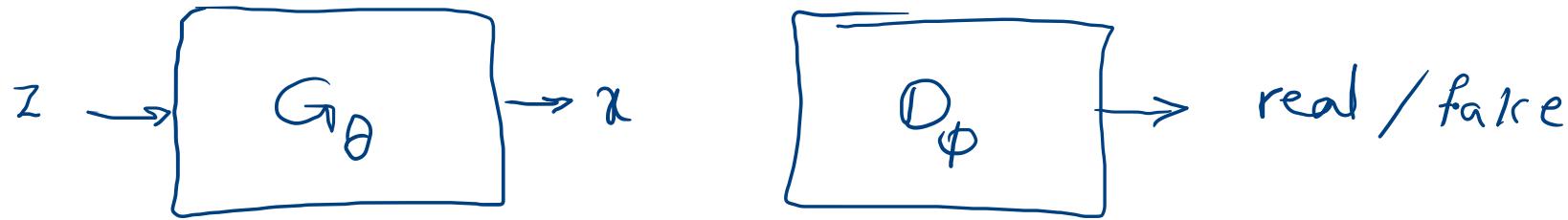


GAN

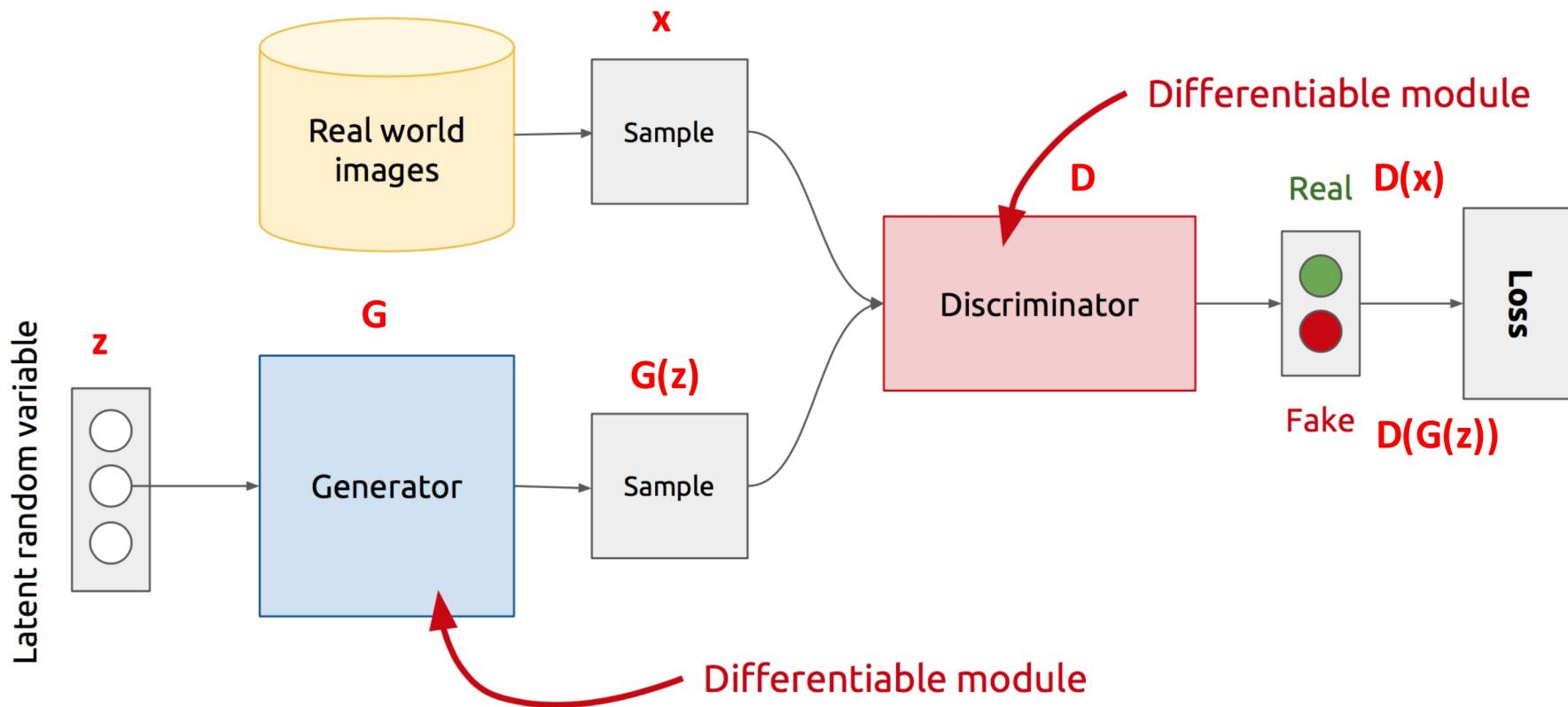
$$D = \{x_1, \dots, x_n\}$$



$$z \sim N(0, I)$$

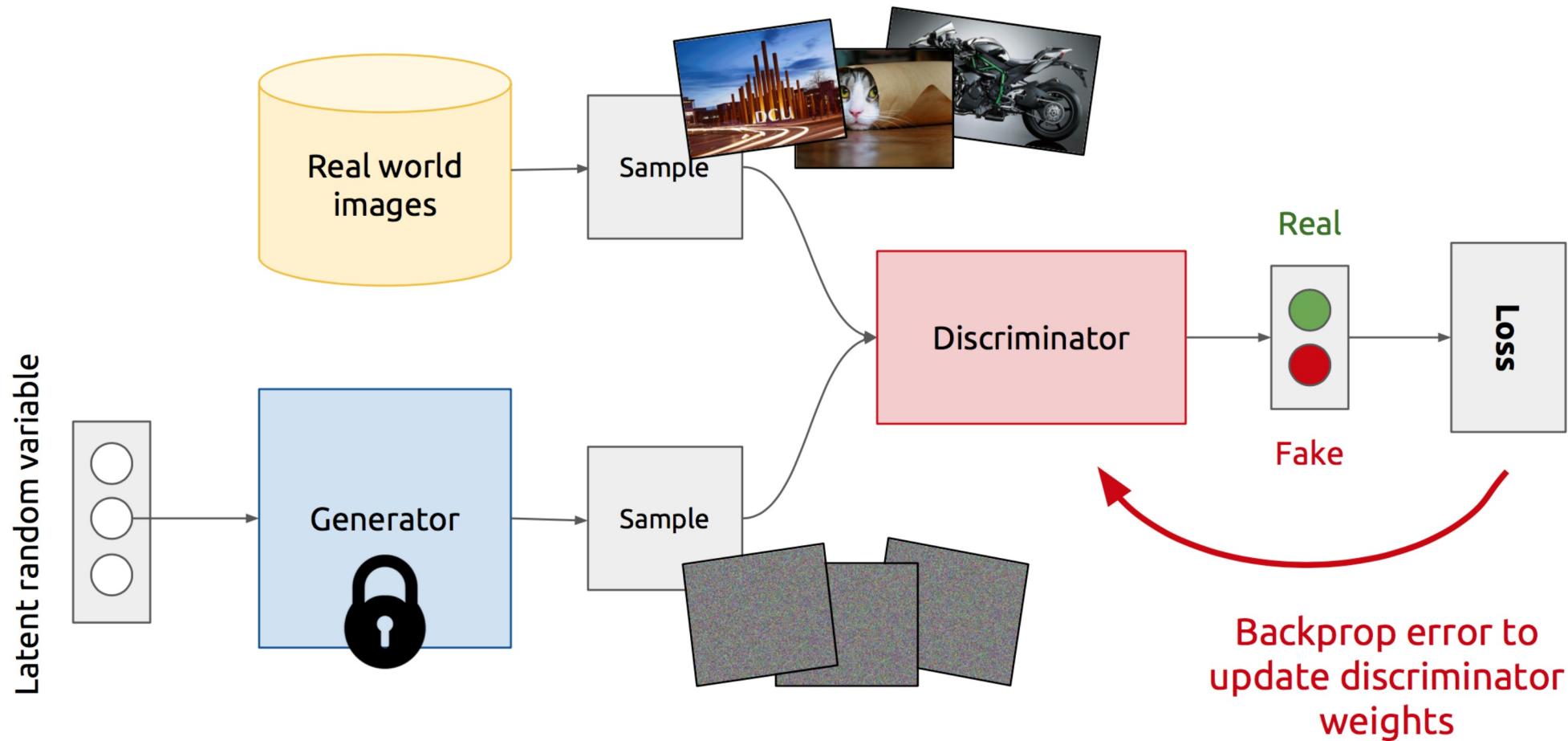
$$\max_{G_\theta} \min_{D_\phi} V(\phi, \theta) = E_{x \sim P_{\text{data}}} [\log D(x)] + E_{x \sim P_\theta} [\log (1 - D(x))]$$

GAN's Architecture

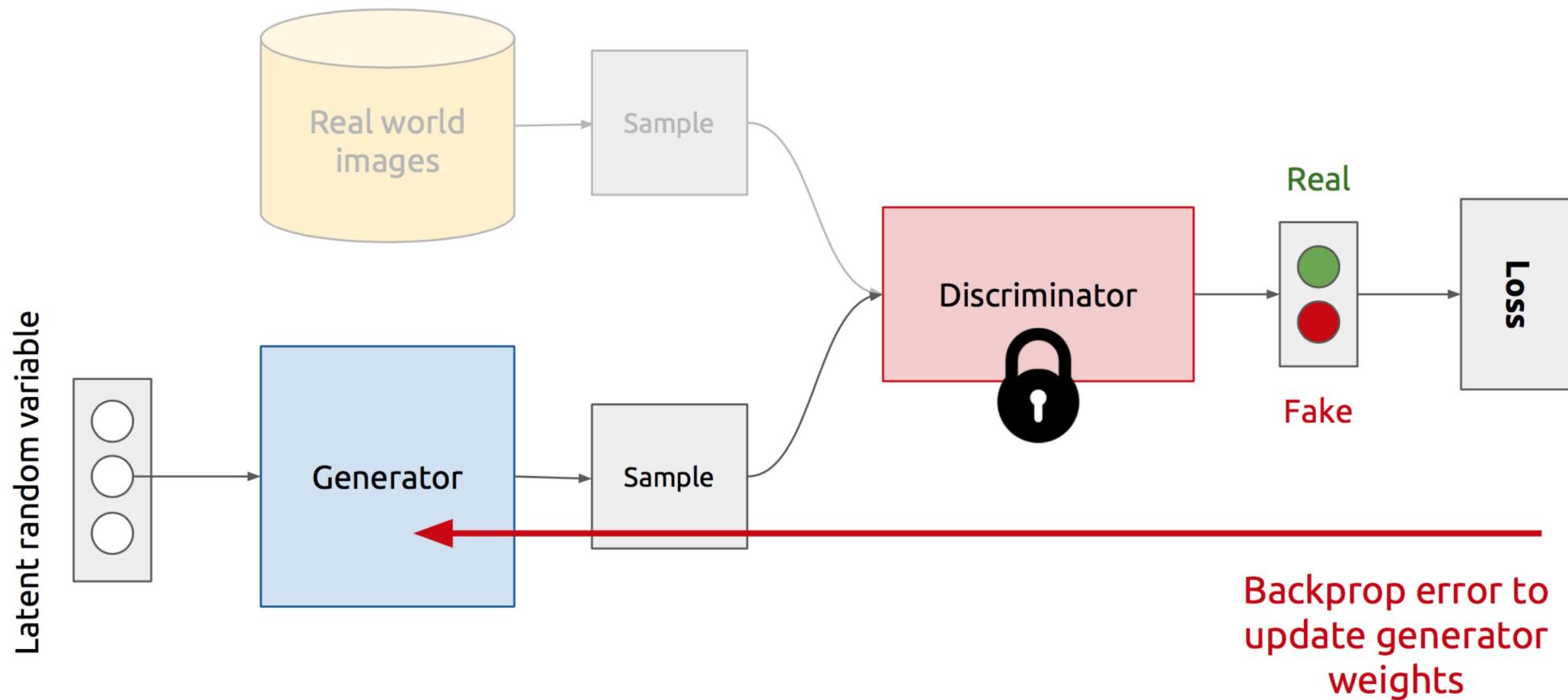


- z is some random noise (Gaussian/Uniform).
- z can be thought as the latent representation of the image.

Training Discriminator



Training Generator



Generative Adversarial Networks

Stefano Ermon

Stanford University

Lecture 10

Recap.

- Likelihood-free training
- Training objective for GANs

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

- With the optimal discriminator D_G^* , we see GAN minimizes a scaled and shifted Jensen-Shannon divergence

$$\min_G \underbrace{2D_{JSD}[p_{\text{data}}, p_G]}_{-\log 4}$$

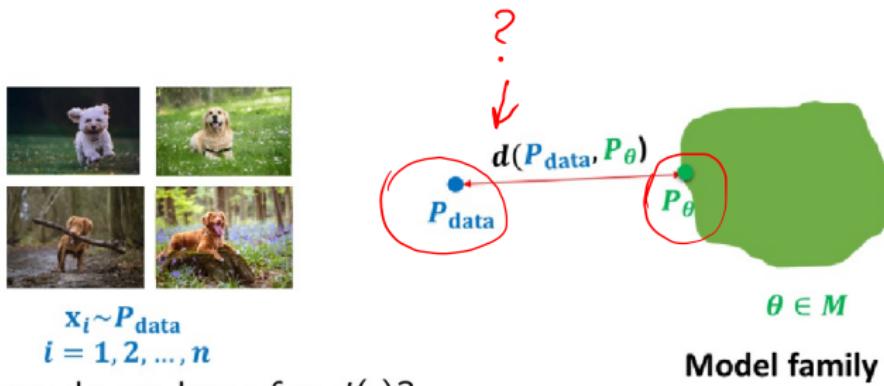
- Parameterize D by ϕ and G by θ . Prior distribution $p(\mathbf{z})$.

$$\min_{\theta} \max_{\phi} E_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\phi}(\mathbf{x})] + E_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_{\phi}(G_{\theta}(\mathbf{z})))]$$

Selected GANs

- <https://github.com/hindupuravinash/the-gan-zoo>
The GAN Zoo: List of all named GANs
- Today
 - Rich class of likelihood-free objectives via f -GANs
 - Wasserstein GAN
 - Inferring latent representations via BiGAN
 - Application: Image-to-image translation via CycleGANs

Beyond KL and Jenson-Shannon Divergence



What choices do we have for $d(\cdot)$?

- KL divergence: Autoregressive Models, Flow models
- (scaled and shifted) Jensen-Shannon divergence (approximately): original GAN objective

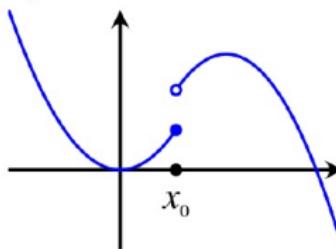
f divergences

- Given two densities p and q , the f -divergence is given by

$$D_f(p, q) = E_{\mathbf{x} \sim q} \left[f \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right]$$

where f is any convex, lower-semicontinuous function with $f(1) = 0$.

- Convex: Line joining any two points lies above the function
- Lower-semicontinuous: function value at any point \mathbf{x}_0 is close to $f(\mathbf{x}_0)$ or greater than $f(\mathbf{x}_0)$

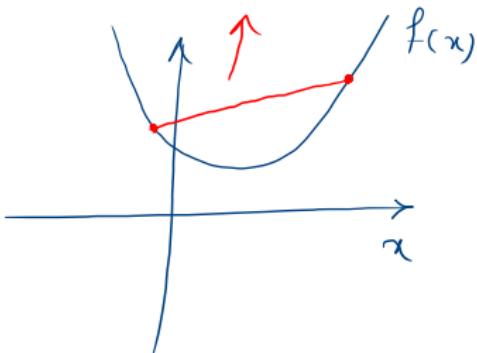


- Jensen's inequality: $E_{\mathbf{x} \sim q}[f(p(\mathbf{x})/q(\mathbf{x}))] \geq f(E_{\mathbf{x} \sim q}[p(\mathbf{x})/q(\mathbf{x})]) = f(\int q(x)p(x)/q(x)) = f(\int p(x)) = f(1) = 0$
- Example: KL divergence with $f(u) = u \log u$

$$D(p, q) = E_q \left[f\left(\frac{p(x)}{q(x)}\right) \right] \geq 0$$

① $D(p, q) \geq 0$

② $D(p, q) = 0 \iff p = q$



Jensen inequality,

$$E[f(x)] \geq f(E[x])$$

$$D(p, q) = E_q \left[f\left(\frac{p(x)}{q(x)}\right) \right] \geq f\left(E_q \left[\frac{p(x)}{q(x)}\right]\right)$$

$$= f\left(\int g(x) \frac{p(x)}{g(x)} dx\right) = 0$$

$$D(p, q) = E_q \left[f\left(\frac{p(x)}{q(x)}\right) \right]$$

$$f(u) = u \log u$$

$$D(p, q) = E_q \left[\frac{p(x)}{q(x)} \log \frac{p(x)}{q(x)} \right] = \int q(x) \frac{p(x)}{q(x)} \log \frac{p(x)}{q(x)} dx$$

$$= \int p(x) \log \frac{p(x)}{q(x)} dx = KL(p \parallel q)$$

$$f(u) = -\log u \rightarrow \frac{-1}{u} \rightarrow \frac{1}{u^2}$$

$$D(p, q) = KL(q \parallel p)$$

f divergences

Many more f-divergences!

Name	$D_f(P\ Q)$	Generator $f(u)$
Total variation	$\frac{1}{2} \int p(x) - q(x) dx$	$\frac{1}{2} u - 1 $
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} dx$	$u \log u$
Reverse Kullback-Leibler	$\int q(x) \log \frac{q(x)}{p(x)} dx$	$-\log u$
Pearson χ^2	$\int \frac{(q(x) - p(x))^2}{p(x)} dx$	$(u - 1)^2$
Neyman χ^2	$\int \frac{(p(x) - q(x))^2}{q(x)} dx$	$\frac{(1-u)^2}{u}$
Squared Hellinger	$\int \left(\sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx$	$(\sqrt{u} - 1)^2$
Jeffrey	$\int (p(x) - q(x)) \log \left(\frac{p(x)}{q(x)} \right) dx$	$(u - 1) \log u$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx$	$-(u + 1) \log \frac{1+u}{2} + u \log u$
Jensen-Shannon-weighted	$\int p(x)\pi \log \frac{p(x)}{\pi p(x)+(1-\pi)q(x)} + (1-\pi)q(x) \log \frac{q(x)}{\pi p(x)+(1-\pi)q(x)} dx$	$\pi u \log u - (1 - \pi + \pi u) \log(1 - \pi + \pi u)$
GAN	$\int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx - \log(4)$	$u \log u - (u + 1) \log(u + 1)$
α -divergence ($\alpha \notin \{0, 1\}$)	$\frac{1}{\alpha(\alpha-1)} \int \left(p(x) \left[\left(\frac{q(x)}{p(x)} \right)^\alpha - 1 \right] - \alpha(q(x) - p(x)) \right) dx$	$\frac{1}{\alpha(\alpha-1)} (u^\alpha - 1 - \alpha(u - 1))$

Source: Nowozin et al., 2016

Training with f -divergences

- Given p_{data} and p_θ , we could minimize $D_f(p_\theta, p_{data})$ or $D_f(p_{data}, p_\theta)$ as learning objectives. Non-negative, and zero if $p_\theta = p_{data}$
- However, it depends on the density ratio which is unknown

f -GAN

$$D_f(p_\theta, p_{data}) = \underbrace{E_{\mathbf{x} \sim p_{data}}}_{\text{approx w. samples}} \left[f \left(\underbrace{\frac{p_\theta(\mathbf{x})}{p_{data}(\mathbf{x})}}_{\text{unknown ratio}} \right) \right]$$

$$D_f(p_{data}, p_\theta) = \underbrace{E_{\mathbf{x} \sim p_\theta}}_{\text{approx w. samples}} \left[f \left(\underbrace{\frac{p_{data}(\mathbf{x})}{p_\theta(\mathbf{x})}}_{\text{unknown ratio}} \right) \right]$$

- To use f -divergences as a two-sample test objective for likelihood-free learning, we need to be able to estimate the objective using only samples (e.g., training data and samples from the model)

Towards Variational Divergence Minimization

- Fenchel conjugate: For any function $f(\cdot)$, its convex conjugate is

$$f^*(t) = \sup_{u \in \text{dom}_f} (ut - f(u))$$

where dom_f is the domain of the function f

- f^* is convex (pointwise supremum of convex functions is convex) and lower semi-continuous.
- Let f^{**} be the Fenchel conjugate of f^*

$$f^{**}(u) = \sup_{t \in \text{dom}_{f^*}} (tu - f^*(t))$$

- $f^{**} \leq f$. Proof: By definition, for all t, u

$$f^*(t) \geq ut - f(u) \quad \text{or equivalently} \quad f(u) \geq ut - f^*(t)$$

$$f(u) \geq \sup_t (ut - f^*(t)) = f^{**}(u)$$

- Strong Duality: $f^{**} = f$ when $f(\cdot)$ is convex, lower semicontinuous.

f -GAN: Variational Divergence Minimization

- We obtain a lower bound to an f -divergence via Fenchel conjugate

$$\begin{aligned} D_f(p, q) &= E_{\mathbf{x} \sim q} \left[f \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right] = E_{\mathbf{x} \sim q} \left[f^{**} \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right] \\ f^{**} &\stackrel{\text{def}}{=} E_{\mathbf{x} \sim q} \left[\sup_{t \in \text{dom}_{f^*}} \left(t \frac{p(\mathbf{x})}{q(\mathbf{x})} - f^*(t) \right) \right] \\ &= E_{\mathbf{x} \sim q} \left[T^*(x) \frac{p(\mathbf{x})}{q(\mathbf{x})} - f^*(T^*(x)) \right] \\ &= \int_{\mathcal{X}} q(x) \left[T^*(x) \frac{p(x)}{q(x)} - f^*(T^*(x)) \right] dx \\ &= \int_{\mathcal{X}} [T^*(x)p(x) - f^*(T^*(x))q(x)] dx \\ &= \sup_{\mathcal{T}} \int_{\mathcal{X}} [T(x)p(x) - f^*(T(x))q(x)] dx \\ &\geq \sup_{T \in \mathcal{T}} \int_{\mathcal{X}} (T(x)p(x) - f^*(T(x))q(x)) dx \\ &= \sup_{T \in \mathcal{T}} (E_{\mathbf{x} \sim p} [T(\mathbf{x})] - E_{\mathbf{x} \sim q} [f^*(T(\mathbf{x}))]) \end{aligned}$$

where $\mathcal{T} : \mathcal{X} \mapsto \mathbb{R}$ is an arbitrary class of functions

- Note:** Lower bound is likelihood-free w.r.t. p and q

f -GAN: Variational Divergence Minimization

- Variational lower bound

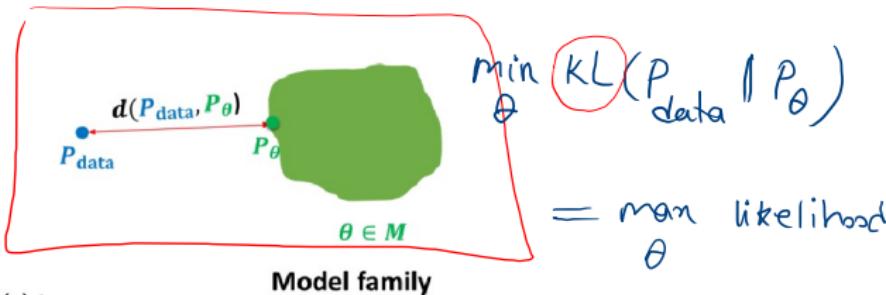
$$D_f(p, q) \geq \sup_{T \in \mathcal{T}} (E_{x \sim p} [T(x)] - E_{x \sim q} [f^*(T(x))])$$

- Choose any f -divergence
- Let $p = p_{\text{data}}$ and $q = p_G$
- Parameterize T by ϕ and G by θ
- Consider the following f -GAN objective

$$\min_{\theta} \max_{\phi} F(\theta, \phi) = E_{x \sim p_{\text{data}}} [T_{\phi}(x)] - E_{x \sim p_{G_{\theta}}} [f^*(T_{\phi}(x))]$$

- Generator G_{θ} tries to minimize the divergence estimate and discriminator T_{ϕ} tries to tighten the lower bound
- Substitute any f -divergence and optimize the f -GAN objective

Beyond KL and Jenson-Shannon Divergence



What choices do we have for $d(\cdot)$?

- KL divergence: Autoregressive Models, Flow models
- (scaled and shifted) Jensen-Shannon divergence (approximately): via the original GAN objective
- Any other f -divergence (approximately): via the f -GAN objective

Problems with GANs

- **Probability Distribution is Implicit**
 - Not straightforward to compute $P(X)$.
 - Thus **Vanilla GANs** are only good for Sampling/Generation.
- **Training is Hard**
 - • Non-Convergence
 - • Mode-Collapse

Training Problems

- Non-Convergence
- Mode-Collapse

GAN's formulation

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

- It is formulated as a **minimax game**, where:
 - The Discriminator is trying to maximize its reward $V(D, G)$
 - The Generator is trying to minimize Discriminator's reward (or maximize its loss)

$$V(D, G) = \boxed{\mathbb{E}_{x \sim p(x)} [\log D(x)]} + \boxed{\mathbb{E}_{z \sim q(z)} [\log(1 - D(G(z)))]}$$

- The Nash equilibrium of this particular game is achieved at:

$$\left\{ \begin{array}{l} \bullet P_{data}(x) = P_{gen}(x) \quad \forall x \\ \bullet D(x) = \frac{1}{2} \quad \forall x \end{array} \right.$$

- Deep Learning models (in general) involve a single player

- The player tries to maximize its reward (minimize its loss).
- Use SGD (with Backpropagation) to find the optimal parameters.
- SGD has convergence guarantees (under certain conditions).
- **Problem:** With non-convexity, we might converge to local optima.

$$\min_G L(G)$$

$$\min_{\theta} \max_{\phi} V(\theta, \phi)$$

- GANs instead involve two (or more) players

- Discriminator is trying to maximize its reward.
- Generator is trying to minimize Discriminator's reward.

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

- SGD was not designed to find the Nash equilibrium of a game.
- **Problem:** We might not converge to the Nash equilibrium at all.

Non-Convergence

$$\min_x \max_y V(x, y)$$

Let $V(x, y) = \underline{xy}$

- • State 1:

x > 0	y > 0	V > 0
-------	-------	-------
- State 2:

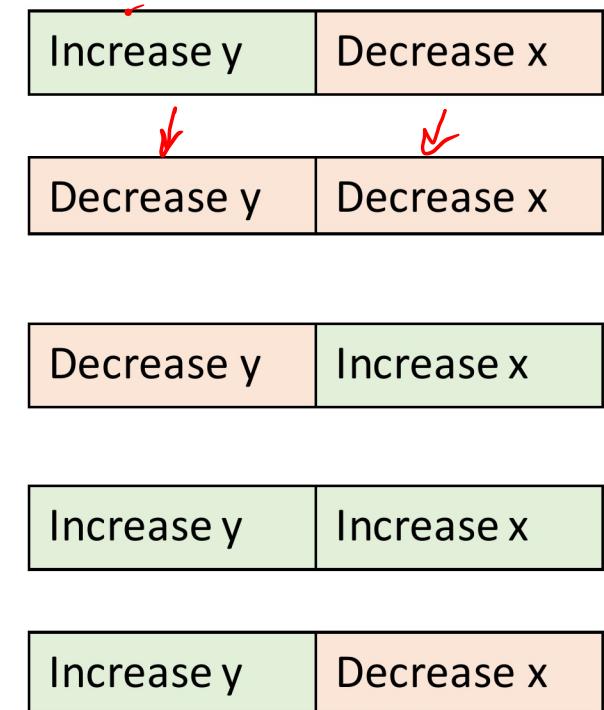
x < 0	y > 0	V < 0
-------	-------	-------
- State 3:

x < 0	y < 0	V > 0
-------	-------	-------
- State 4 :

x > 0	y < 0	V < 0
-------	-------	-------
- • State 5:

x > 0	y > 0	V > 0
-------	-------	-------

 == State 1



Vanishing gradient strikes back again...

$$\sigma(a) \rightarrow \sigma'(a) = \sigma(a)(1 - \sigma(a))$$

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

$$V(D, G) = \mathbb{E}_{x \sim p(x)} [\log D(x)] + \boxed{\mathbb{E}_{z \sim q(z)} [\log(1 - D(G(z)))]}$$

$$\nabla_{\theta_G} V(D, G) = \nabla_{\theta_G} \mathbb{E}_{z \sim q(z)} \left[\log \left(1 - \underbrace{D(G(z))}_{\text{red}} \right) \right]$$

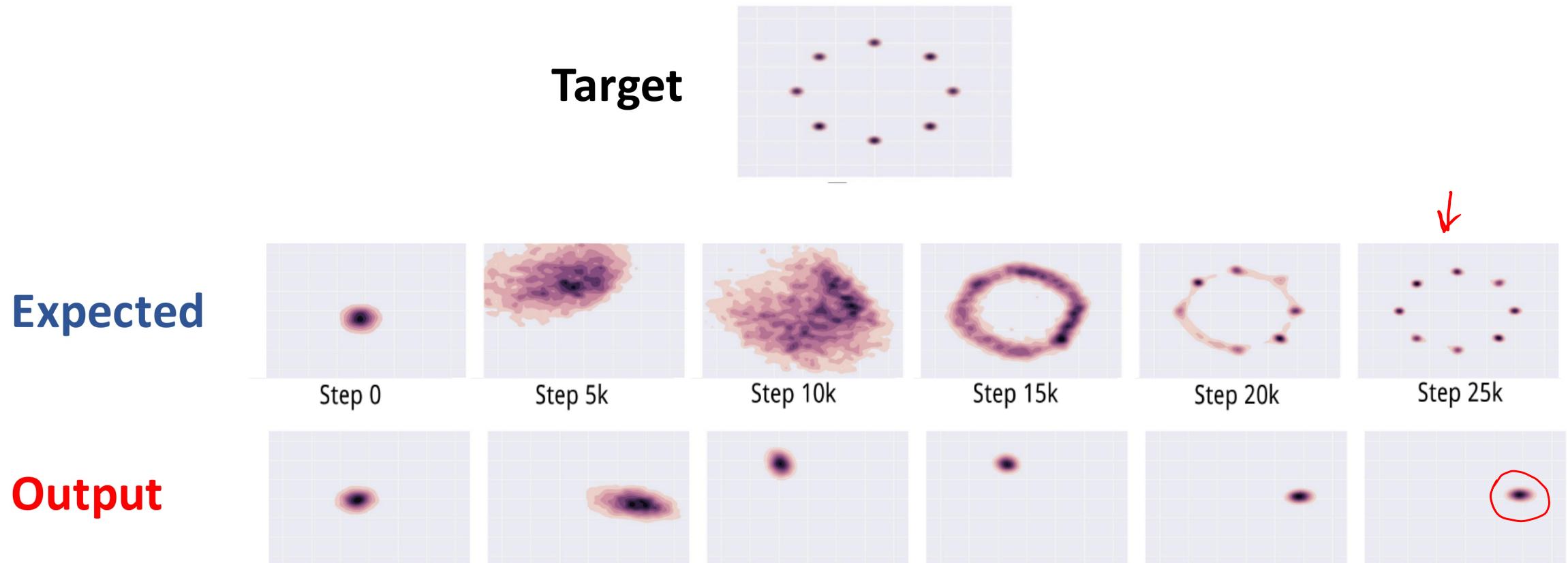
- $\nabla_a \log(1 - \sigma(a)) = \frac{-\nabla_a \sigma(a)}{1 - \sigma(a)} = \frac{-\sigma(a)(1 - \sigma(a))}{1 - \sigma(a)} = \underline{-\sigma(a)} = \underline{-D(G(z))}$
- Gradient goes to 0 if D is confident, i.e. $\boxed{D(G(z)) \rightarrow 0}$
- Minimize $\min \boxed{-\mathbb{E}_{z \sim q(z)} [\log D(G(z))]}$ for **Generator** instead (keep Discriminator as it is)

Problems with GANs

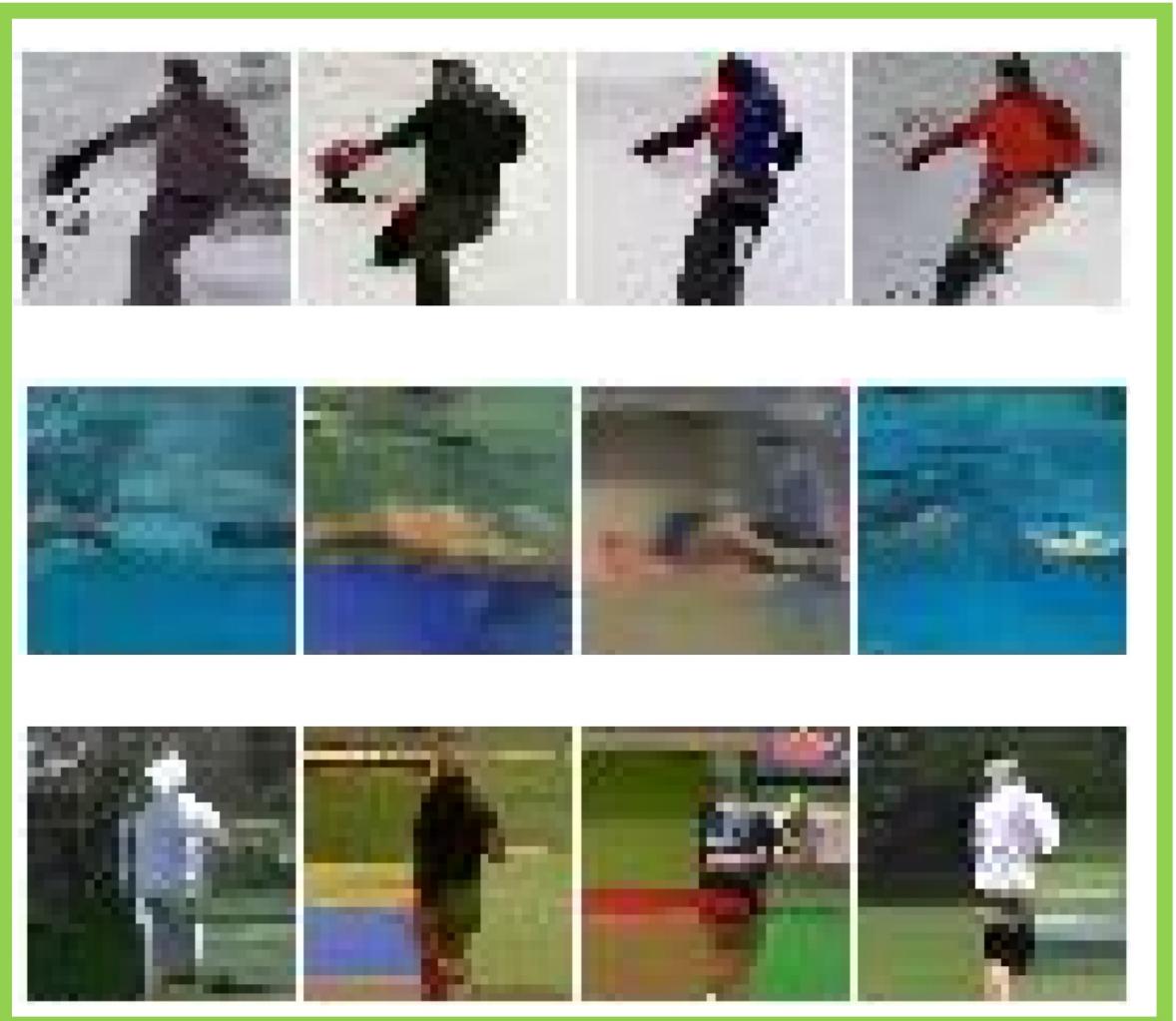
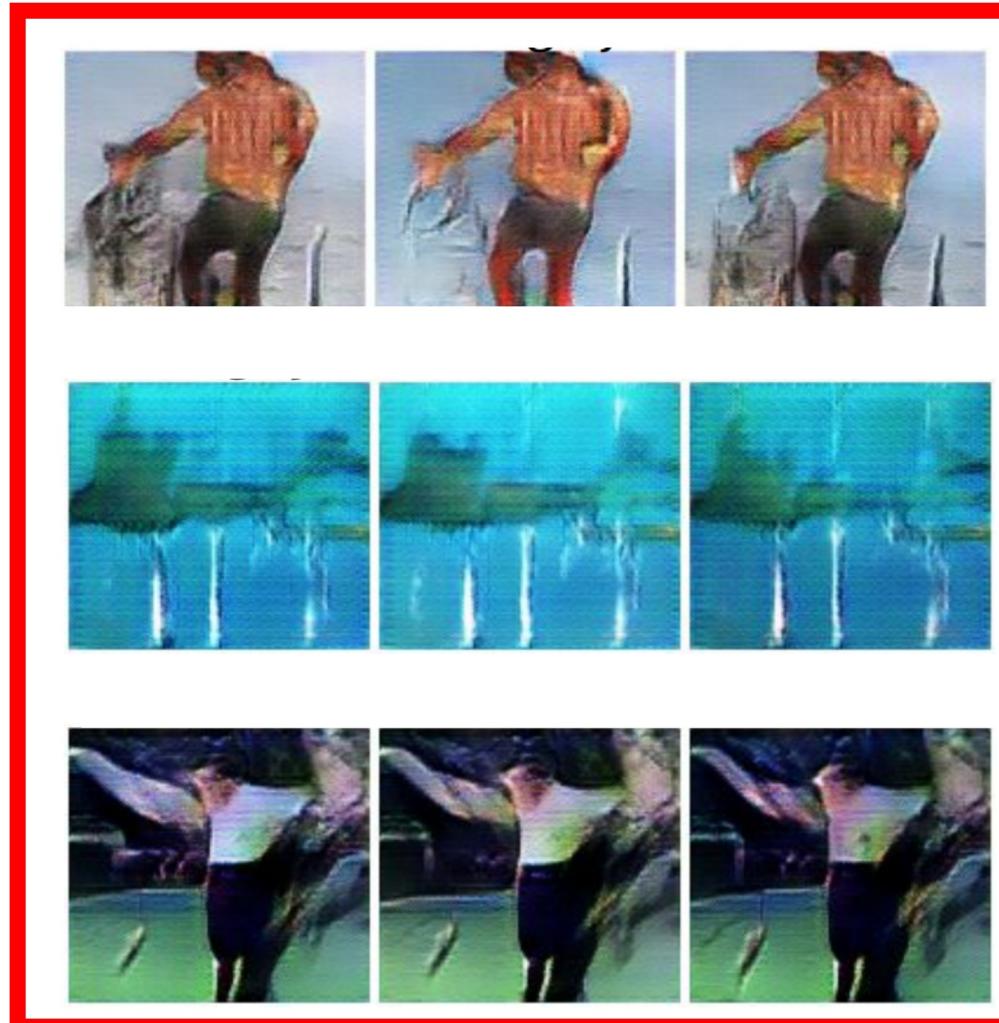
- Non-Convergence
- **Mode-Collapse**

Mode-Collapse

- Generator fails to output diverse samples



Some real examples



Wasserstein GAN: beyond f -divergences

The f -divergence is defined as

$$D_f(p, q) = E_{\mathbf{x} \sim q} \left[f \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) \right]$$

$$\min_{\theta} KL(P \parallel q_{\theta})$$

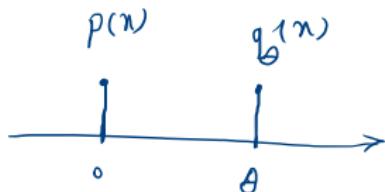
- The support of q has to cover the support of p . Otherwise discontinuity arises in f -divergences.

- Let $p(\mathbf{x}) = \begin{cases} 1, & \mathbf{x} = 0 \\ 0, & \mathbf{x} \neq 0 \end{cases}$, and $q_{\theta}(\mathbf{x}) = \begin{cases} 1, & \mathbf{x} = \theta \\ 0, & \mathbf{x} \neq \theta \end{cases}$.

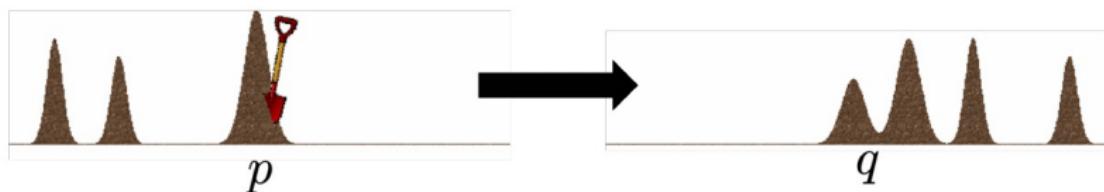
→ • $D_{KL}(p, q_{\theta}) = \begin{cases} 0, & \theta = 0 \\ \infty, & \theta \neq 0 \end{cases}$.

→ • $D_{JS}(p, q_{\theta}) = \begin{cases} 0, & \theta = 0 \\ \log 2, & \theta \neq 0 \end{cases}$.

- We need a “smoother” distance $D(p, q)$ that is defined when p and q have disjoint supports.



Wasserstein (Earth-Mover) distance



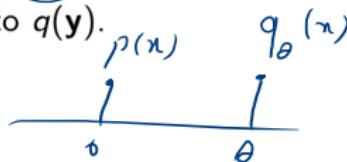
- Wasserstein distance

$$D_w(p, q) = \inf_{\gamma \in \Pi(p, q)} E_{(x, y) \sim \gamma} [\|x - y\|_1],$$

p q
↓ ↓
 $\delta(x, y)$

where $\Pi(p, q)$ contains all joint distributions of (x, y) where the marginal of x is $p(x) = \int \gamma(x, y) dy$, and the marginal of y is $q(y)$

- $\gamma(y | x)$: a probabilistic earth moving plan that warps $p(x)$ to $q(y)$.
- Let $p(x) = \begin{cases} 1, & x = 0 \\ 0, & x \neq 0 \end{cases}$, and $q_\theta(x) = \begin{cases} 1, & x = \theta \\ 0, & x \neq \theta \end{cases}$.
- $D_w(p, q_\theta) = |\theta|.$ ←



Wasserstein GAN (WGAN)

- Kantorovich-Rubinstein duality

$$D_w(p, q) = \sup_{\|f\|_L \leq 1} E_{x \sim p}[f(x)] - E_{x \sim q}[f(x)]$$

$\|f\|_L \leq 1$ means the Lipschitz constant of $f(x)$ is 1. Technically,

$$\forall x, y : |f(x) - f(y)| \leq \|x - y\|_1$$

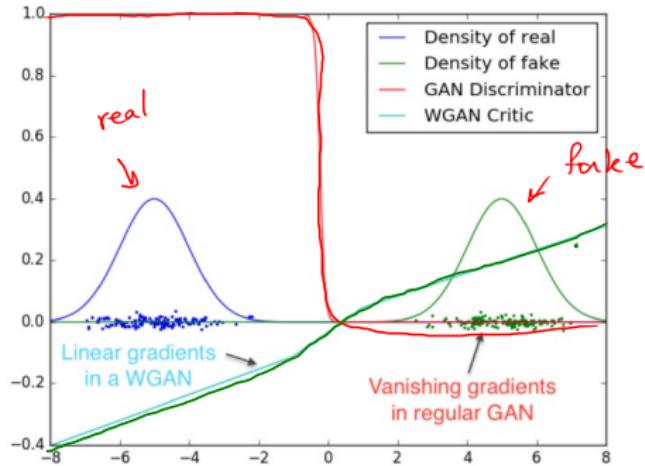
Intuitively, f cannot change too rapidly.

- Wasserstein GAN with discriminator $D_\phi(x)$ and generator $G_\theta(z)$:

$$\min_{\theta} \max_{\phi} E_{x \sim p_{\text{data}}}[D_\phi(x)] - E_{z \sim p(z)}[D_\phi(G_\theta(z))]$$

Lipschitzness of $D_\phi(x)$ is enforced through weight clipping or gradient penalty on $\nabla_x D_\phi(x)$.

Wasserstein GAN (WGAN)



- More stable training, and less mode collapse.

Inferring latent representations in GANs

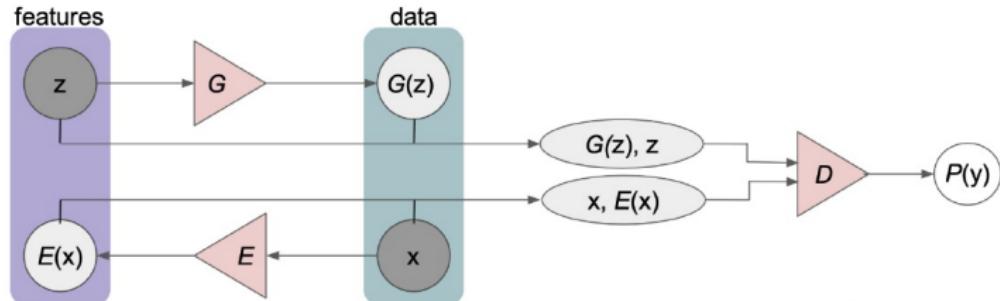
- The generator of a GAN is typically a directed, latent variable model with latent variables \mathbf{z} and observed variables \mathbf{x} . How can we infer the latent feature representations in a GAN?
- Unlike a normalizing flow model, the mapping $G : \mathbf{z} \mapsto \mathbf{x}$ need not be invertible
- Unlike a variational autoencoder, there is no inference network $q(\cdot)$ which can learn a variational posterior over latent variables
- **Solution 1:** For any point \mathbf{x} , use the activations of the prefinal layer of a discriminator as a feature representation
- Intuition: Similar to supervised deep neural networks, the discriminator would have learned useful representations for \mathbf{x} while distinguishing real and fake \mathbf{x}

$\mathbf{z} \rightarrow \mathbf{x}$

Inferring latent representations in GANs

- If we want to directly infer the latent variables \mathbf{z} of the generator, we need a different learning algorithm
- A regular GAN optimizes a two-sample test objective that compares samples of \mathbf{x} from the generator and the data distribution
- **Solution 2:** To infer latent representations, we will compare samples of \mathbf{x}, \mathbf{z} from the joint distributions of observed and latent variables as per the model and the data distribution
- For any \mathbf{x} generated via the model, we have access to \mathbf{z} (sampled from a simple prior $p(\mathbf{z})$)
- For any \mathbf{x} from the data distribution, the \mathbf{z} is however unobserved (latent). Need an encoder!

Bidirectional Generative Adversarial Networks (BiGAN)



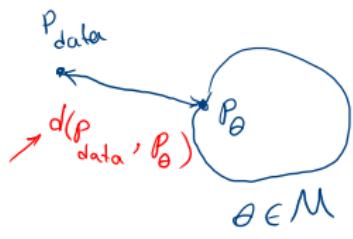
- In a BiGAN, we have an encoder network E in addition to the generator network G
- The encoder network only observes $\mathbf{x} \sim p_{\text{data}}(\mathbf{x})$ during training to learn a mapping $E : \mathbf{x} \mapsto \mathbf{z}$
- As before, the generator network only observes the samples from the prior $\mathbf{z} \sim p(\mathbf{z})$ during training to learn a mapping $G : \mathbf{z} \mapsto \mathbf{x}$

f-GAN

f-Divergence

$$D(p, q) = \mathbb{E}_q \left[f\left(\frac{p(x)}{q(x)}\right)\right]$$

- ① convex
- ② lower-semicontinuous
- ③ $f(1) = 0$



$$\hat{\theta} = \arg \max_{\theta} KL(p_{\text{data}} \parallel p_{\theta})$$

$$= \arg \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)]$$

NF

$$\approx \arg \max_{\theta} \frac{1}{n} \sum_{i=1}^n \log p_{\theta}(x_i)$$



$$\hat{\theta} = \arg \max_{\theta} \text{ELBO}$$

← VAE

$$\text{GAN: } d(P_{\text{data}}, P_{\theta}) = 2 \text{JSD}(P_{\text{data}}, P_{\theta}) - \log 4$$

Wasserstein GAN (WGAN)

$$W(p, q) = \inf_{\gamma \in \underline{\Pi}(p, q)} E \left[\|x - y\|_1 \right]$$

$$\begin{matrix} p(x) \\ q(y) \end{matrix}$$

$$\boxed{\begin{aligned} q(y) &= \int \gamma(x, y) dx \\ p(x) &= \int \gamma(x, y) dy \end{aligned}}$$

$$W(p, q) = \sup_{\|f\|_L \leq 1} E_p[f(x)] - E_q[f(x)]$$

$$W(P_{\text{data}}, P_\theta) \simeq \sup_f \frac{1}{n} \sum_{i=1}^n f(x_i) - \frac{1}{n} \sum_{i=1}^n f(z_i)$$

$x_i \sim P_{\text{data}}$ $z_i \sim P_\theta$

$$\min_\theta W(P_{\text{data}}, P_\theta)$$

$$W(P_{\text{data}}, P_\theta) = \max_\phi \frac{1}{n} \sum_{i=1}^n D_\phi(x_i) - \frac{1}{n} \sum_{i=1}^n D_\phi(G_\theta(z_i))$$

$$\boxed{\min_\theta \max_\phi \frac{1}{n} \sum_{i=1}^n D_\phi(x_i) - \frac{1}{n} \sum_{i=1}^n D_\phi(G_\theta(z_i)) + \lambda \underbrace{\|\nabla_x D_\phi(x)\|}$$

$\mathbb{R} \rightarrow D_\phi \rightarrow \text{real/false}$

$$\frac{\partial D}{\partial z} = \frac{\partial D}{\partial \phi} \frac{\partial \phi}{\partial z}$$

$$0 \leq \phi(z) \leq 1$$

$$\phi'(z) = \phi(z)(1-\phi(z))$$

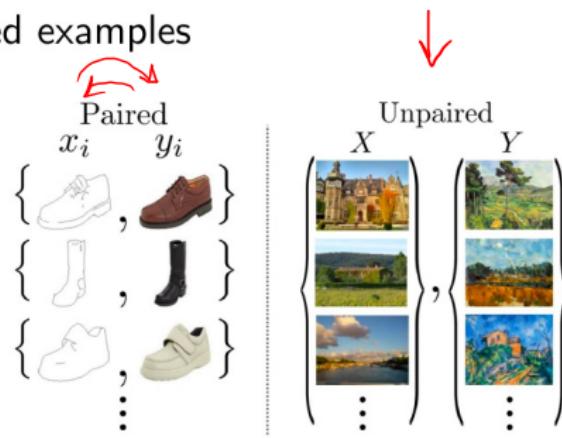
L-Lipschitz :

$$\|f(x) - f(x')\| \leq L \|x - x'\|$$

$$\frac{\|f(x) - f(x + \Delta x)\|}{\|\Delta x\|} \leq L$$

Translating across domains

- Image-to-image translation: We are given images from two domains, \mathcal{X} and \mathcal{Y}
- Paired vs. unpaired examples

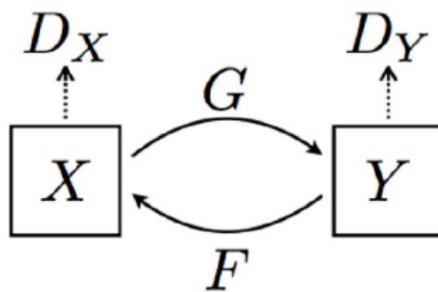


Source: Zhu et al., 2016

- Paired examples can be expensive to obtain. Can we translate from $\mathcal{X} \leftrightarrow \mathcal{Y}$ in an unsupervised manner?

CycleGAN: Adversarial training across two domains

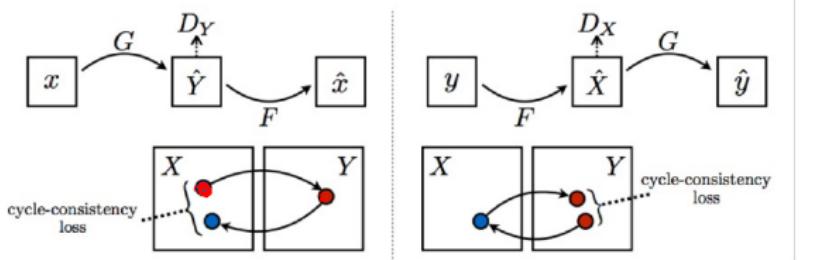
- To match the two distributions, we learn two parameterized conditional generative models $G : \mathcal{X} \mapsto \mathcal{Y}$ and $F : \mathcal{Y} \mapsto \mathcal{X}$
- G maps an element of \mathcal{X} to an element of \mathcal{Y} . A discriminator D_Y compares the observed dataset Y and the generated samples $\hat{Y} = G(X)$
- Similarly, F maps an element of \mathcal{Y} to an element of \mathcal{X} . A discriminator D_X compares the observed dataset X and the generated samples $\hat{X} = F(Y)$



Source: Zhu et al., 2016

CycleGAN: Cycle consistency across domains

- **Cycle consistency:** If we can go from X to \hat{Y} via G , then it should also be possible to go from \hat{Y} back to X via F
 - $F(G(X)) \approx X$
 - Similarly, vice versa: $G(F(Y)) \approx Y$

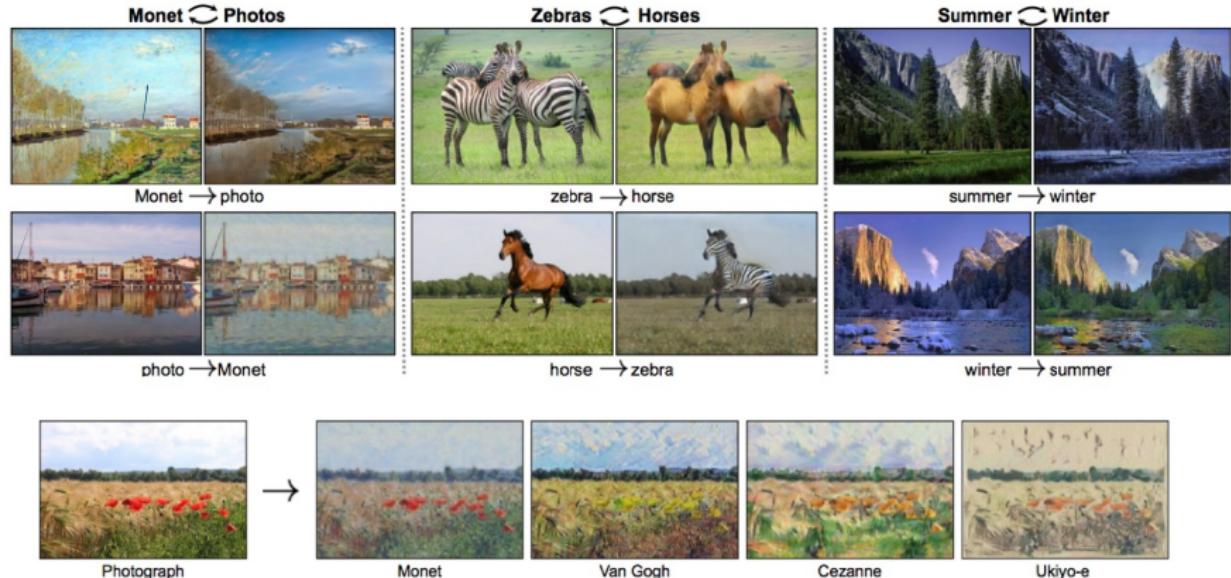


Source: Zhu et al., 2016

- Overall loss function

$$\begin{aligned} & \min_{F, G, D_X, D_Y} \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \mathcal{L}_{\text{GAN}}(F, D_X, X, Y) \\ & + \lambda \underbrace{\left(E_X[\|F(G(X)) - X\|_1] + E_Y[\|G(F(Y)) - Y\|_1] \right)}_{\text{cycle consistency}} \end{aligned}$$

CycleGAN in practice



Source: Zhu et al., 2016

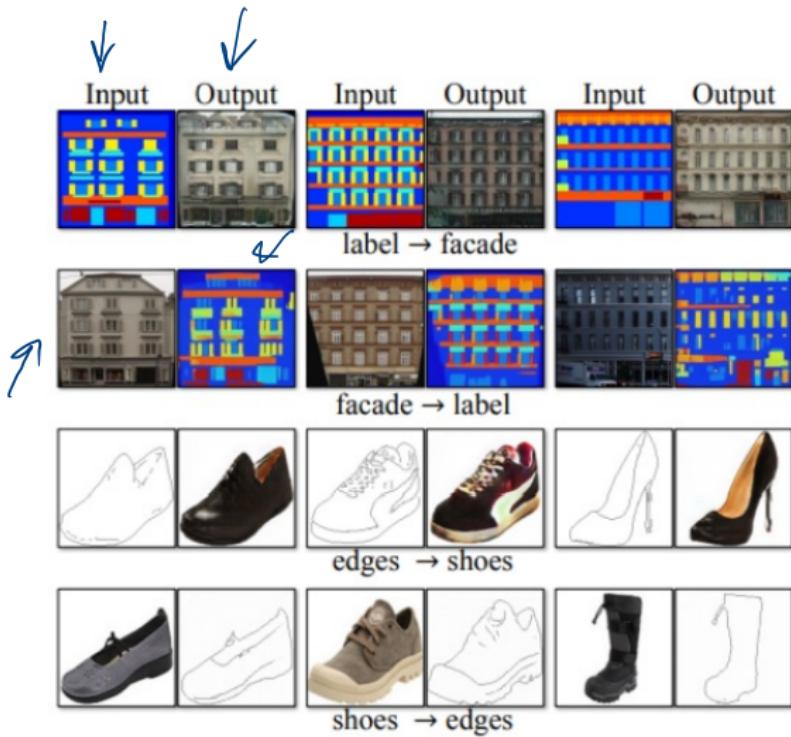


Figure 8: Example results of CycleGAN on paired datasets used in “pix2pix” [22] such as architectural labels↔photos and edges↔shoes.

AlignFlow (Grover et al.)

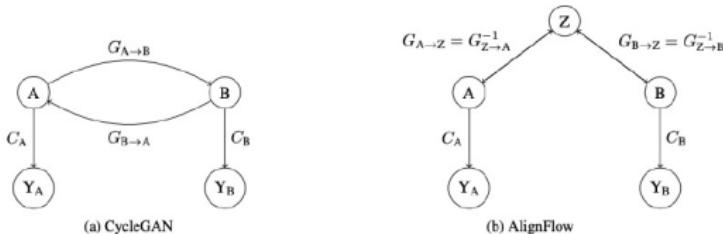


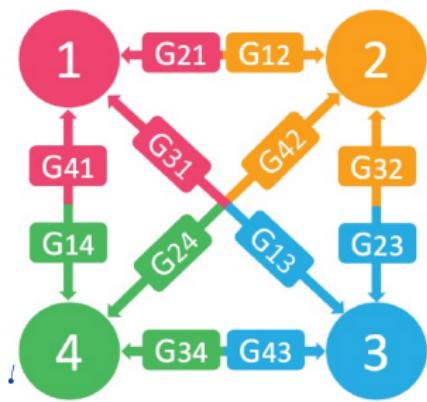
Figure 1: CycleGAN v.s. AlignFlow for unpaired cross-domain translation. Unlike CycleGAN, AlignFlow specifies a single invertible mapping $G_{A \rightarrow Z} \circ G_{B \rightarrow Z}^{-1}$ that is exactly cycle-consistent, represents a shared latent space Z between the two domains, and can be trained via both adversarial training and exact maximum likelihood estimation. Double-headed arrows denote invertible mappings. Y_A and Y_B are random variables denoting the output of the critics used for adversarial training.

- What if G is a flow model?
 - No need to parameterize F separately! $F = G^{-1}$
 - Can train via MLE and/or adversarial learning!
 - Exactly cycle-consistent
- {
- $F(G(X)) = X$
 - $G(F(Y)) = Y$

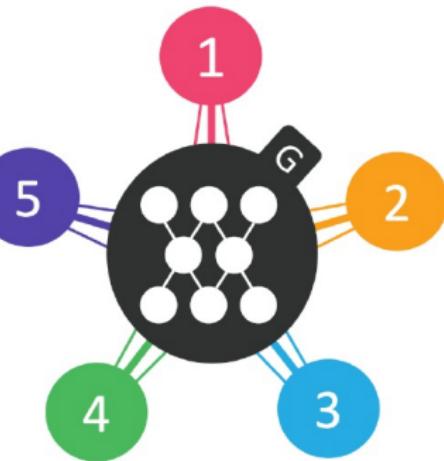
StarGAN (Choi et al.)

- What if there are multiple domains?

(a) Cross-domain models

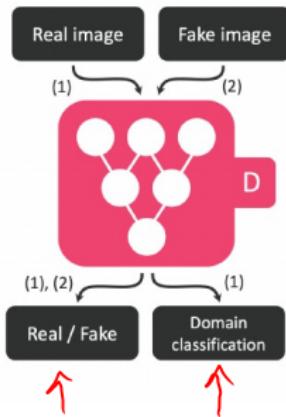


(b) StarGAN

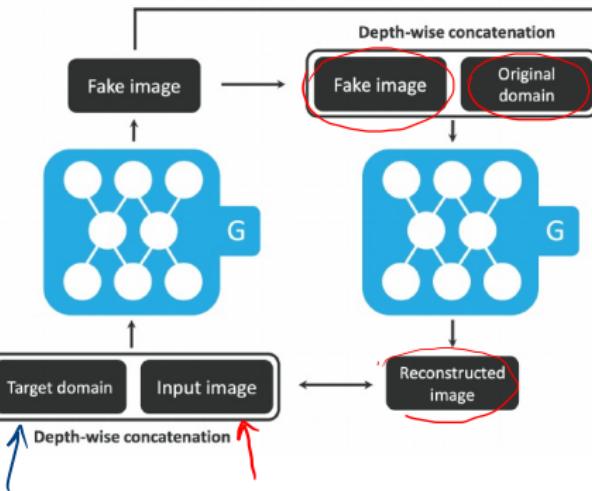


StarGAN (Choi et al.)

(a) Training the discriminator

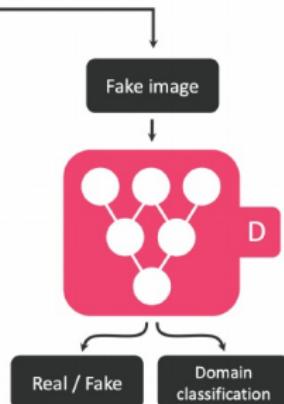


(b) Original-to-target domain

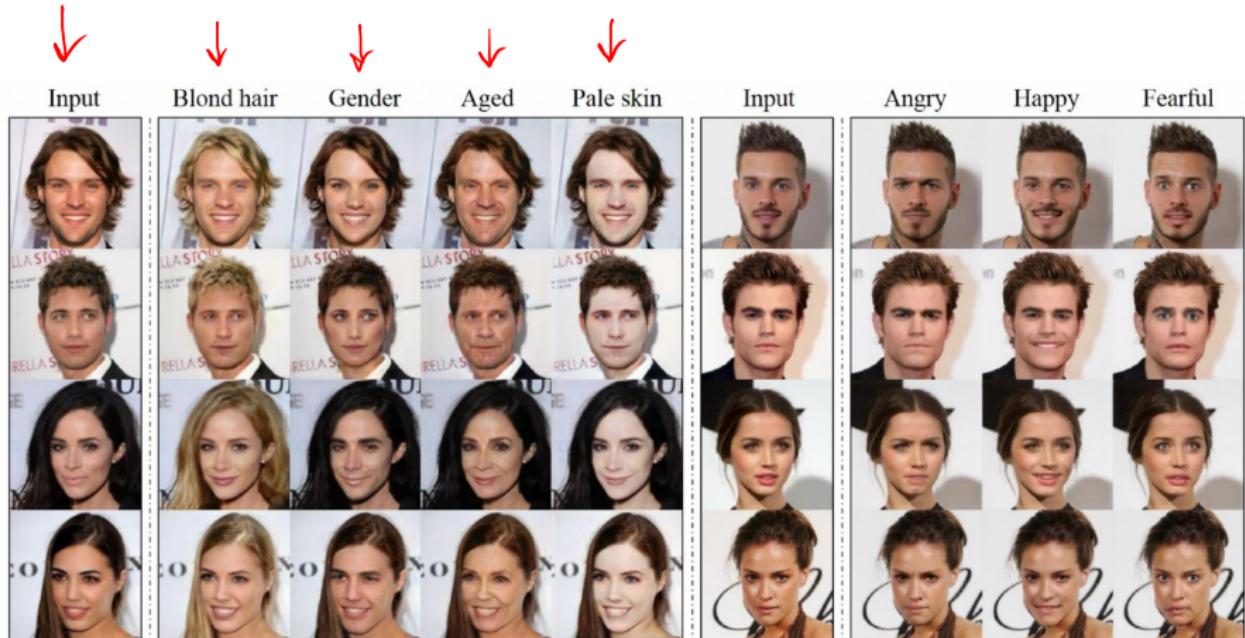


(c) Target-to-original domain

(d) Fooling the discriminator



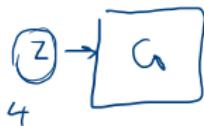
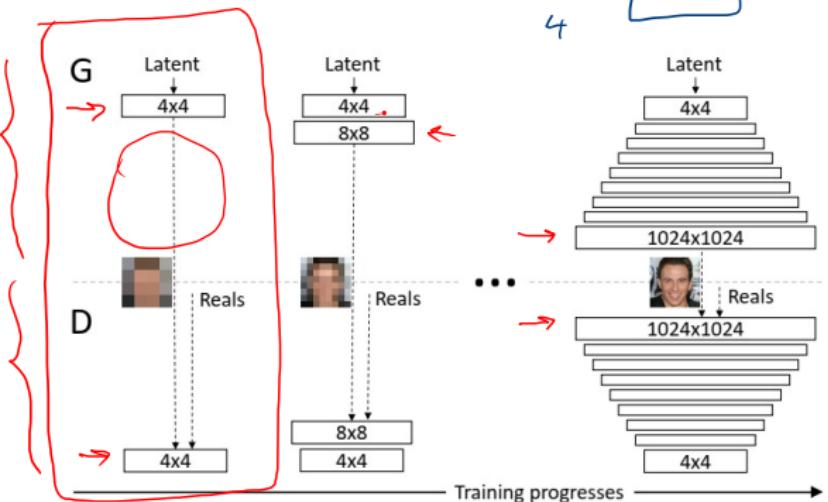
StarGAN (Choi et al.)



Summary of Generative Adversarial Networks

- Key observation: Samples and likelihoods are not correlated in practice
- Two-sample test objectives allow for learning generative models only via samples (likelihood-free)
- Wide range of two-sample test objectives covering f -divergences and Wasserstein distances (and more)
- Latent representations can be inferred via BiGAN
- Cycle-consistent domain translations via CycleGAN, AlignFlow and StarGAN.

Progressive GAN



style GAN 1, 2, 3



Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout the process. Here $N \times N$ refers to convolutional layers operating on $N \times N$ spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. One the right we show six example images generated using progressive growing at 1024×1024 .

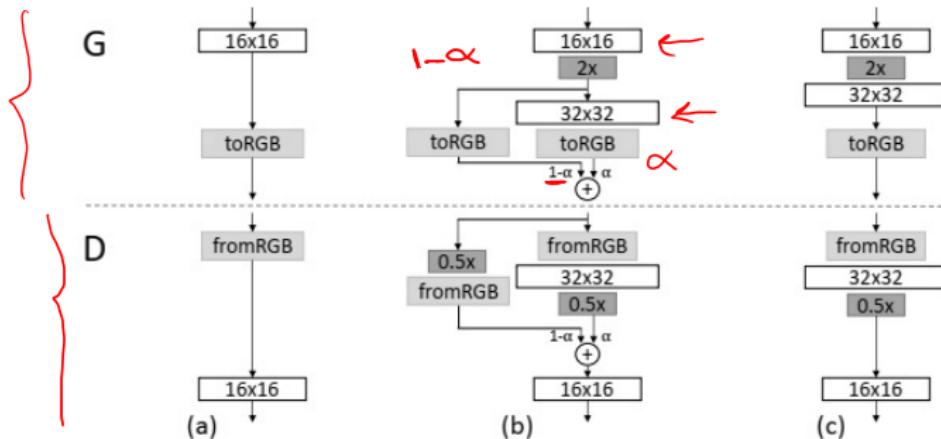


Figure 2: When doubling the resolution of the generator (G) and discriminator (D) we fade in the new layers smoothly. This example illustrates the transition from 16×16 images (a) to 32×32 images (c). During the transition (b) we treat the layers that operate on the higher resolution like a residual block, whose weight α increases linearly from 0 to 1. Here $2\times$ and $0.5\times$ refer to doubling and halving the image resolution using nearest neighbor filtering and average pooling, respectively. The $\boxed{\text{toRGB}}$ represents a layer that projects feature vectors to RGB colors and $\boxed{\text{fromRGB}}$ does the reverse; both use 1×1 convolutions. When training the discriminator, we feed in real images that are downsampled to match the current resolution of the network. During a resolution transition, we interpolate between two resolutions of the real images, similarly to how the generator output combines two resolutions.

StyleGAN (1, 2, 3)

