# Mode Missing Problem

$$\min_G \max_D V(G, D) \neq \max_D \min_G V(G, D)$$

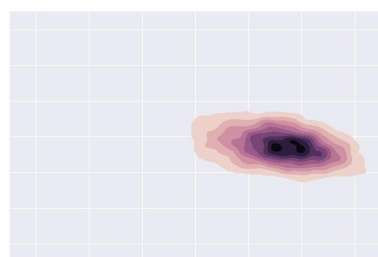
- ***D*** in inner loop: convergence to correct distribution
- ***G*** in inner loop: place all mass on most likely point



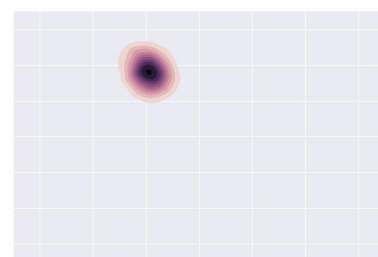
Target



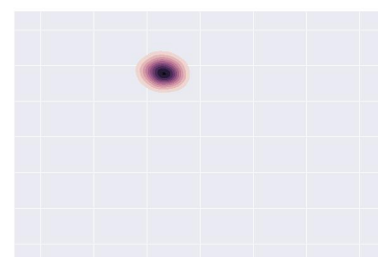
Step 0



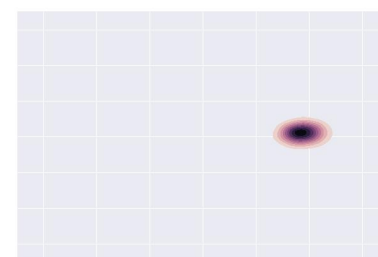
Step 5k



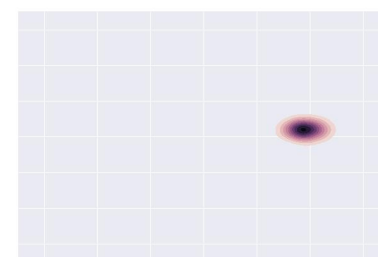
Step 10k



Step 15k



Step 20k



Step 25k

# Mode Regularized GANs

- Regularized GANs

- for encoder ***E***:  $\mathbb{E}_{x \sim p_d} [\lambda_1 d(x, G \circ E(x)) + \lambda_2 \log D(G \circ E(x))]$

- for generator ***G***:

$$-\mathbb{E}_z [\log D(G(z))] + \mathbb{E}_{x \sim p_d} [\lambda_1 d(x, G \circ E(x)) + \lambda_2 \log D(G \circ E(x))]$$

- for discriminator ***D***: same as vanilla GAN

# Mode Regularized GANs

- Regularized GANs

- for encoder ***E***:  $\mathbb{E}_{x \sim p_d} [\lambda_1 d(x, G \circ E(x)) + \lambda_2 \log D(G \circ E(x))]$

- for generator ***G***:

$$-\mathbb{E}_z [\log D(G(z))] + \mathbb{E}_{x \sim p_d} [\lambda_1 d(x, G \circ E(x)) + \lambda_2 \log D(G \circ E(x))]$$

- for discriminator ***D***: same as vanilla GAN

- But it still suffers from gradient vanishing!

- because ***D*** is still comparing between real data and fake data



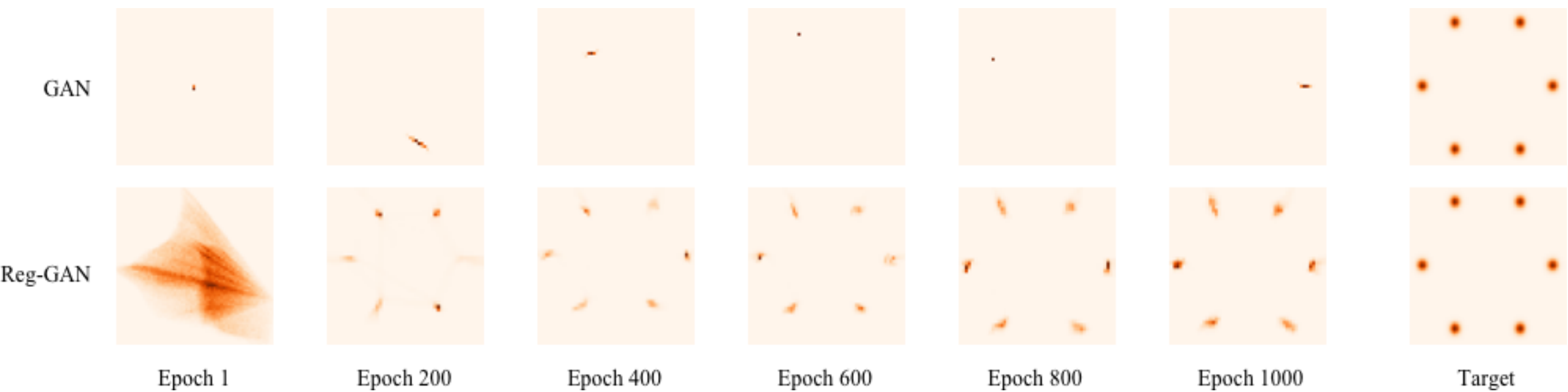
# Mode Regularized GANs

- Manifold-Diffusion GANs (MDGAN):
  - Manifold-step:
    - Try to match the generation manifold and the real data manifold
  - Diffusion-step:
    - Try to distribute the probability mass on the generation manifold fairly according to the real data distribution

# Mode Regularized GANs

- Manifold-Diffusion GANs (MDGAN):
  - Manifold-step:
    - Try to match the generation manifold and the real data manifold
  - Diffusion-step:
    - Try to distribute the probability mass on the generation manifold fairly according to the real data distribution
- ***D*** is firstly comparing between real data and the encoded data — much harder!

# Mode Regularized GANs



# Mode Regularized GANs

MDGAN



Regularized  
-GAN



ALI



VAEGAN



DCGAN



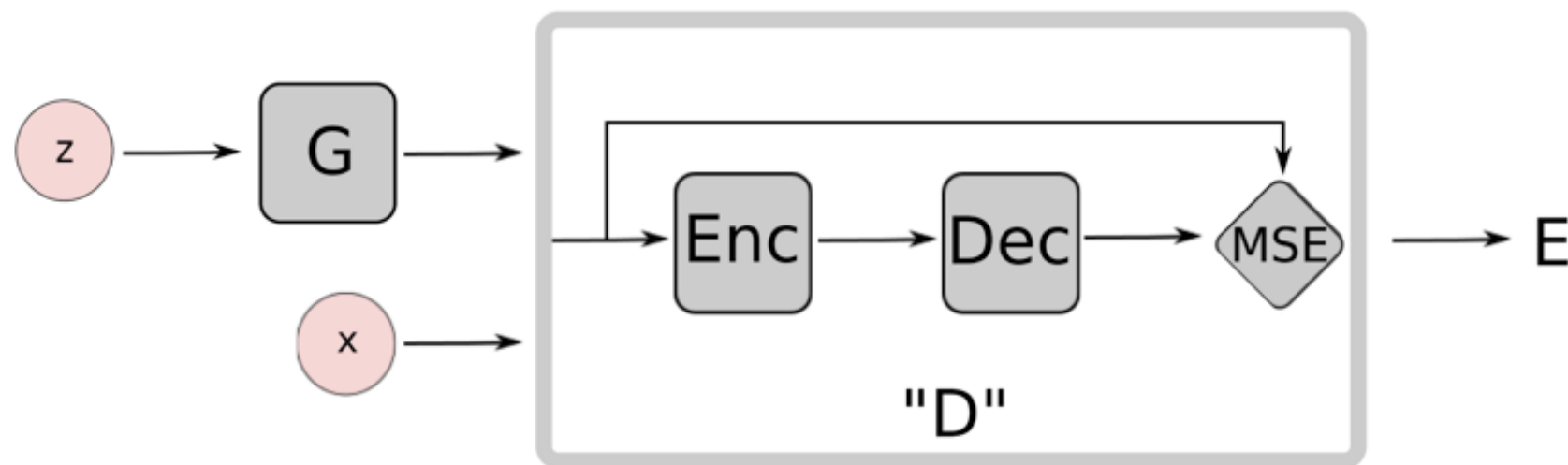


# Solution 2: Encoder-incorporated

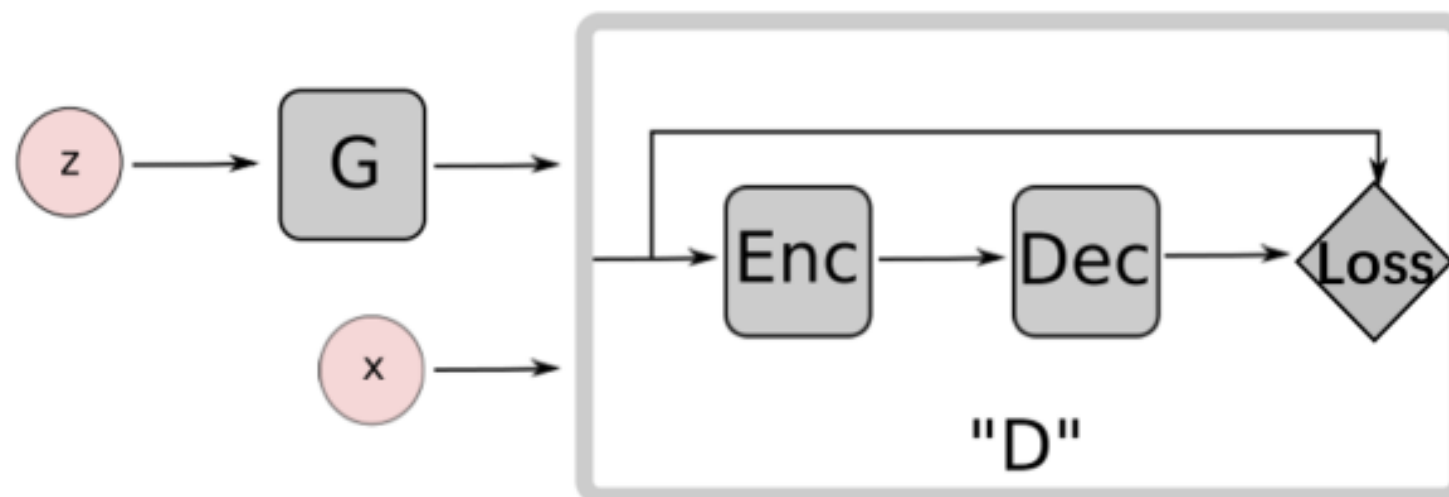
- Mode Regularized GANs (Che et al., 2017)
- Energy-based GANs (Zhao et al., 2017)
- Boundary Equilibrium GANs (Berthelot et al., 2017)
- etc.

# Solution 2: Encoder-incorporated

- Energy-based GANs (Zhao et al., 2017)



- Boundary Equilibrium GANs (Berthelot et al., 2017)

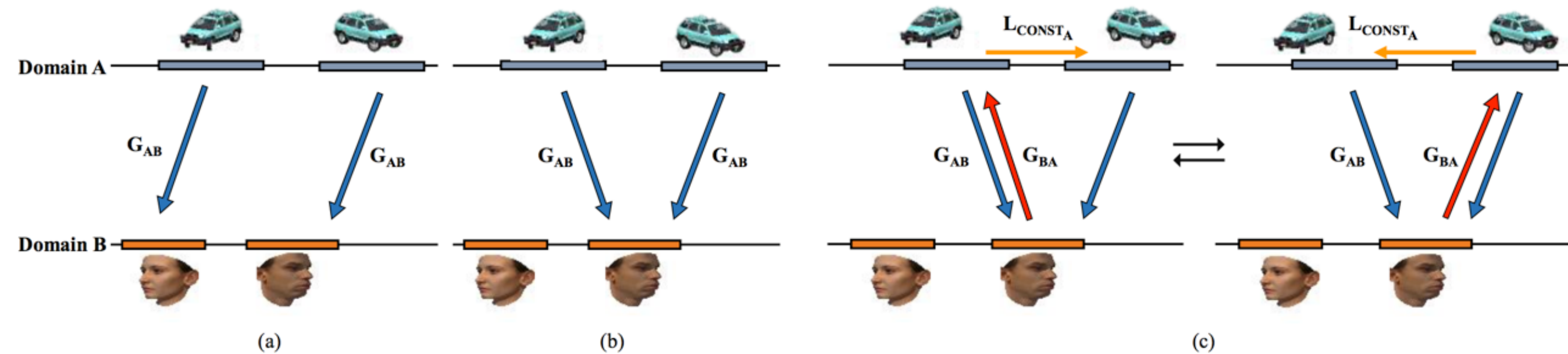


(Zhao et al., 2017)  
(Berthelot et al., 2017)

## Solution 2: *\*Noisy Input*

- Add noise to input (both real data and fake data) before passing into D (Arjovsky & Bottou, 2017, *Theorem 3.2*)
- Add noise to layers in D and G (Zhao et al., 2017)
- Instance Noise (Sønderby et al., 2017)
- All these are indeed “enforcing”  $P_r$  and  $P_g$  to overlap

# Review Mode Missing Problem

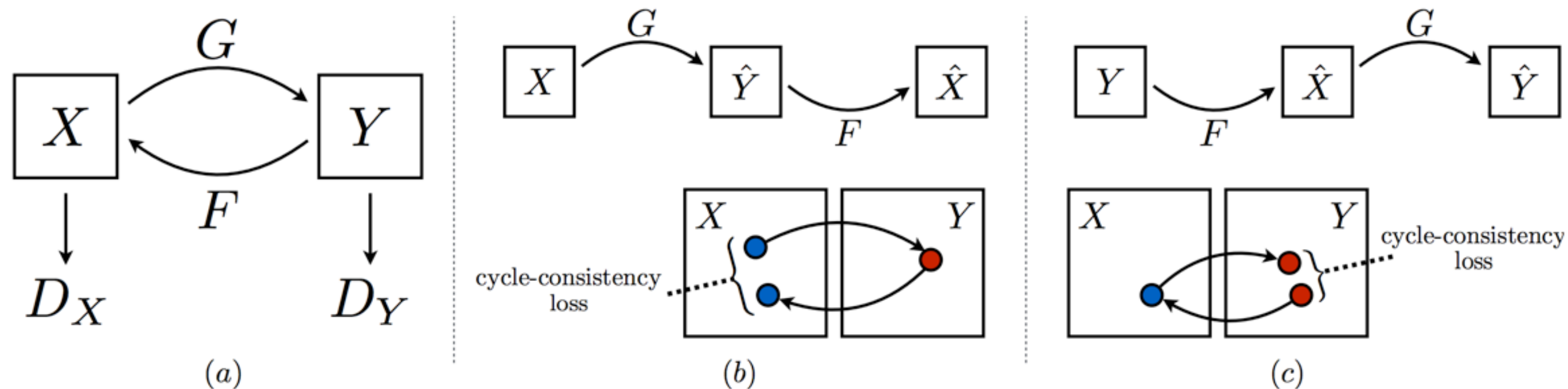


*Figure 3.* Illustration of our models on simplified one dimensional domains. (a) ideal mapping from domain A to domain B in which the two domain A modes map to two different domain B modes, (b) GAN model failure case, (c) GAN with reconstruction model failure case.



# Solution 2: Encoders-incorporated

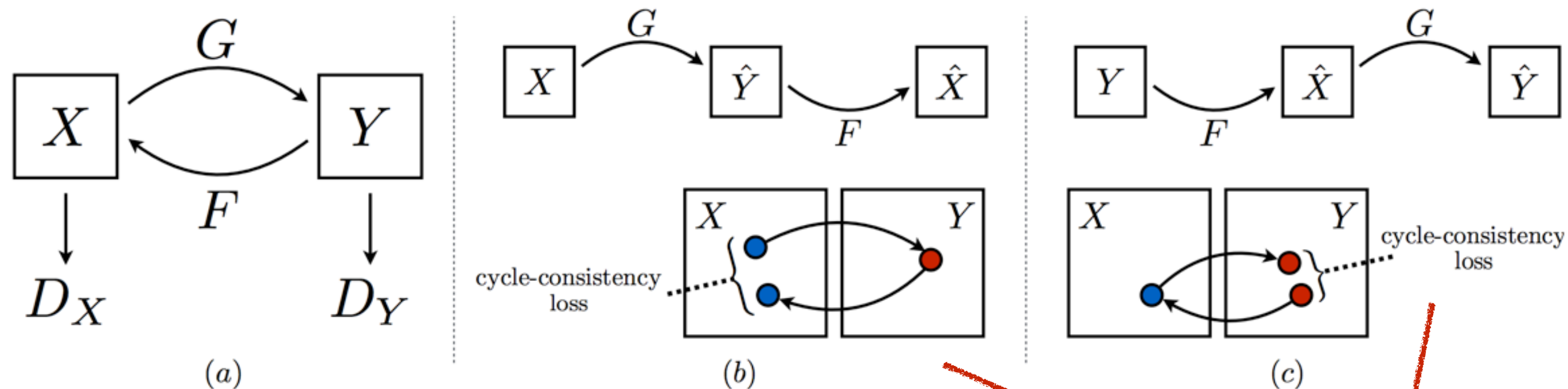
- CycleGAN (Zhu et al., 2017)



- DiscoGAN (Kim et al., 2017)
- DualGAN (Yi et al., 2017)

# Solution 2: Encoders-incorporated

- CycleGAN (Zhu et al., 2017)



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

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]. \end{aligned}$$





(<https://junyanz.github.io/CycleGAN/>)

Monet ↔ Photos



Monet → photo



photo → Monet

Zebras ↔ Horses



zebra → horse



horse → zebra

Summer ↔ Winter



summer → winter



winter → summer



Photograph



Monet



Van Gogh



Cezanne



Ukiyo-e



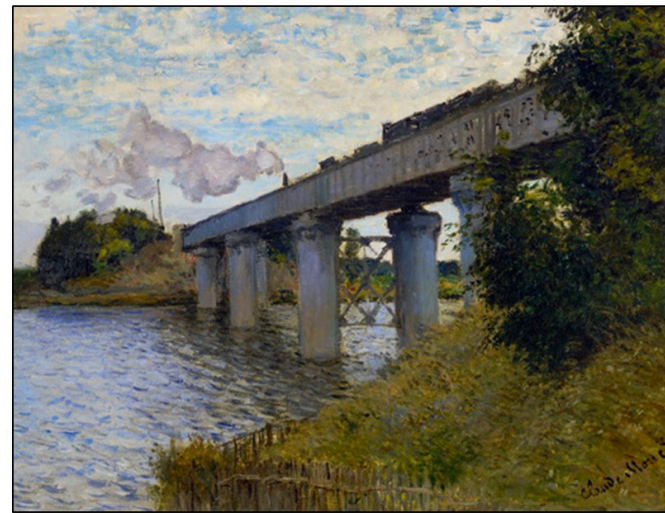
Input



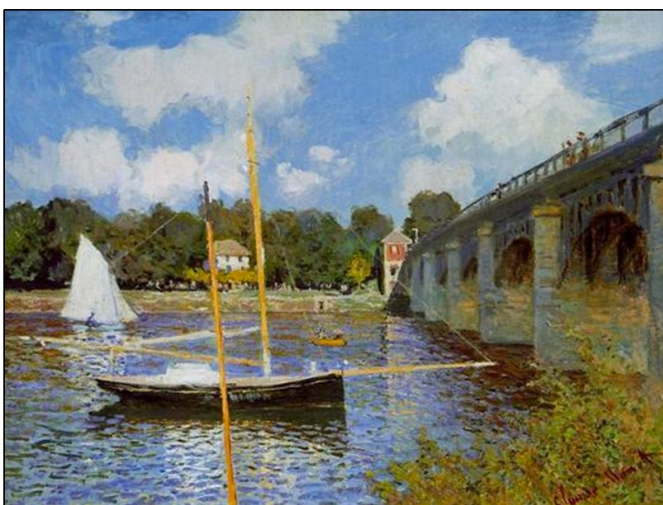
Output



Input



Output





Input



Monet



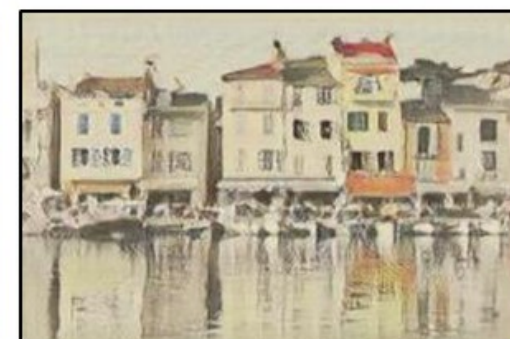
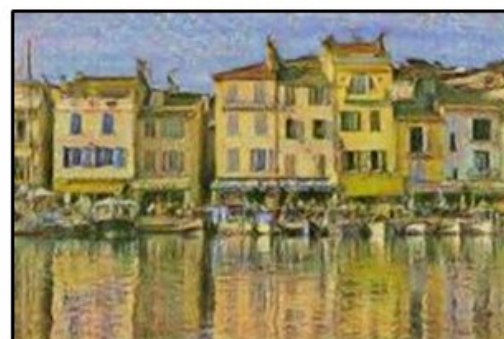
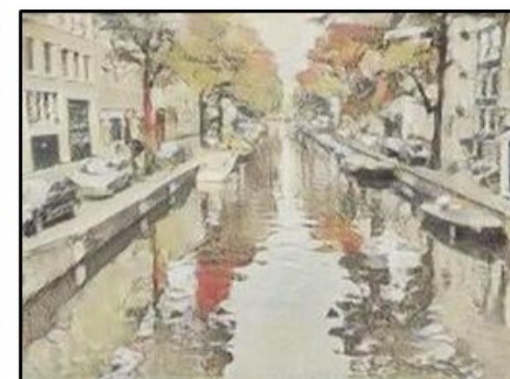
Van Gogh



Cezanne



Ukiyo-e





Input



Output



Input

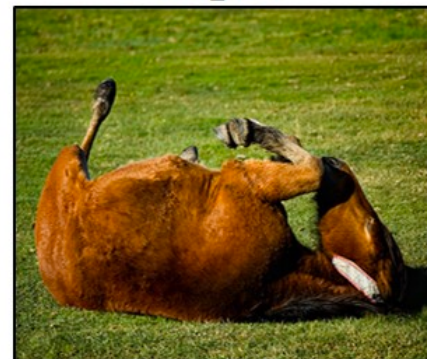


Output

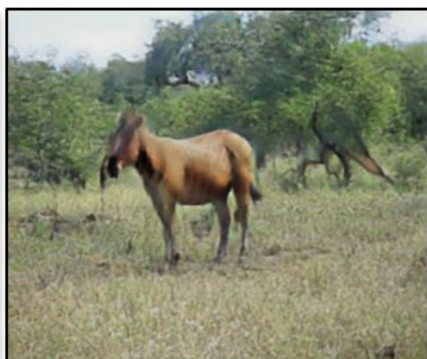


horse → zebra

Input



Output



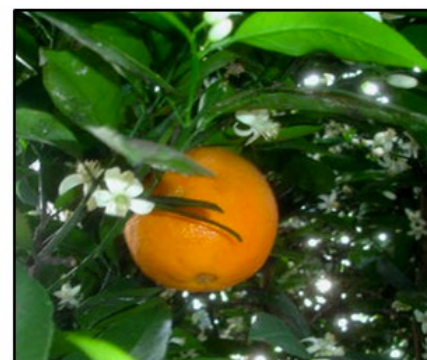
zebra → horse



apple → orange



orange → apple





# Solution 2: Encoders-incorporated

- CycleGAN (Zhu et al., 2017)

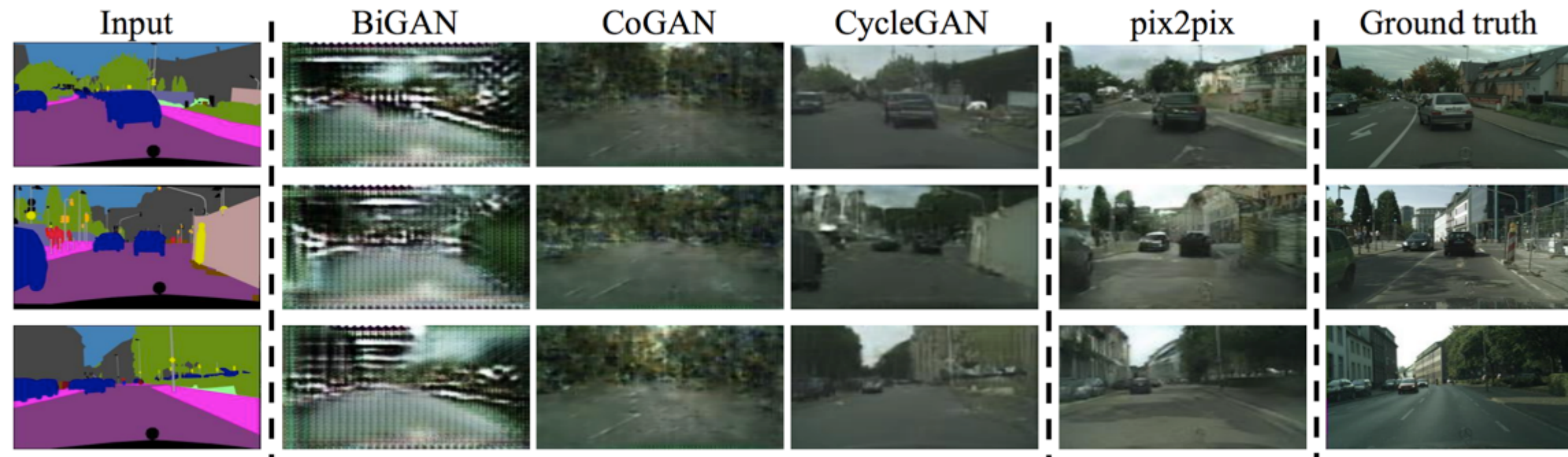
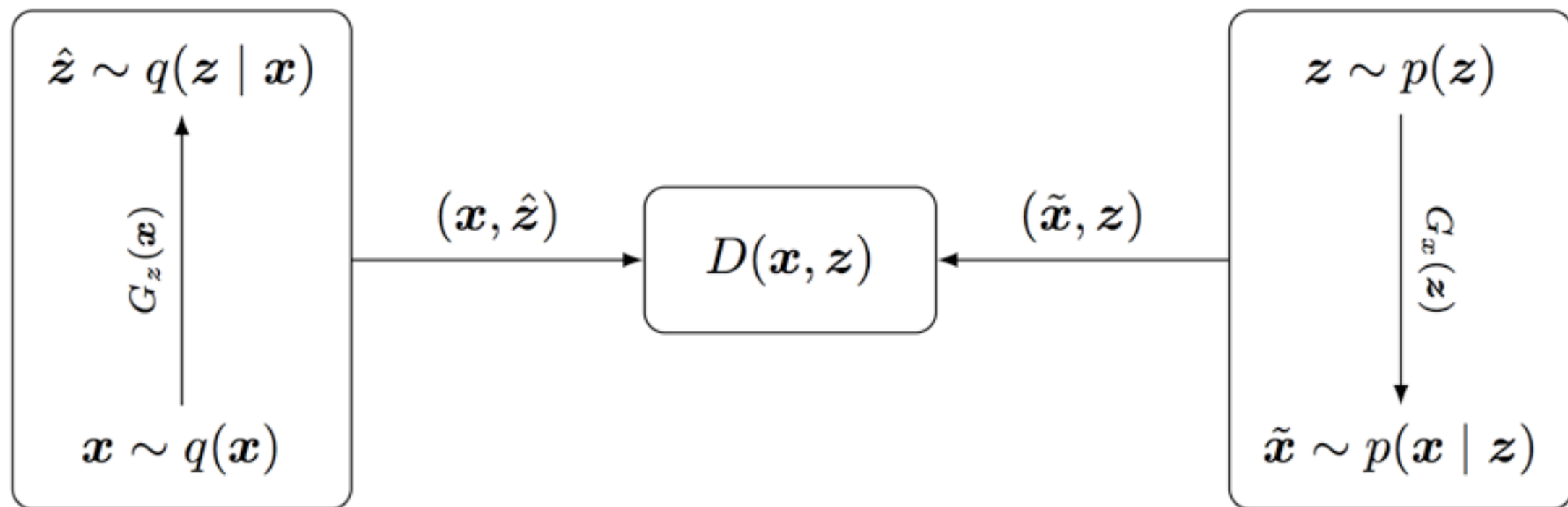


Figure 5: Different methods for mapping labels $\leftrightarrow$ photos trained on cityscapes. From left to right: input, BiGAN [5, 6], CoupledGAN [27], CycleGAN (ours), pix2pix [18] trained on paired data, and ground truth.

# Solution 2: Encoders-incorporated

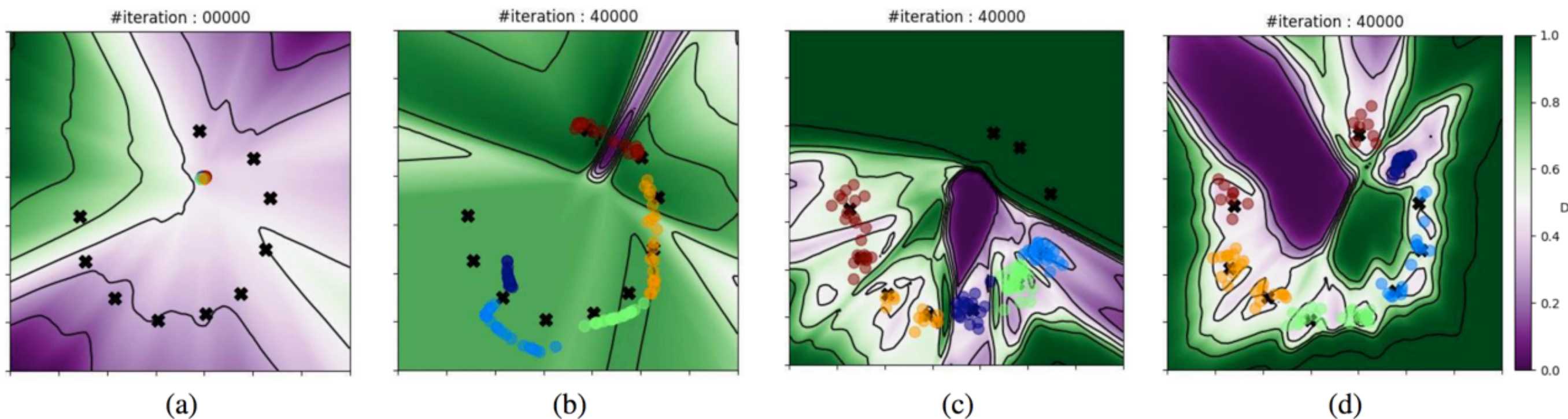
- CycleGAN (Zhu et al., 2017)
- BiGAN (Donahue et al., 2017; Dumoulin et al., 2017)
  - $G: Z \rightarrow X$  +  $F: X \rightarrow Z$





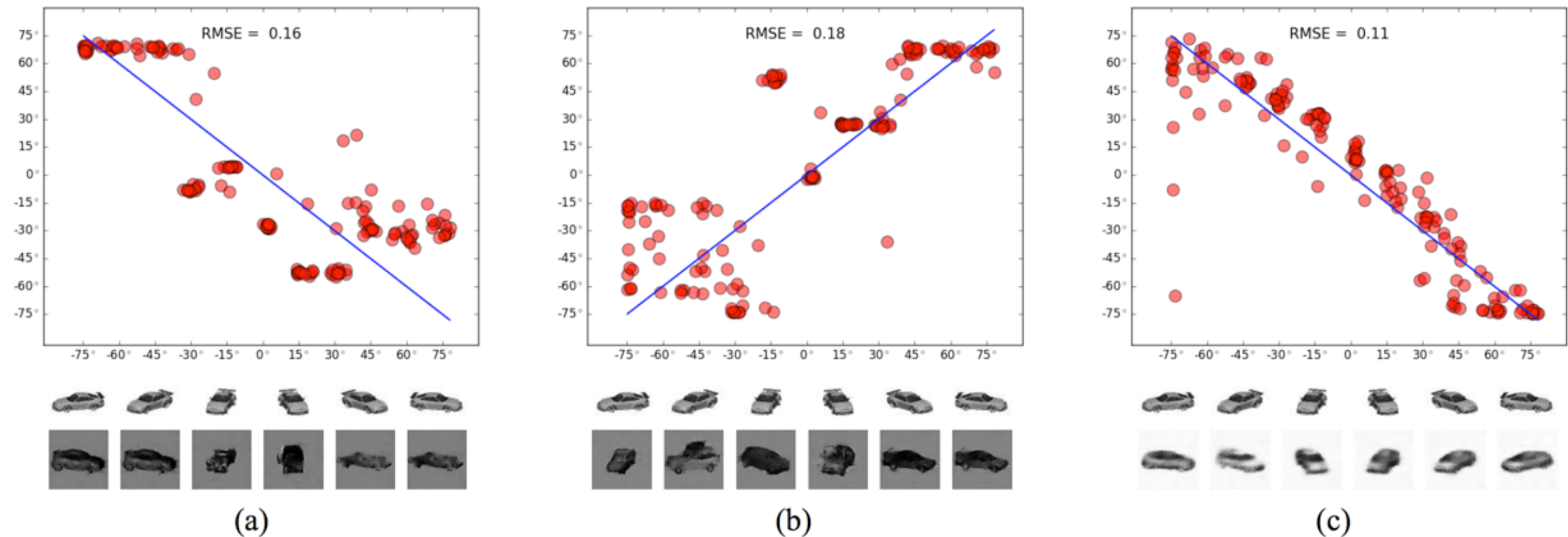
# Solution 2: Encoders-incorporated

- DiscoGAN (Kim et al., 2017)



# Solution 2: Encoders-incorporated

- DiscoGAN (Kim et al., 2017)



# Difficulty 3

Tackled!

- **minimizing** the *KL divergence* only **is biased**:

$$KL(\mathbb{P}_{g_\theta} || \mathbb{P}_r) - 2JSD(\mathbb{P}_{g_\theta} || \mathbb{P}_r)]$$

- because *KL divergence* is asymmetric, and thus it is not equally treated when **G** generates an unreal sample and when **G** fails to generate real sample
- Therefore, **G** will generate too many few-mode (less diverse) but real samples , a safer strategy



# Content

- Generative Adversarial Networks
  - Basics and Attractiveness
  - Difficulties
- Solution 1: Partial and Fine-grained Guidance
- Solution 2: Encoder-incorporated
- **Solution 3: Wasserstein Distance**

# Solution 3: Wasserstein Distance

- Wasserstein GANs (Arjovsky et al., 2017)
- Wasserstein-1 Distance (Earth-Mover Distance):

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

# Solution 3: Wasserstein Distance

- Wasserstein GANs (Arjovsky et al., 2017)
- Wasserstein-1 Distance (Earth-Mover Distance):

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

- *Why is it superior to KL and JS divergence?*



# Solution 3: Wasserstein Distance

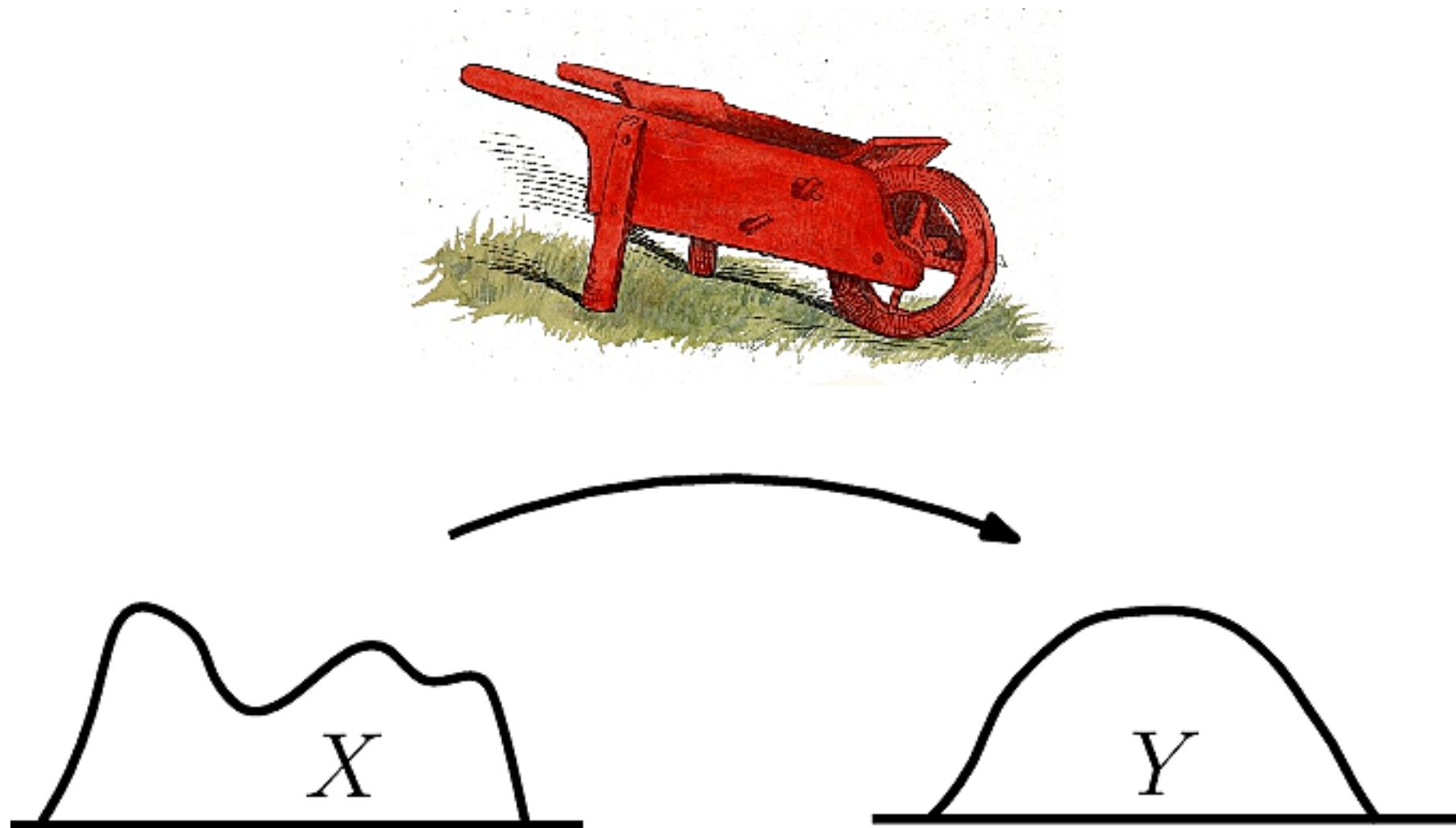
- Wasserstein-1 Distance (Earth-Mover Distance):

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

where  $\Pi(\mathbb{P}_r, \mathbb{P}_g)$  denotes the set of all joint distributions  $\gamma(x, y)$  whose marginals are respectively  $\mathbb{P}_r$  and  $\mathbb{P}_g$ . Intuitively,  $\gamma(x, y)$  indicates how much “mass” must be transported from  $x$  to  $y$  in order to transform the distributions  $\mathbb{P}_r$  into the distribution  $\mathbb{P}_g$ . The EM distance then is the “cost” of the optimal transport plan.

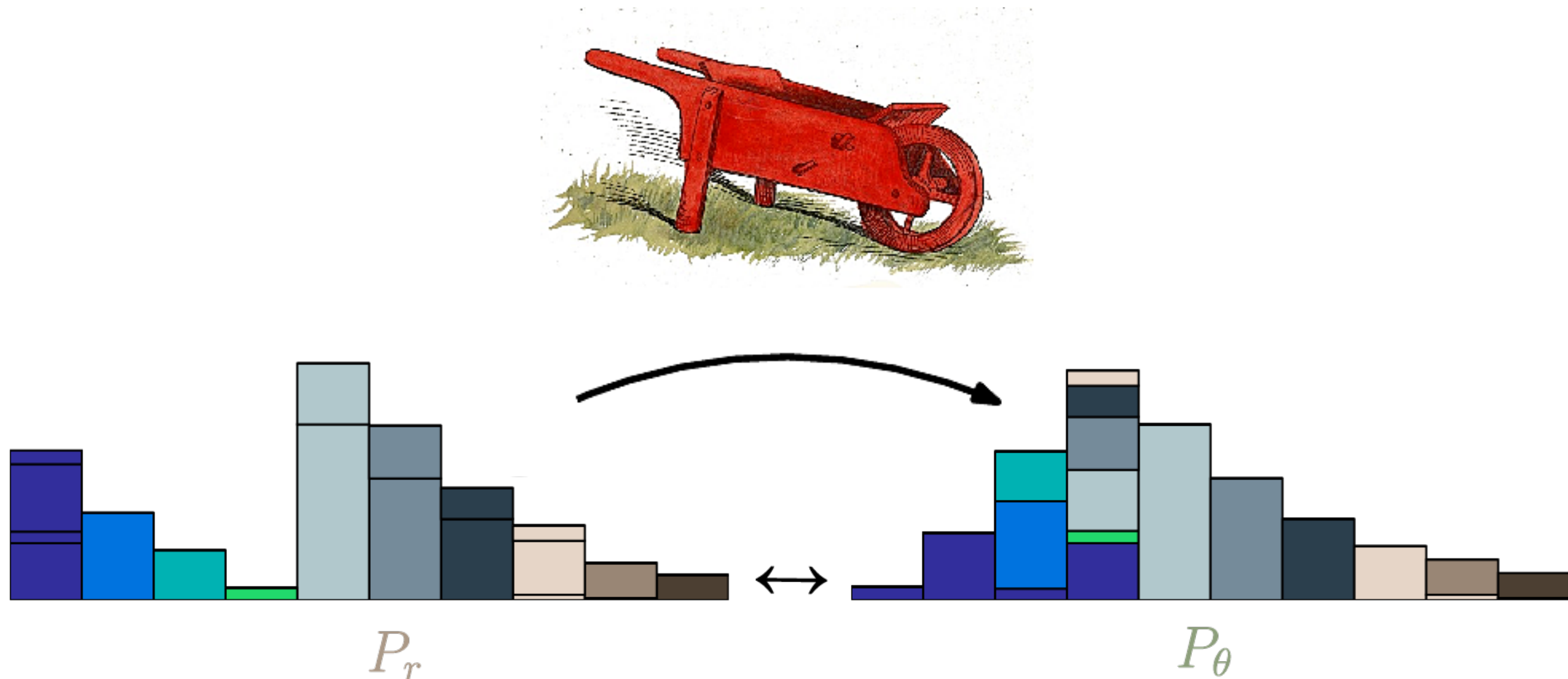
# Solution 3: Earth Move Distance

- Wasserstein-1 Distance (Earth-Mover Distance):



# Solution 3: Earth Move Distance

- Wasserstein-1 Distance (Earth-Mover Distance):





# Solution 3: Wasserstein Distance

- Wasserstein-1 Distance (Earth-Mover Distance):

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

- The distance is shown to have the desirable property that under mild assumptions
  - It is continuous everywhere and
  - differentiable almost everywhere.

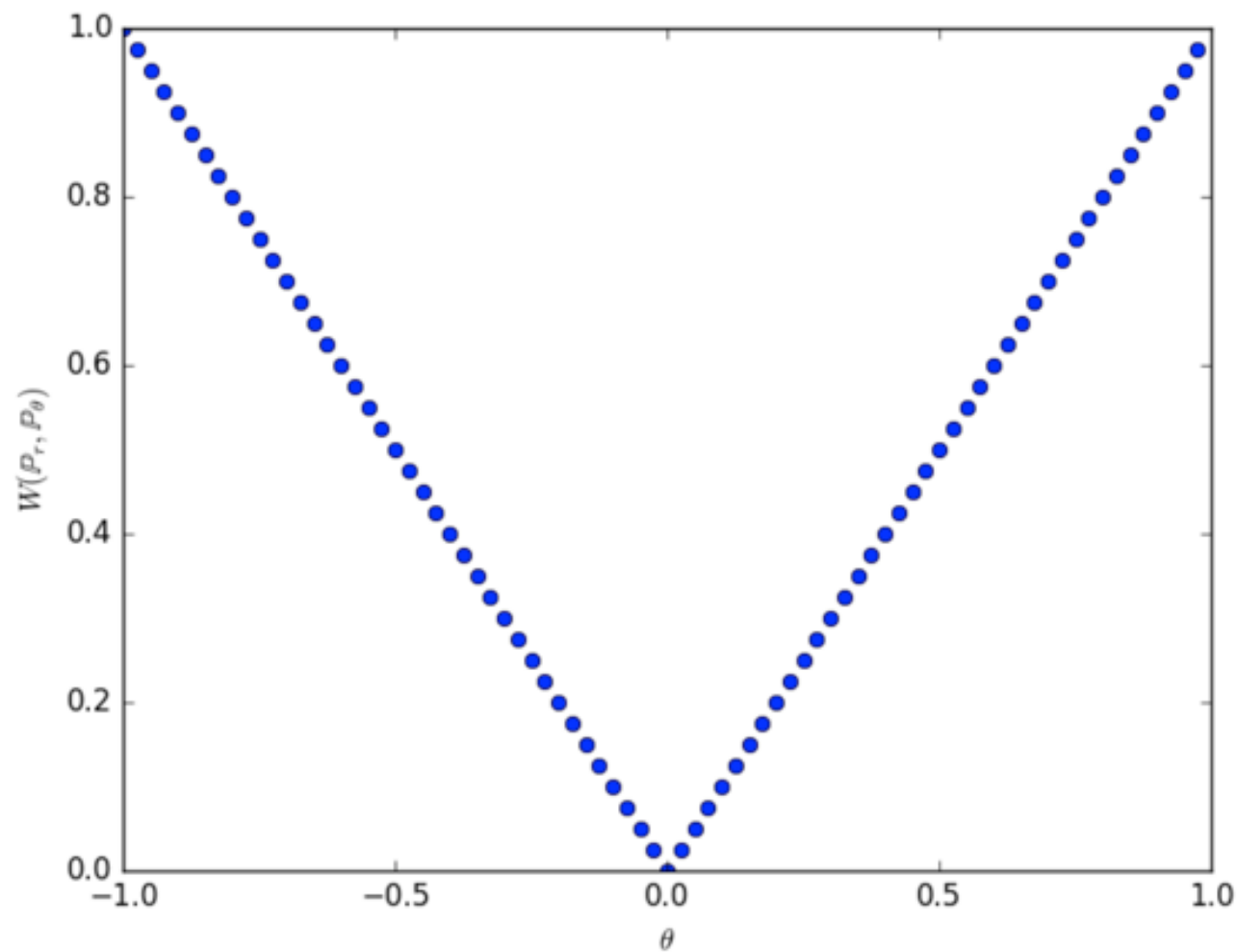
# Solution 3: Wasserstein Distance

- Wasserstein-1 Distance (Earth-Mover Distance):

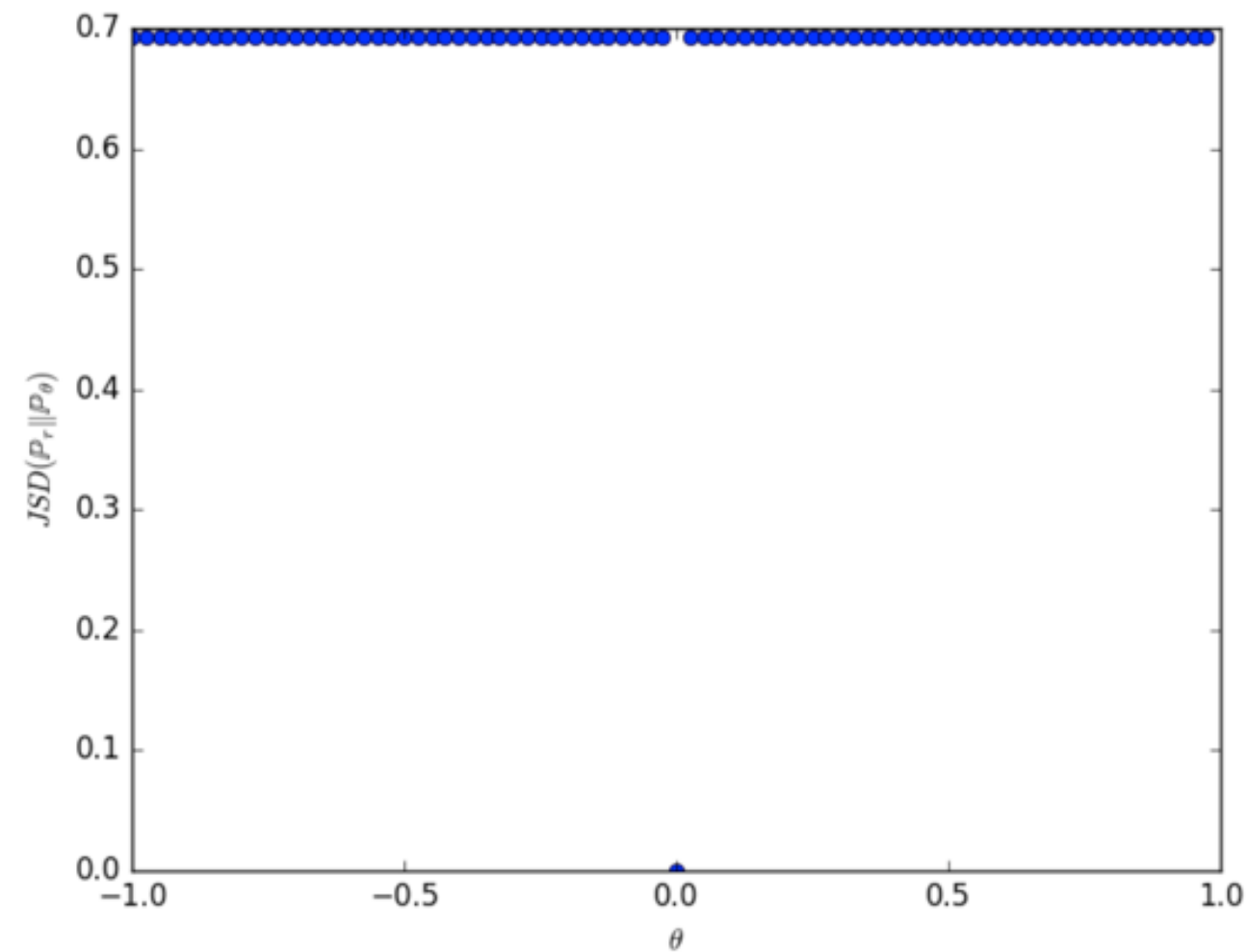
$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

- The distance is shown to have the desirable property that under mild assumptions
  - *And most importantly, it can reflect the distance of two distributions even if they do not overlap, and thus can provide meaningful gradients*

# Solution 3: Wasserstein Distance



Wasserstein Distance



JS Divergence



# Solution 3: Wasserstein Distance

- Wasserstein-1 Distance (Earth-Mover Distance):

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

- By applying the Kantorovich-Rubinstein duality (Villani, 2008), Wasserstein GANs becomes:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

# Wasserstein GANs

- This new value function of WGAN gives rise to the additional requirement that the discriminator must lie within in the space of 1-Lipschitz functions:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

- In other words,  **$\mathcal{D}$**  is the set of 1-Lipschitz functions
- Lipschitz continuity

# Lipschitz Continuity

- real-value function:  $f: R \rightarrow R$
- positive constant:  $K$

$$|f(x_1) - f(x_2)| \leq K|x_1 - x_2|$$

- In other words, a Lipschitz continuous function has bounded first derivative. Intuitively, the slope of a KK-Lipschitz function never exceeds KK, for a more general definition of slope.

$$d_Y(f(x_1), f(x_2)) \leq K d_X(x_1, x_2)$$

# Wasserstein GANs

- This new value function of WGAN gives rise to the additional requirement that the discriminator must lie within in the space of 1-Lipschitz functions:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

- To satisfy this requirement, WGAN enforces the weights of ***D*** lie within a compact space  $[-c, c]$  by applying **weight clipping**



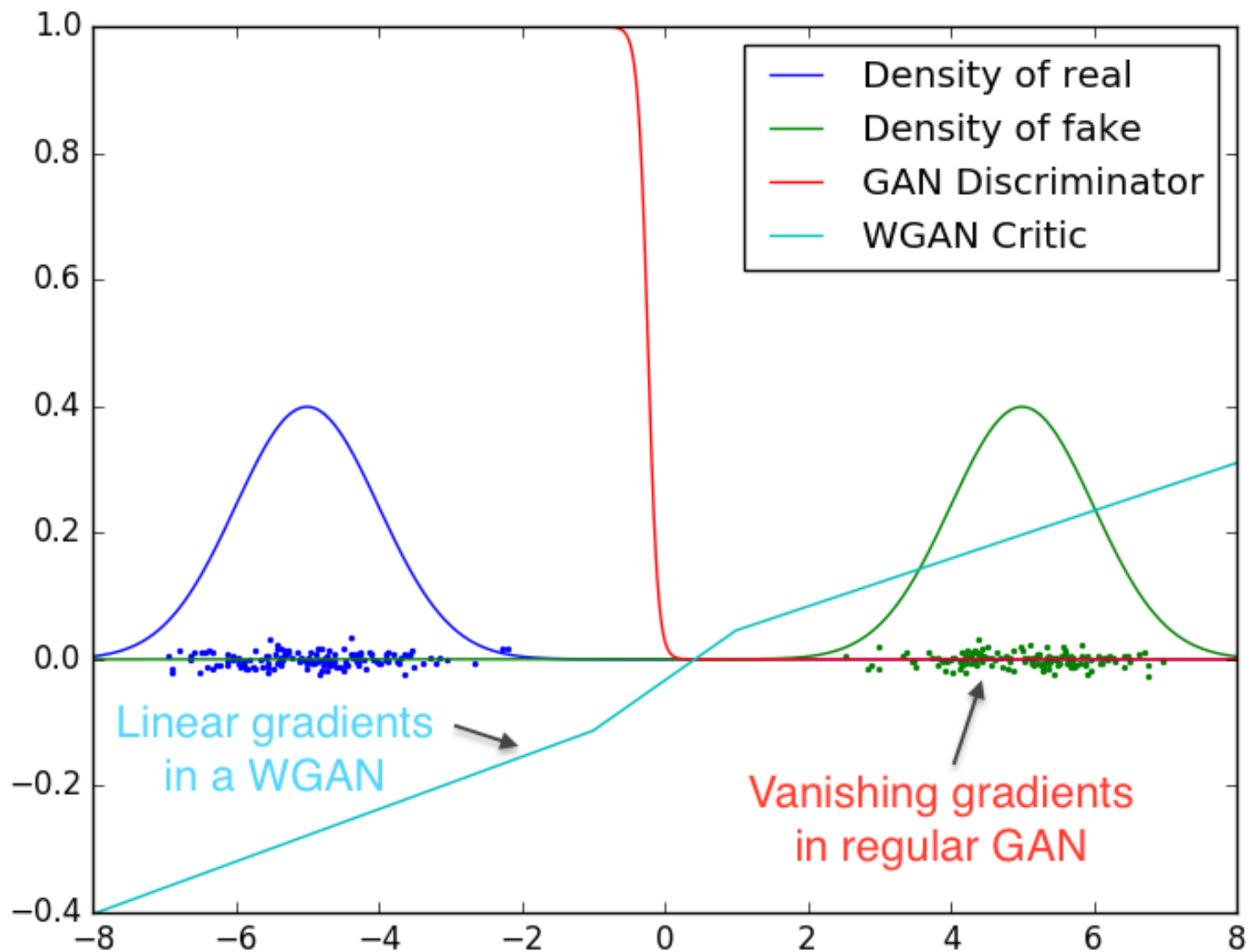
# Wasserstein GANs

- This new value function of WGAN gives rise to the additional requirement that the discriminator must lie within in the space of 1-Lipschitz functions:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})] - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})]$$

- Also, WGAN removes the sigmoid layer in ***D*** because by using Wasserstein distance, ***D*** in WGAN is doing regression rather than classification

# Wasserstein GANs



# Difficulty 2

**Tackled!**

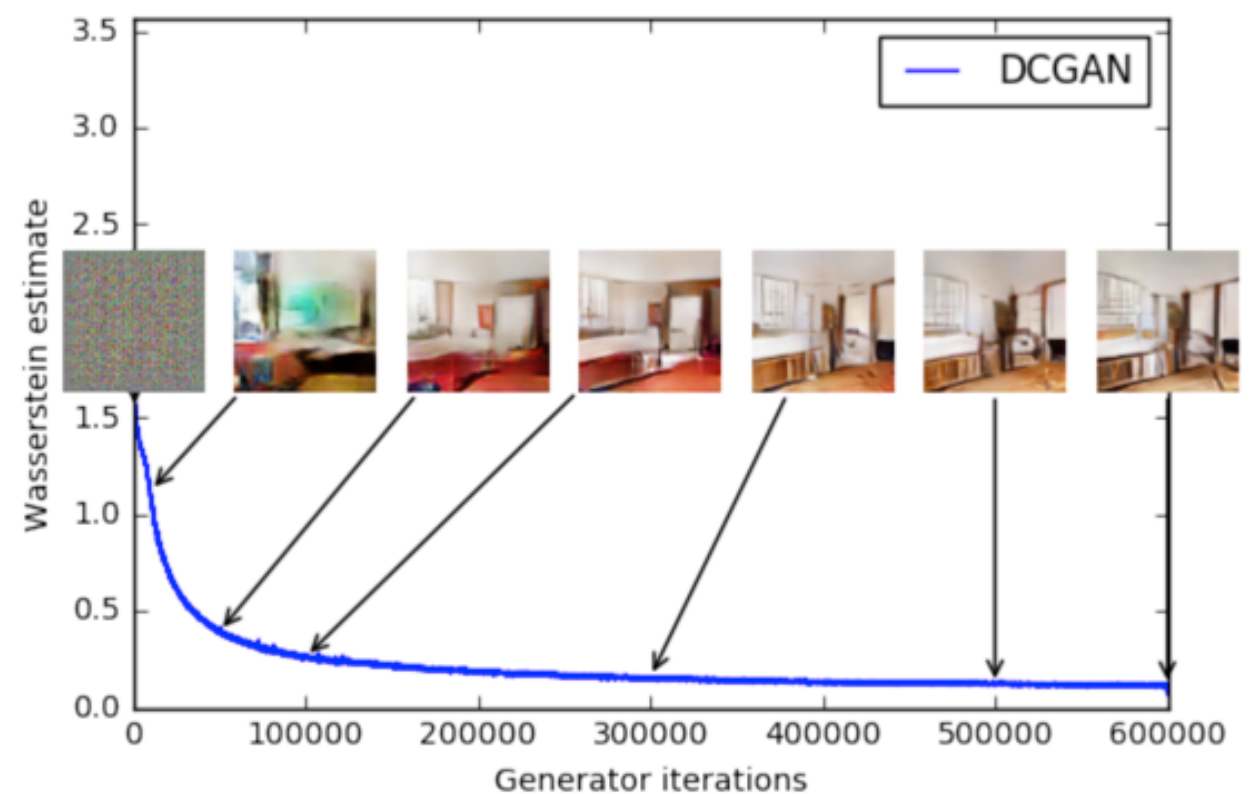
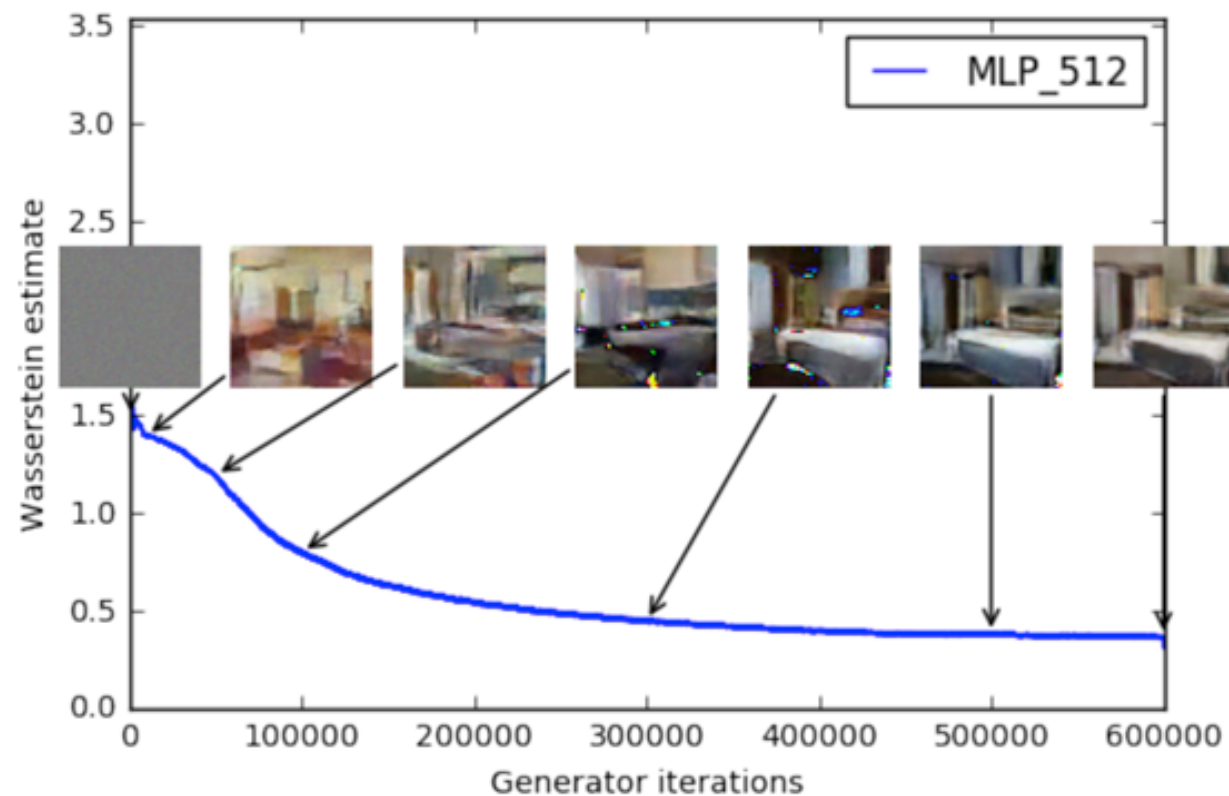
- when:

$$L(D^*, g_\theta) = 2JSD(\mathbb{P}_r \parallel \mathbb{P}_g) - 2 \log 2$$

- The JS divergence for the two distributions  $P_r$  and  $P_g$  is (almost) always  $\log 2$  because  $P_r$  and  $P_g$  hardly can overlap (Arjovsky & Bottou, 2017, Theorem 2.1~2.3)
- This results in vanishing gradient in theory!

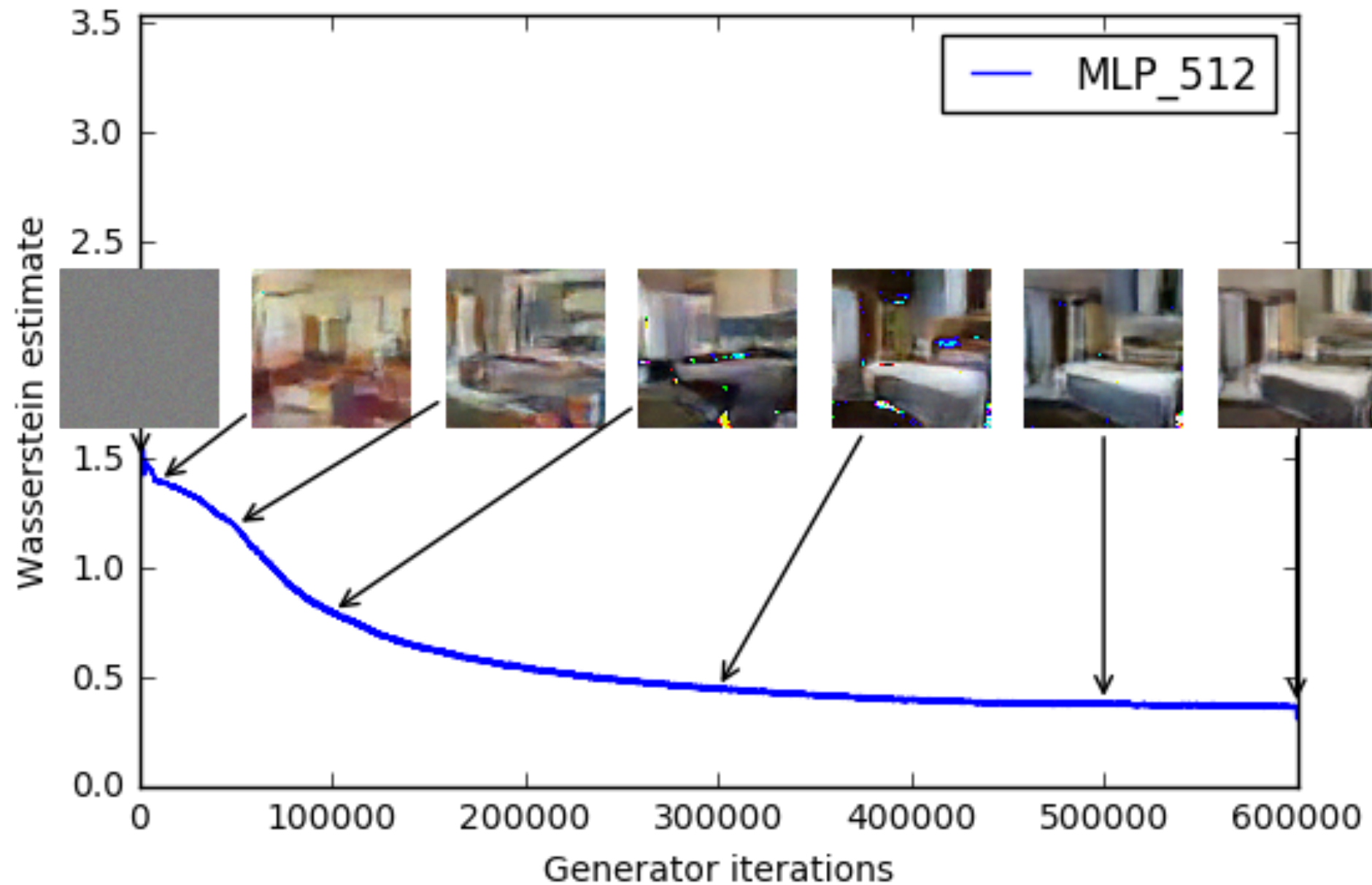
# Wasserstein GANs

- This new value function of WGAN seems correlate with the quality of the generated samples:





# Wasserstein GANs



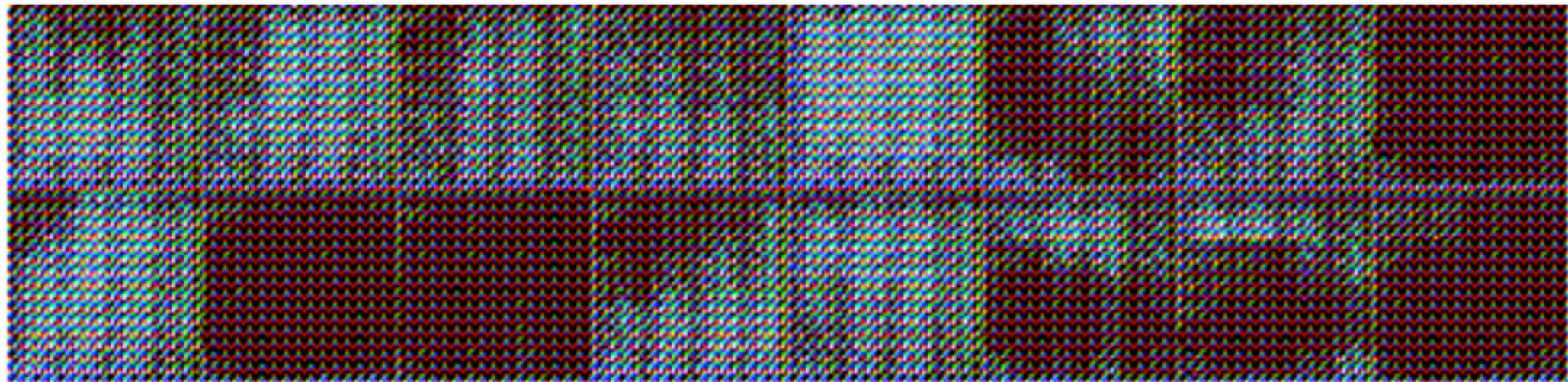
# Wasserstein GANs



*Top:* WGAN with the same DCGAN architecture. *Bottom:* DCGAN



# Wasserstein GANs



*Top:* WGAN with DCGAN architecture, no batch norm. *Bottom:* DCGAN, no batch norm.



# Wasserstein GANs

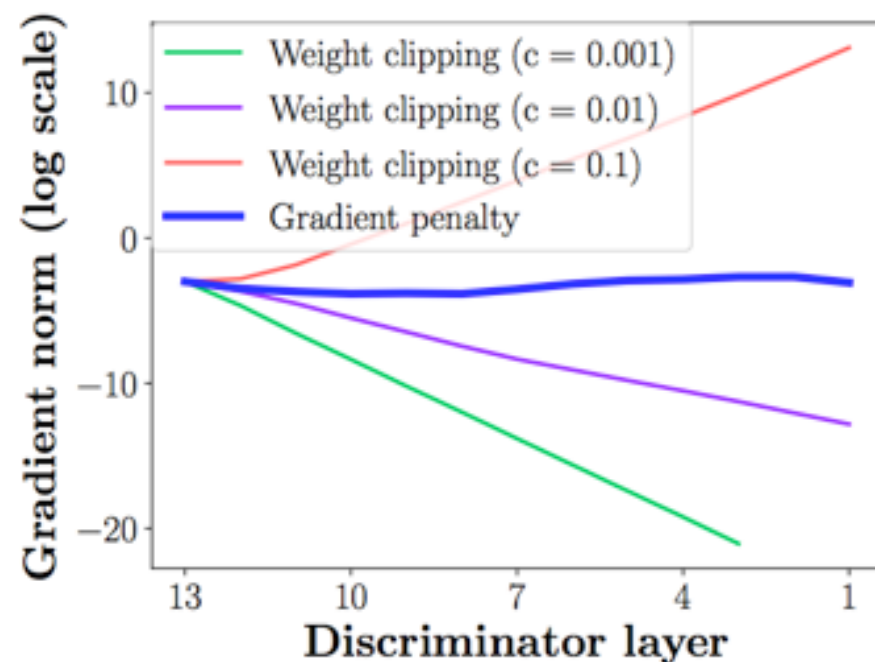


*Top: WGAN with MLP architecture. Bottom: Standard GAN, same architecture.*

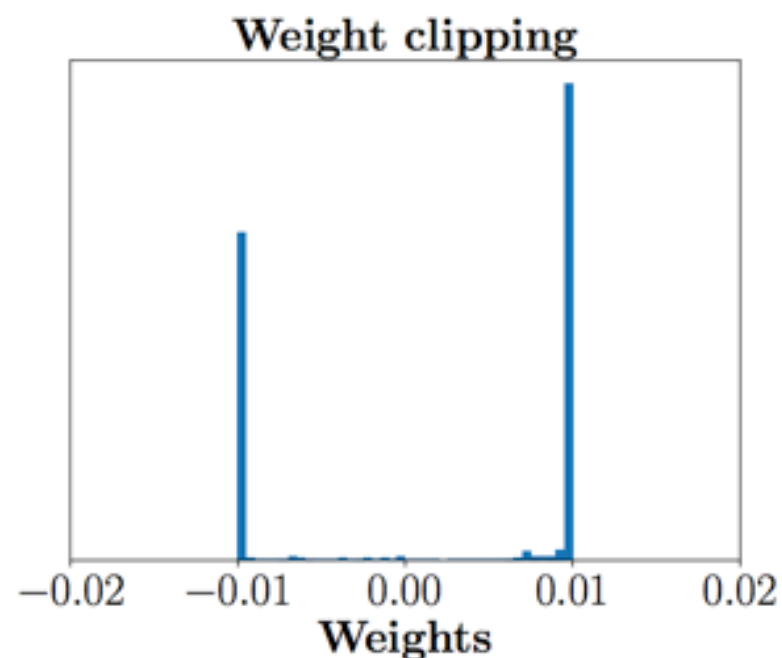


# Improved Wasserstein GANs

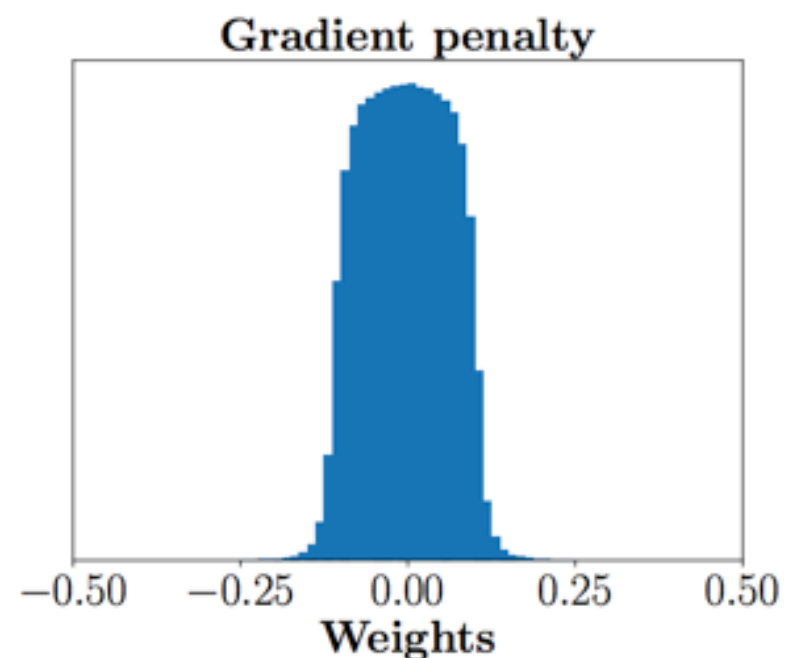
- The drawbacks of weight clipping



(a)



(b)



- bias ***D*** toward much simpler functions

# Improved Wasserstein GANs

- The drawbacks of weight clipping

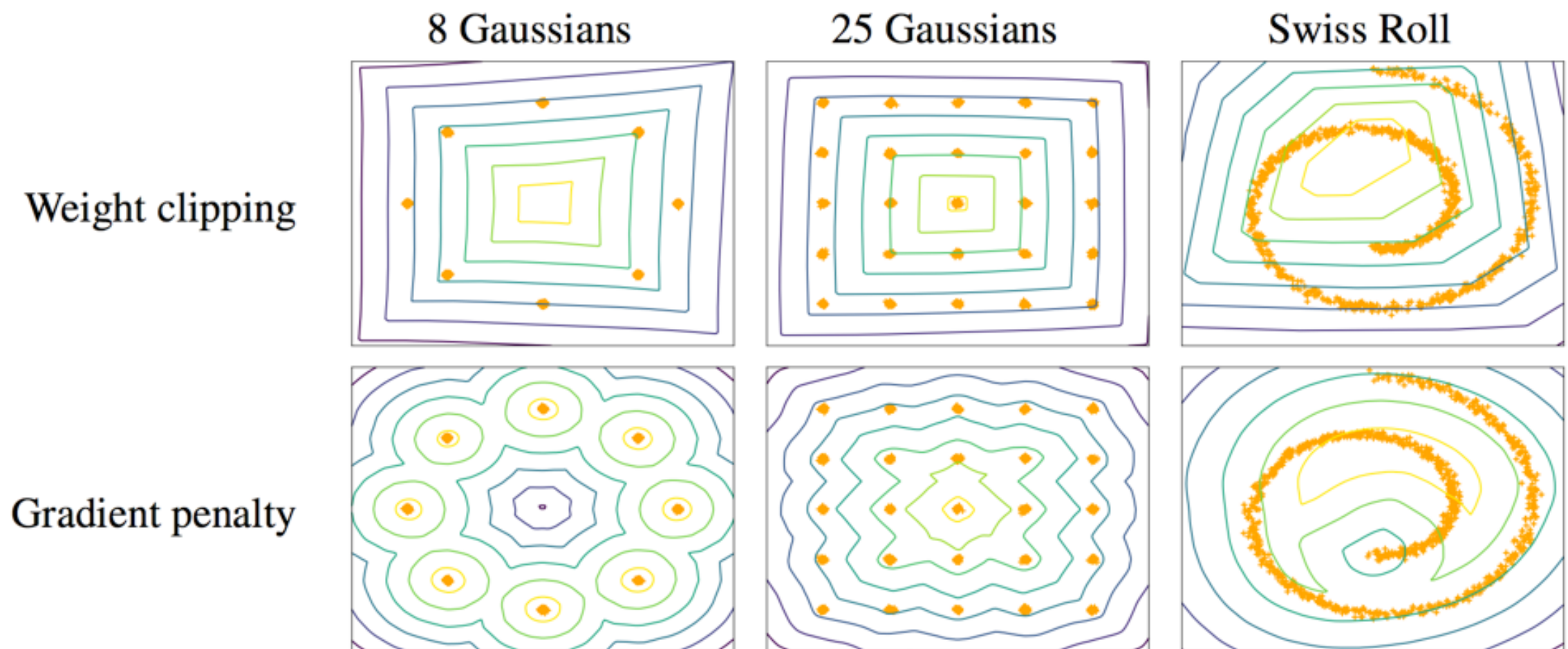
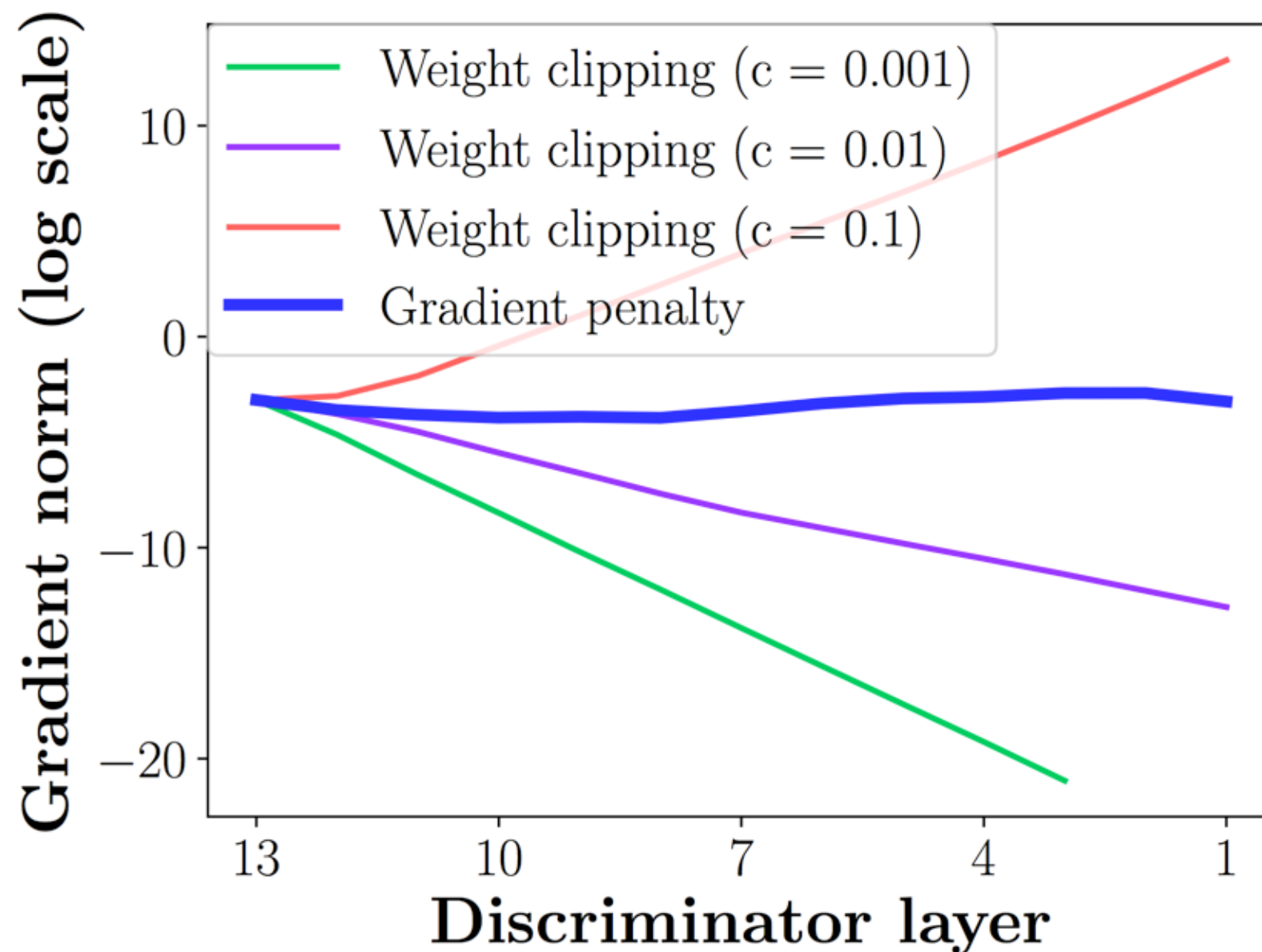


Figure 1: Value surfaces of WGAN critics trained to optimality on toy datasets. Critics trained with weight clipping fail to capture higher moments of the data distribution. The ‘generator’ is held fixed at the real data plus Gaussian noise.

# Improved Wasserstein GANs

- The drawbacks of weight clipping



# Improved Wasserstein GANs

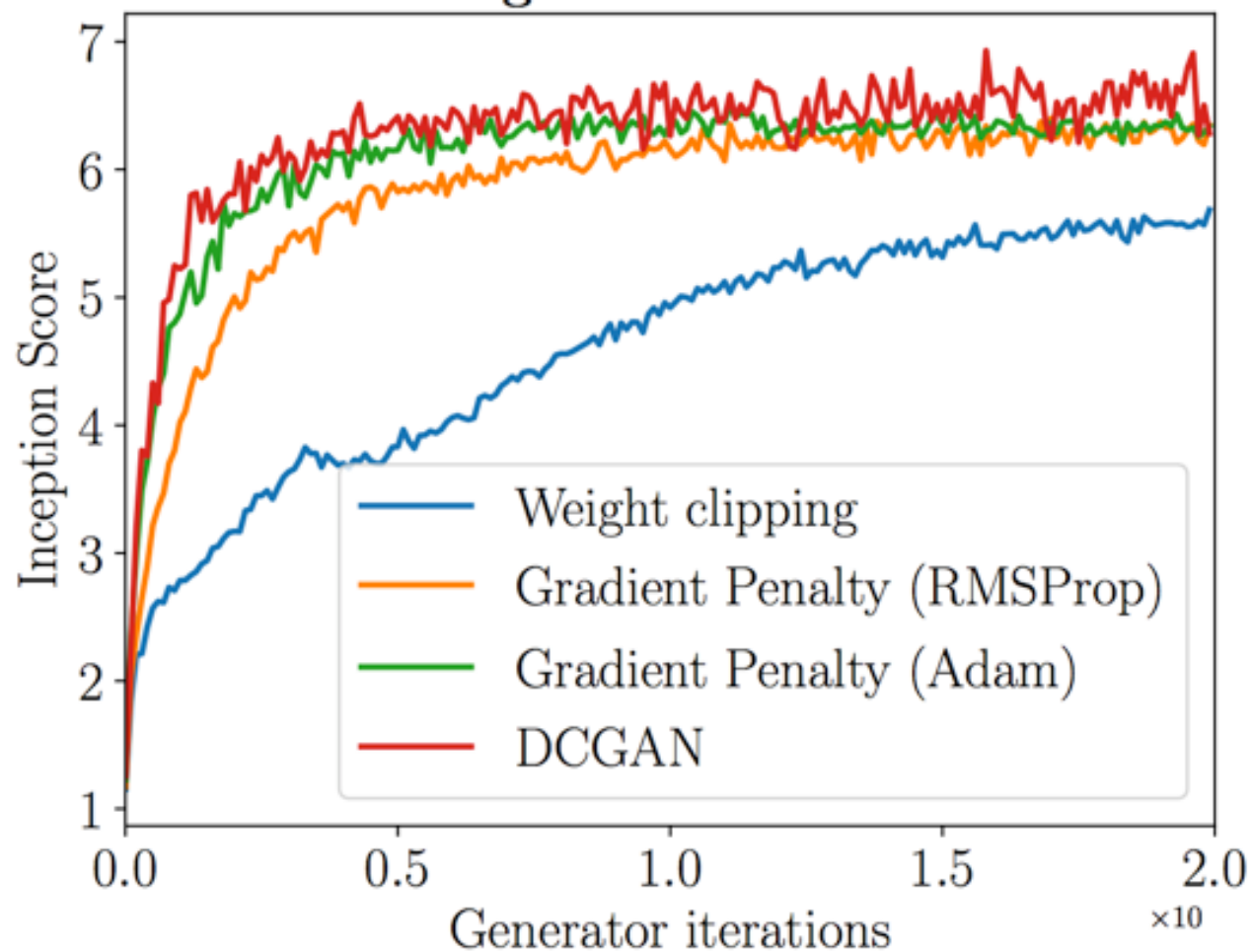
- The optimal ***D*** under WGAN:
  - has gradients with norm 1 almost everywhere under  $P_r$  and  $P_g$ .
- The objective of improved WGAN-GP:

$$L = \underbrace{\mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}} [(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2]}_{\text{Our gradient penalty}}$$

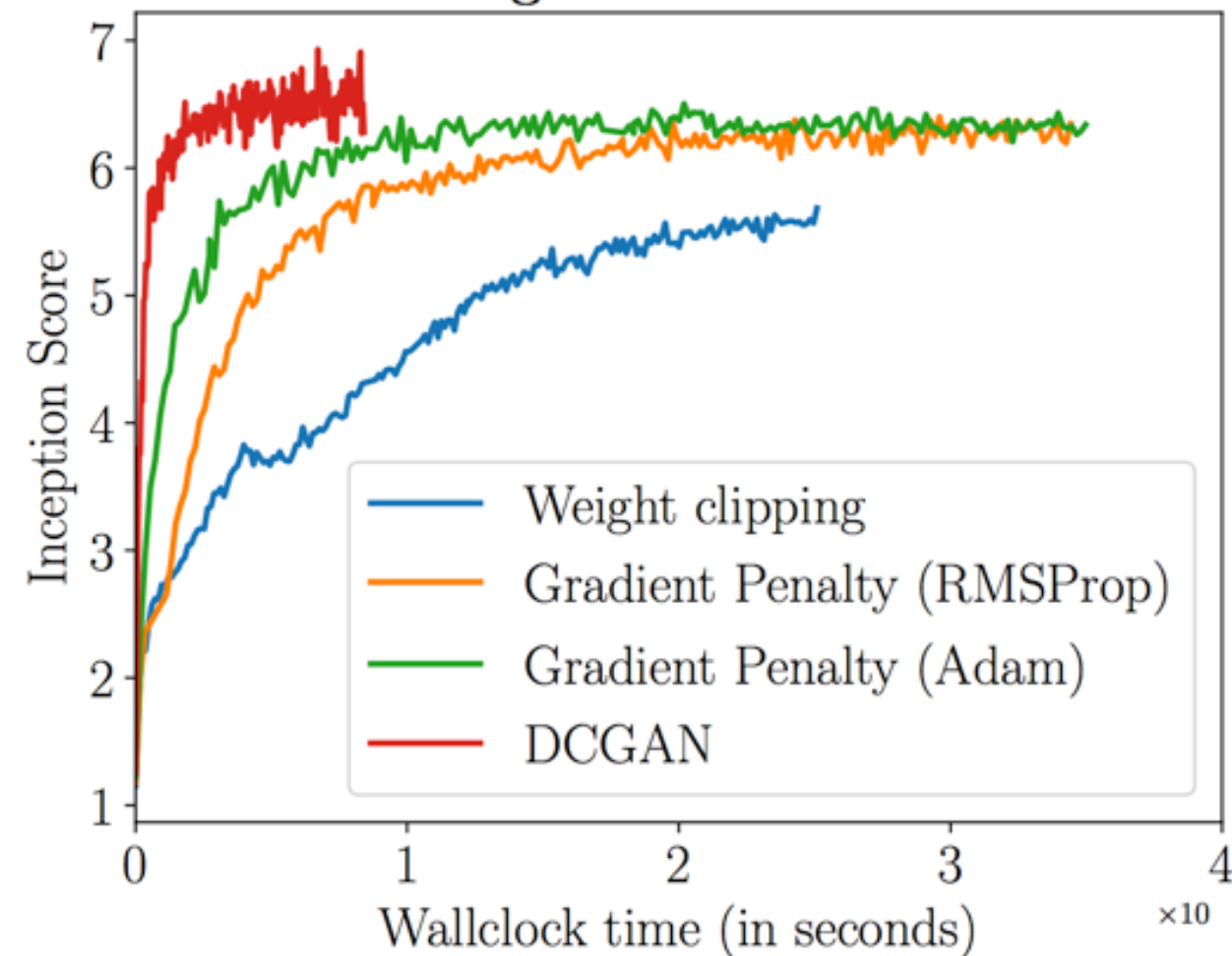


# Improved Wasserstein GANs

Convergence on CIFAR-10



Convergence on CIFAR-10



# Improved Wasserstein GANs

**DCGAN**

**LSGAN**

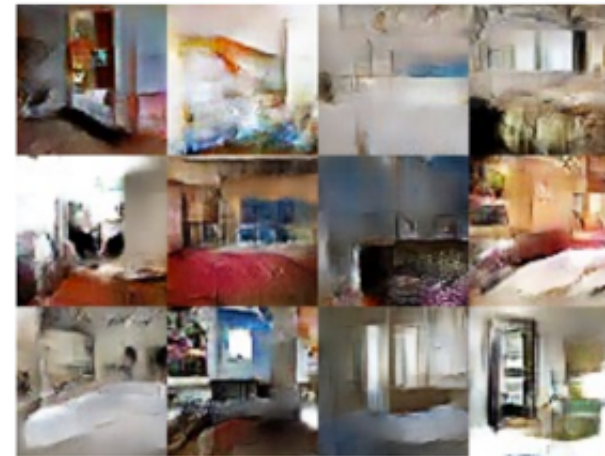
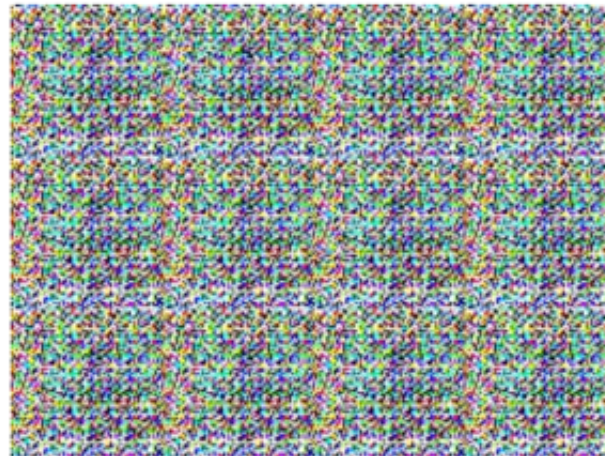
**WGAN (clipping)**

**WGAN-GP (ours)**

Baseline ( $G$ : DCGAN,  $D$ : DCGAN)



$G$ : No BN and a constant number of filters,  $D$ : DCGAN





# Take-home Messages

- Try WGAN-GP
- Try noisy input
- Try specific architecture (with careful analysis of the certain problem)
- Try different type of the noise
- Checklist here: <https://github.com/soumith/ganhacks>

Thanks for your attention!  
Any questions?





# References

- Arjovsky and Bottou, “Towards Principled Methods for Training Generative Adversarial Networks”. ICLR 2017.
- Goodfellow et al., “Generative Adversarial Networks”. ICLR 2014.
- Che et al., “Mode Regularized Generative Adversarial Networks”. ICLR 2017.
- Zhao et al., “Energy-based Generative Adversarial Networks”. ICLR 2017.
- Berthelot et al., “BEGAN: Boundary Equilibrium Generative Adversarial Networks”. arXiv preprint 2017.
- Sønderby, et al., “Amortised MAP Inference for Image Super-Resolution”. ICLR 2017.
- Arjovsky et al., “Wasserstein GANs”. ICML 2017.
- Villani, Cedric. “Optimal transport: old and new”, volume 338. Springer Science & Business Media, 2008

# References

- Jun-Yan Zhu\*, Taesung Park\*, Phillip Isola, Alexei A. Efros. “Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks”. arXiv preprint 2017.
- Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, Jiwon Kim. “Learning to Discover Cross-Domain Relations with Generative Adversarial Networks”. ICML 2017.
- Zili Yi, Hao Zhang, Ping Tan, Minglun Gong. “DualGAN: Unsupervised Dual Learning for Image-to-Image Translation”. arXiv preprint 2017.
- Jeff Donahue, Philipp Krähenbühl, Trevor Darrell. “Adversarial Feature Learning”. ICLR 2017.
- Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, Aaron Courville. “Adversarially Learned Inference”. ICLR 2017.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, Aaron Courville. “Improved Training of Wasserstein GANs”. arXiv preprint 2017.

# References

- Mehdi Mirza, Simon Osindero. “Conditional Generative Adversarial Nets”. arXiv preprint 2014.
- Emily Denton, Soumith Chintala, Arthur Szlam, Rob Fergus. “Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks”. arXiv preprint 2015.
- Alec Radford, Luke Metz, Soumith Chintala. “Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks”. ICLR 2016.
- Scott Reed, Zeynep Akata, Xincheng Yan, Lajanugen Logeswaran, Bernt Schiele, Honglak Lee. “Generative Adversarial Text to Image Synthesis”. ICML 2016.
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen. “Improved Techniques for Training GANs”. arXiv preprint 2016.
- Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaolei Huang, Xiaogang Wang, Dimitris Metaxas. “StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks”. arXiv preprint 2016.
- Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, Alexei A. Efros. “Image-to-Image Translation with Conditional Adversarial Networks”. CVPR 2017.

# References

- Yaniv Taigman, Adam Polyak, Lior Wolf. “Unsupervised Cross-Domain Image Generation”. ICLR 2017.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio. “Generative Adversarial Nets”. NIPS 2014.
- Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman and Alexei A. Efros. “Generative Visual Manipulation on the Natural Image Manifold”, ECCV 2016.
- Huikai Wu, Shuai Zheng, Junge Zhang, Kaiqi Huang. “GP-GAN: Towards Realistic High-Resolution Image Blending”. arXiv preprint 2017.
- Anh Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, Jason Yosinski. “Plug & Play Generative Networks: Conditional Iterative Generation of Images in Latent Space”. CVPR 2017.