

# **TTIC 31230, Fundamentals of Deep Learning**

David McAllester, Winter 2019

## **Generative Adversarial Networks (GANs)**

# Generative Adversarial Nets

## Goodfellow et al., June 2014

Let  $\text{Pop} \uplus P_\Phi$  be the distribution defined flipping an unbiased coin and, if heads, returning  $(1, y)$  with  $y \sim \text{Pop}$  and, if tails, returning  $(-1, y)$  with  $y \sim P_\Phi$ .

$$\Phi^* = \operatorname{argmax}_\Phi \min_\Psi E_{(i,y) \sim (\text{Pop} \uplus P_\Phi)} - \ln P_\Psi(i|y)$$

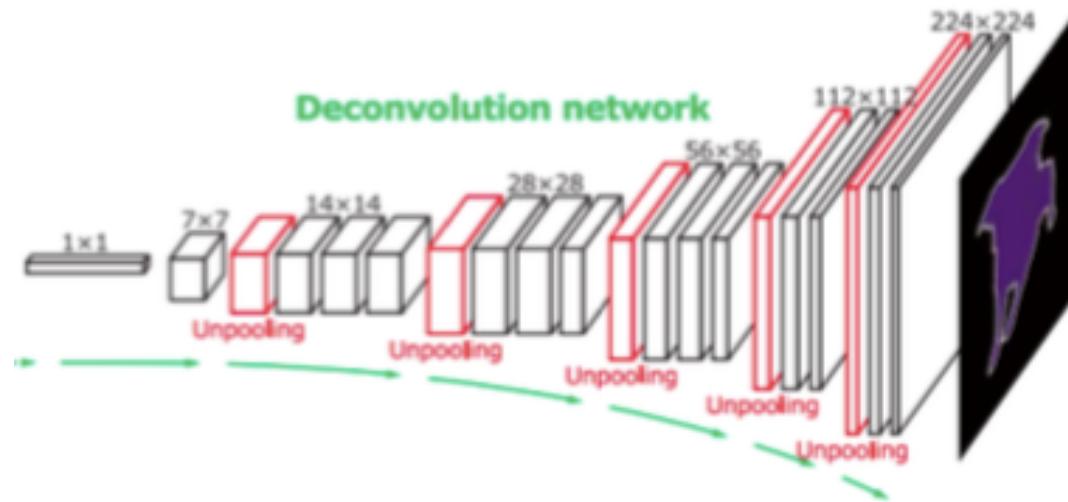
/ \\\  
Generator      Discriminator

Assuming Universality:  $P_{\Phi^*} = \text{Pop}$

# Sampling from the Generator

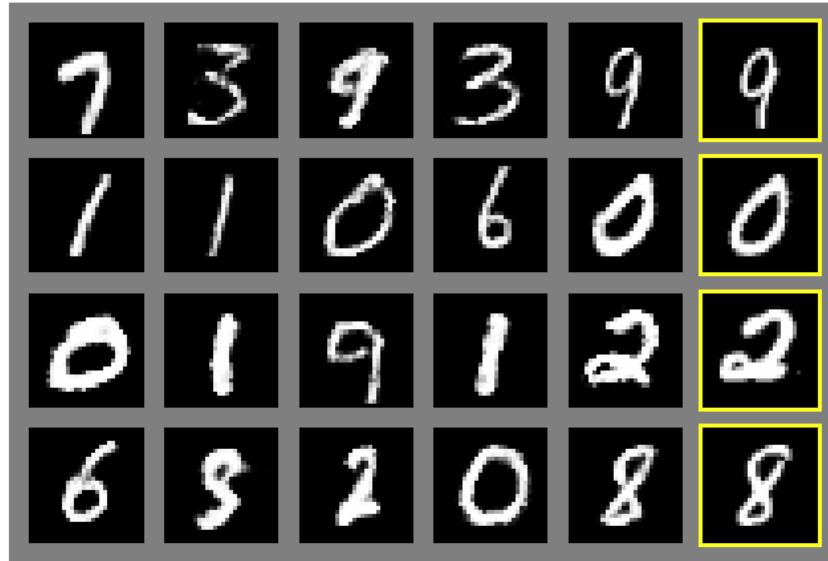
$$z \sim \mathcal{N}(0, I)$$

$$y \sim P_\Phi$$



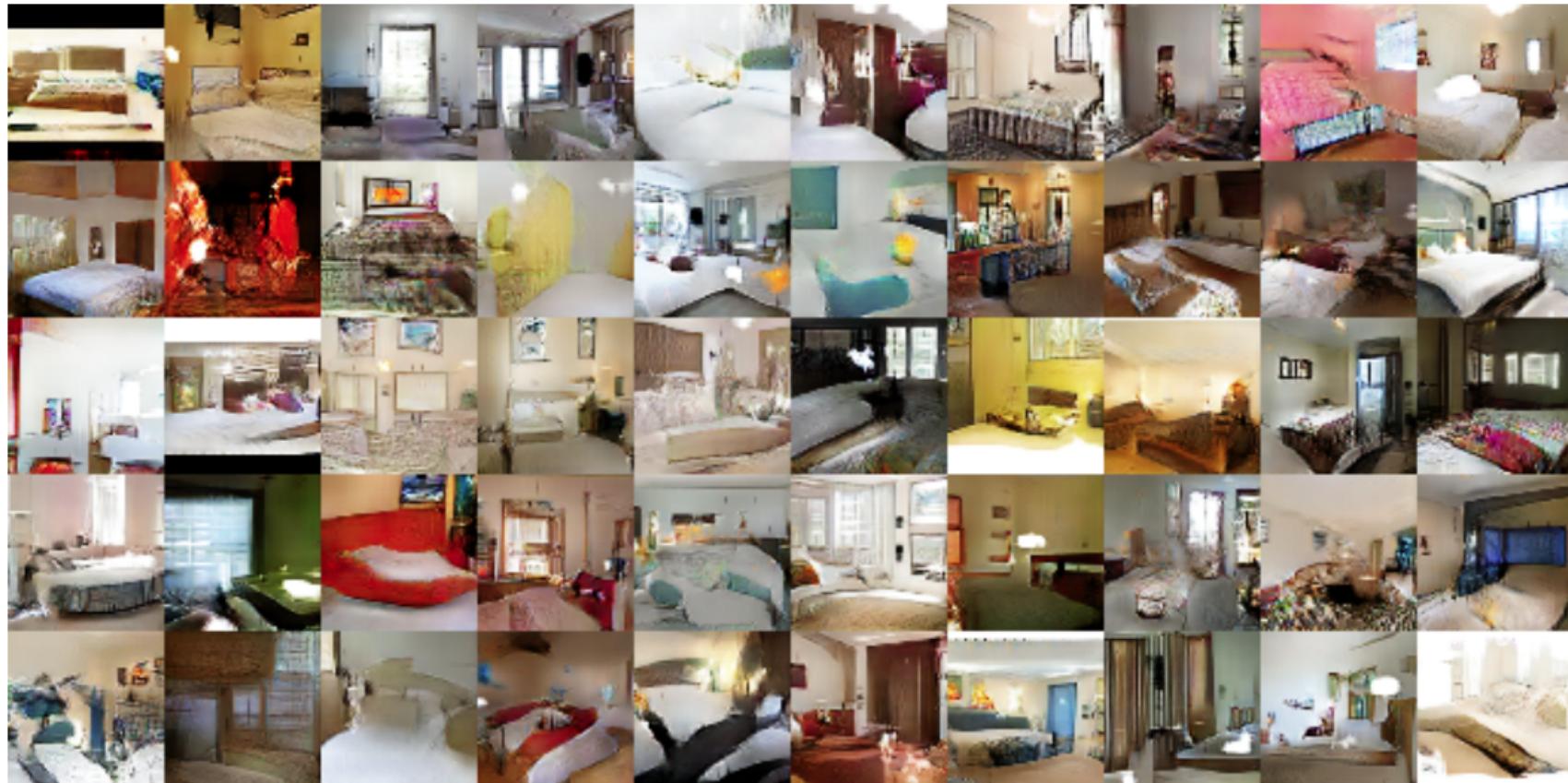
# Generative Adversarial Nets

## Goodfellow et al., June 2014



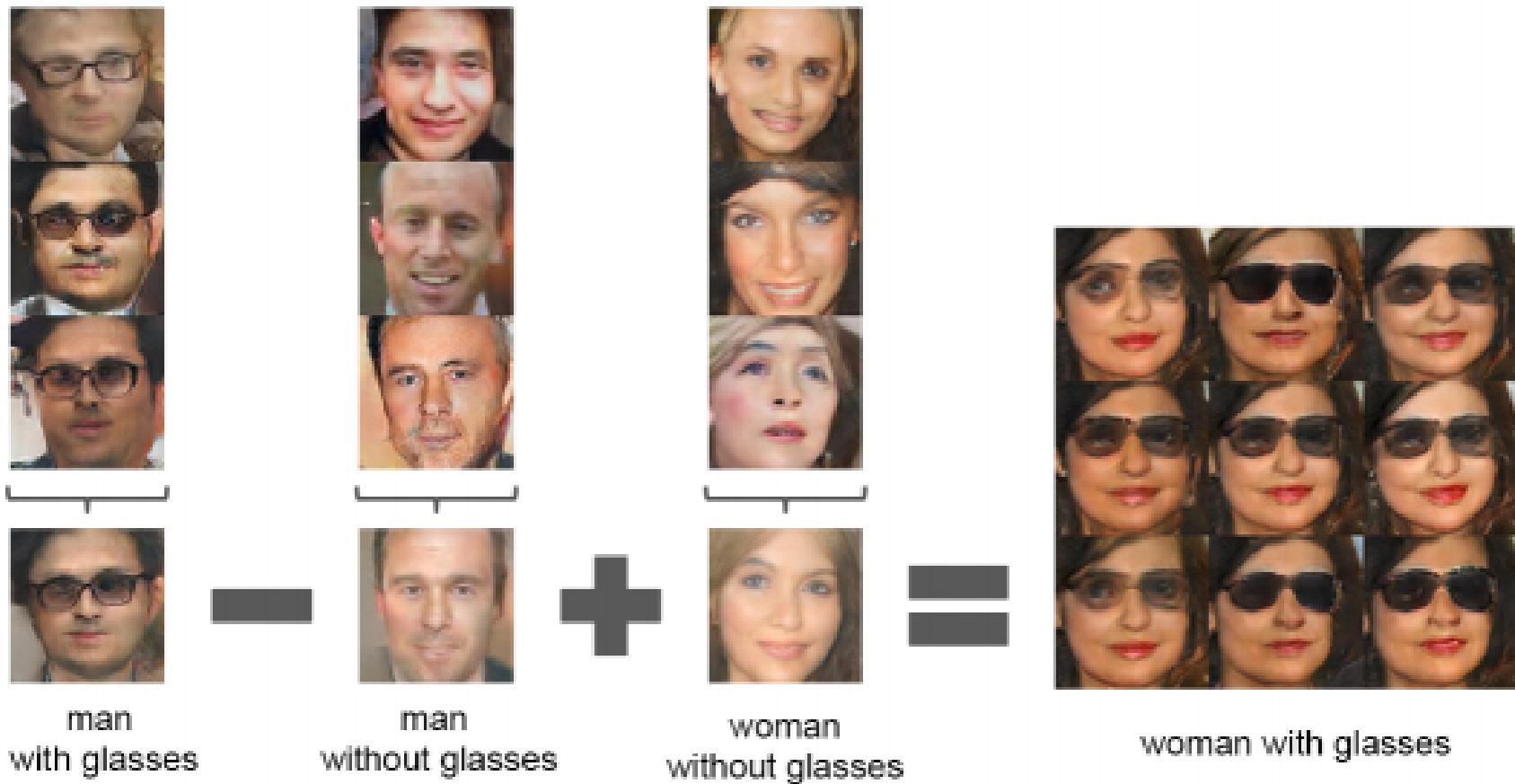
# Unsupervised Representation Learning ... (DC GANS)

## Radford et al., Nov. 2015



# Unsupervised Representation Learning ... (DC GANS)

Radford et al., Nov. 2015



# Interpolated Faces

[Ayan Chakrabarti]



# Early GANs on ImageNet



## Conditional GANs

All distribution modeling methods apply to conditional distributions.

$$\Phi^* = \operatorname{argmax}_{\Phi} \min_{\Psi} E_{i,x,y \sim (\text{Pop} \uplus \text{Pop}(x)P_{\Phi}(y|x))} - \ln P_{\Psi}(i|x, y)$$

If  $x$  is never repeated we can replace  $P_{\Phi}(y|x)$  by a deterministic function  $\hat{y}_{\Phi}(x)$ .

$$\Phi^* = \operatorname{argmax}_{\Phi} \min_{\Psi} E_{i,x,y \sim (\text{Pop} \uplus (\text{Pop}(x);\hat{y}(x)))} - \ln P_{\Psi}(i|x, y)$$

## Discrimination Loss can Augment Distortion Loss

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{x,y \sim \text{Pop}} \text{Dist}(y, \hat{y}_\Phi(x))$$

This is the fundamental equation if  $\text{Dist}(y, \hat{y}_\Phi(x)) = -\ln P_\Phi(y|x)$ .

$$\Phi^* = \operatorname{argmin}_{\Phi} \mathcal{L}_{\text{Dist}}(\Phi) + \lambda \mathcal{L}_{\text{DisC}}(\Phi)$$

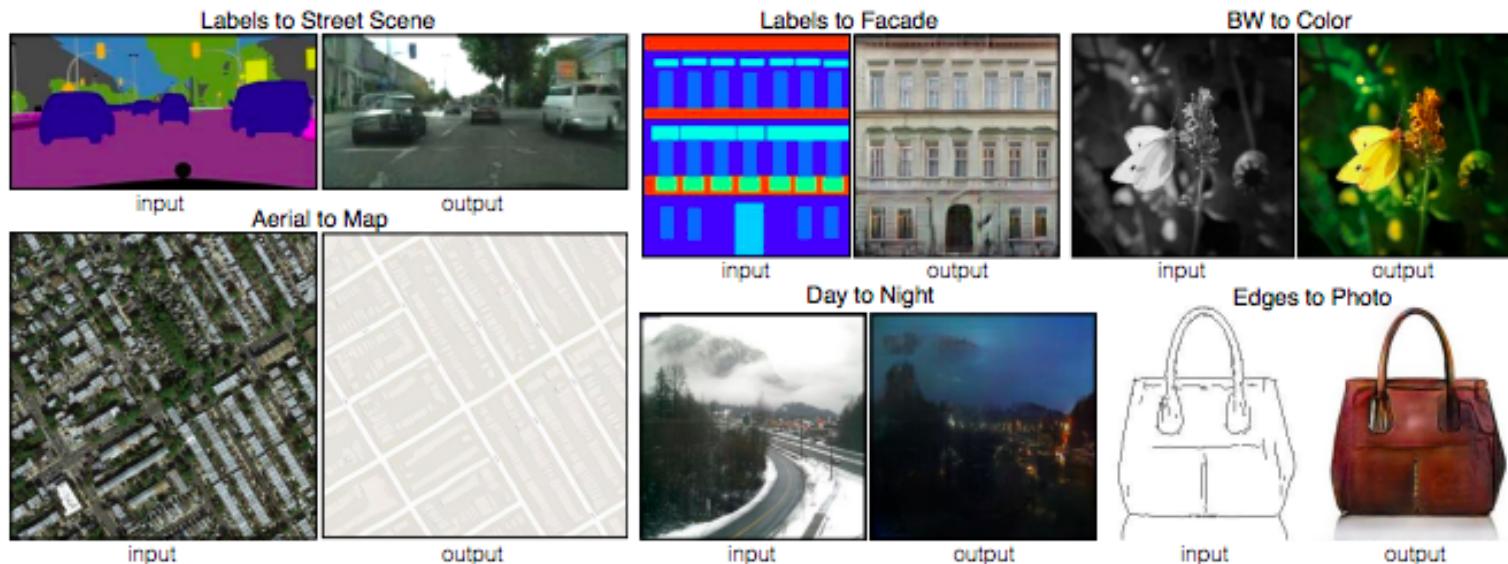
$$\mathcal{L}_{\text{Dist}}(\Phi) = E_{x,y \sim \text{Pop}} \text{Dist}(y, \hat{y}_\Phi(x))$$

$$\mathcal{L}_{\text{DisC}}(\Phi) = \max_{\Psi} E_{i,x,y \sim (\text{Pop} \uplus (\text{Pop}(x); \hat{y}_\Phi(x)))} \ln P_\Psi(i|x, y)$$

# Image-to-Image Translation (Pix2Pix)

Isola et al., Nov. 2016

We assume a corpus of “image translation pairs” such as images paired with semantic segmentations.



# Image-to-Image Translation (Pix2Pix)

Isola et al., Nov. 2016



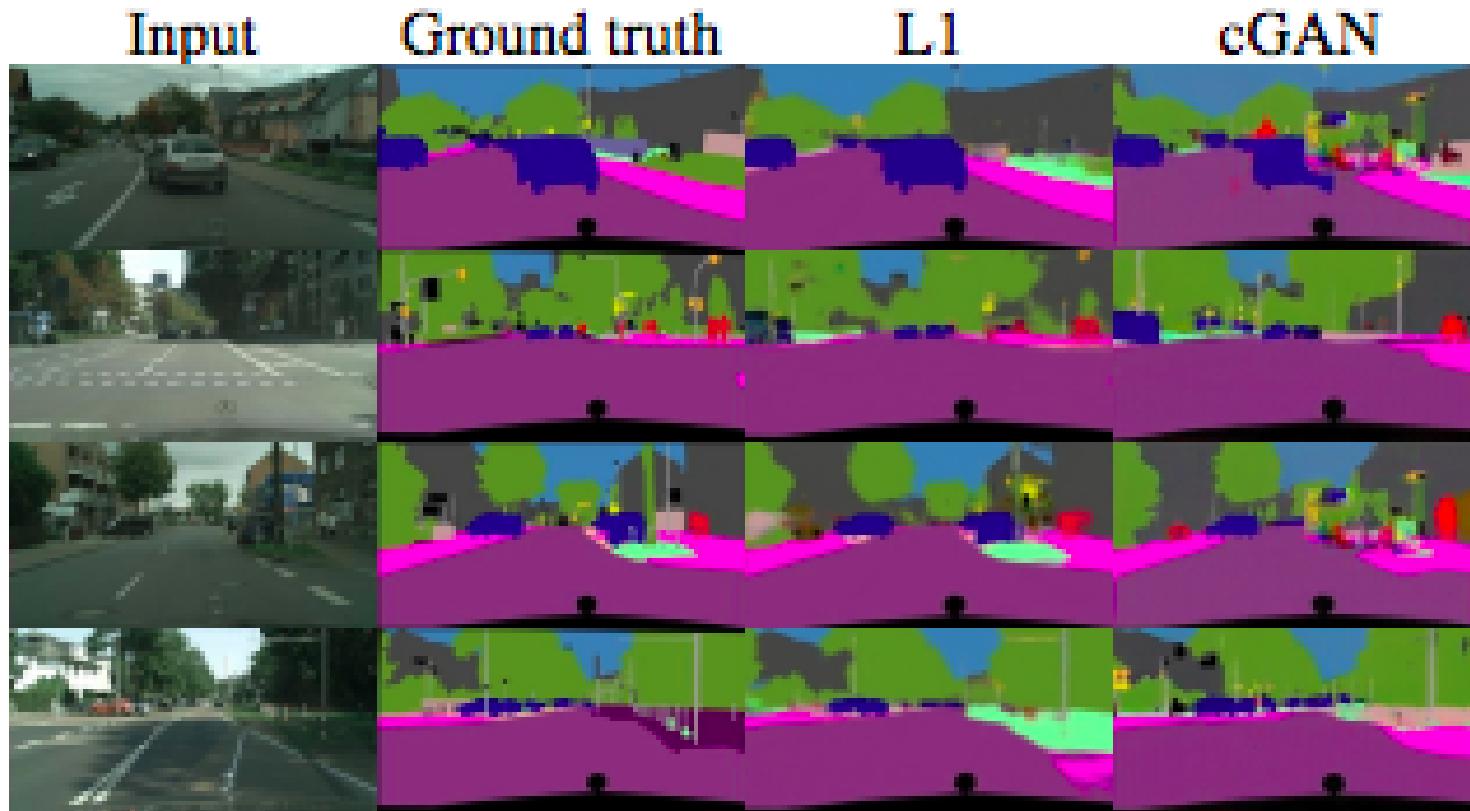
# Arial Photo to Map and Back



# Colorization



# Semantic Segmentation

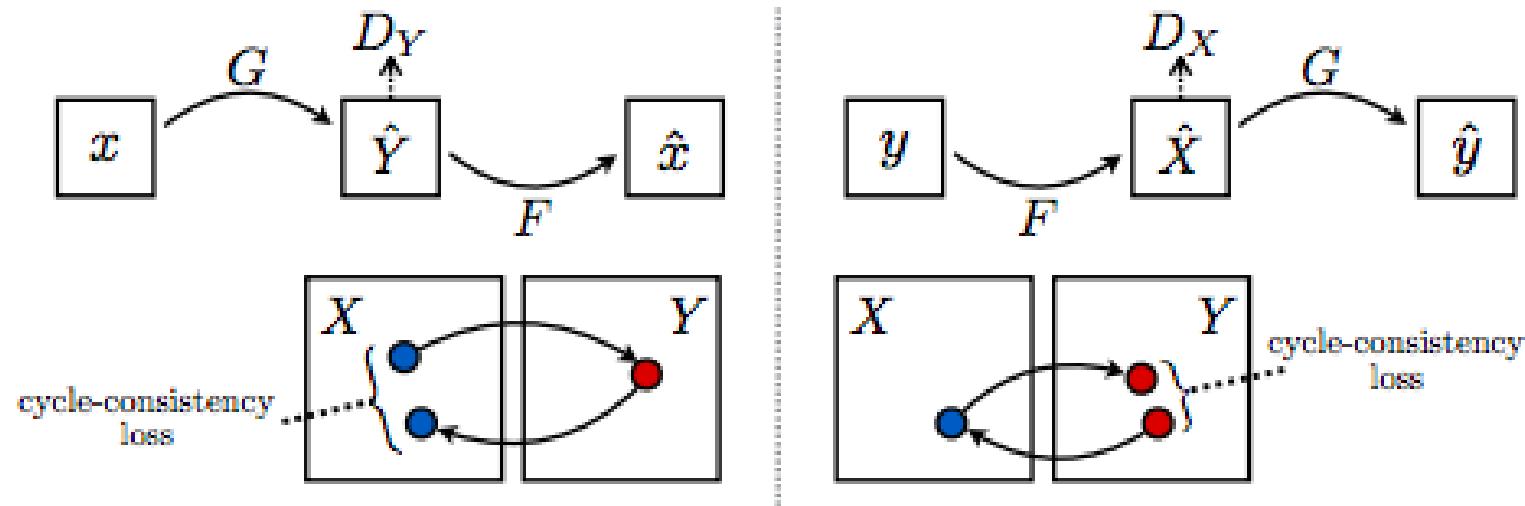


# Unpaired Image-to-Image Translation (Cycle GANs)

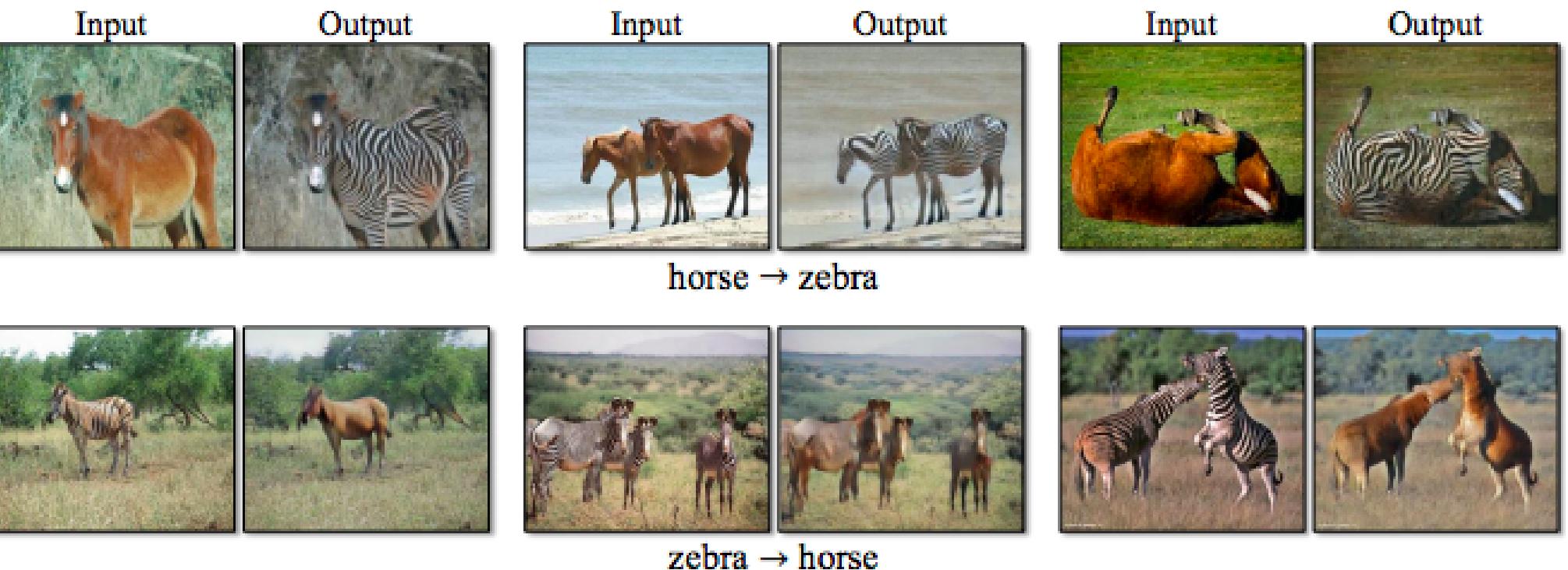
Zhu et al., March 2017

We have two corpora of images, say images of zebras and unrelated images of horses, or photographs and unrelated paintings by Monet.

We want to construct translations between the two classes.



# Cycle Gans



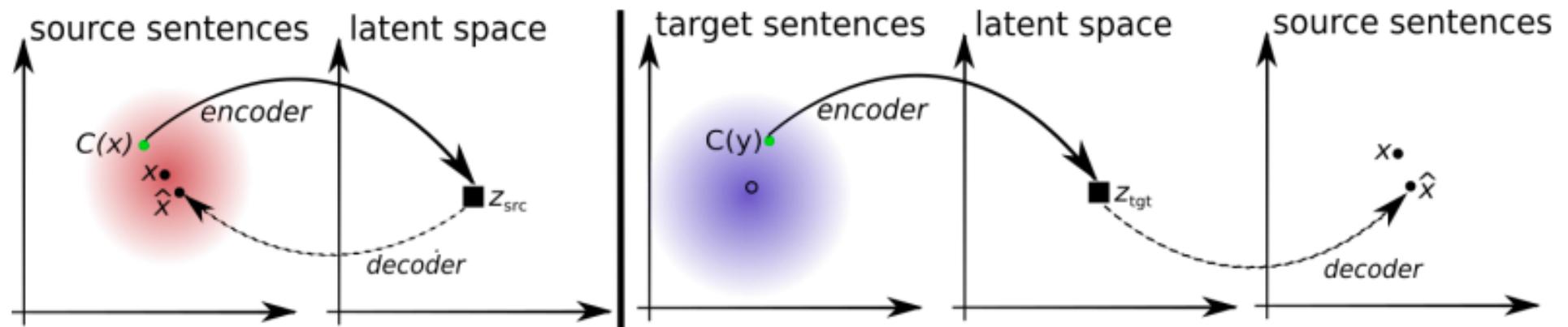
# Cycle Gans



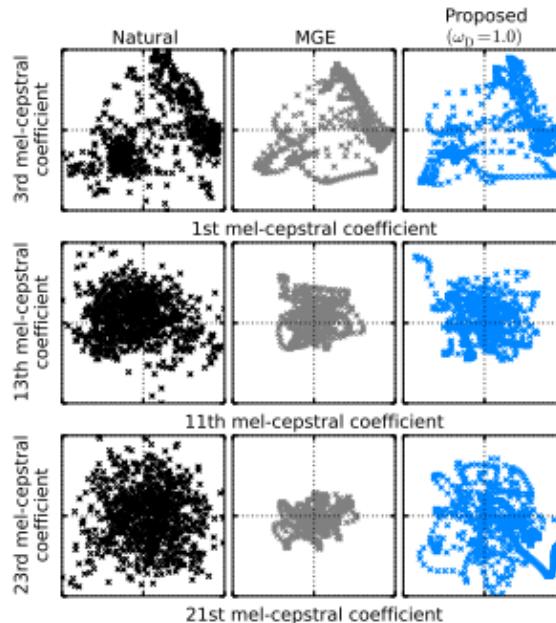
Horse → Zebra

# Unsupervised Machine Translation (UMT)

Lample et al, Oct. 2017, also Artetxe et al., Oct. 2017



# Text to Speech (Saito et al. Sept. 2017)

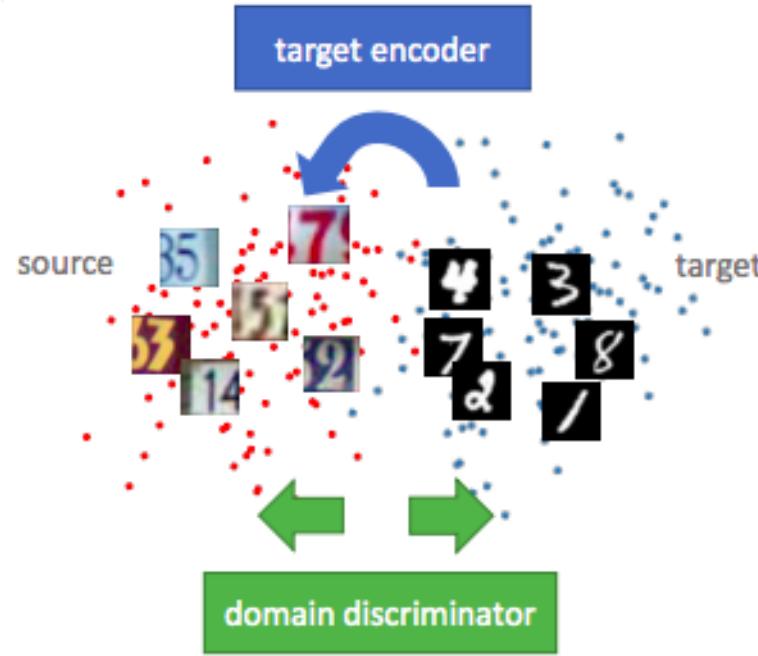


Minimum Generation Error (MGE) uses **perceptual distortion** — a distance between the feature vector of the generated sound wave and the feature vector of the original.

**Perceptual Naturalness** can be enforced by a discriminator.

# Adversarial Discriminative Domain Adaptation

Tzeng et al. Feb. 2017



A GAN is used to map target images to the source feature distribution.

# Progressive Growing of GANs

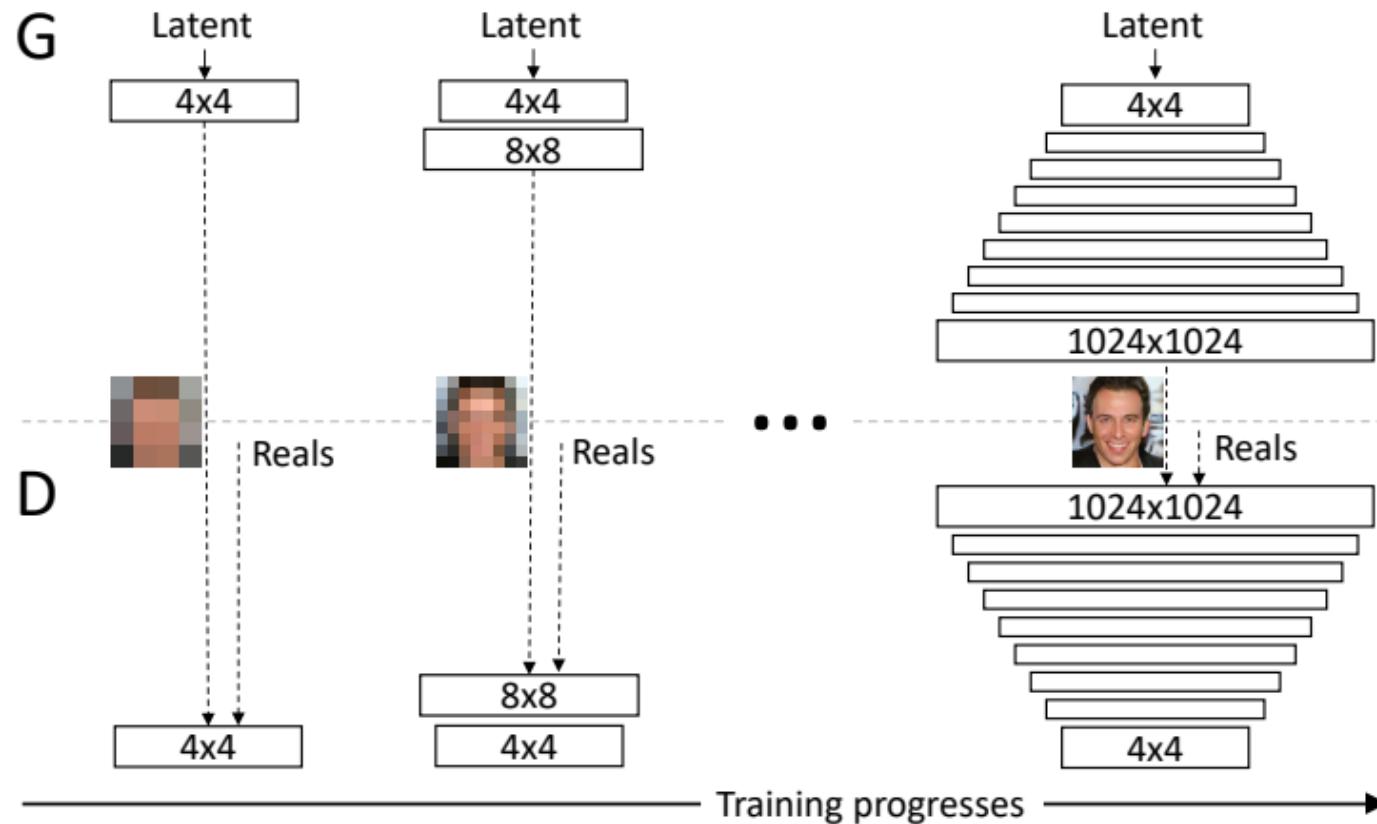
Karras et al., Oct. 2017



Figure 5:  $1024 \times 1024$  images generated using the CELEBA-HQ dataset. See Appendix F for a larger set of results, and the accompanying video for latent space interpolations.

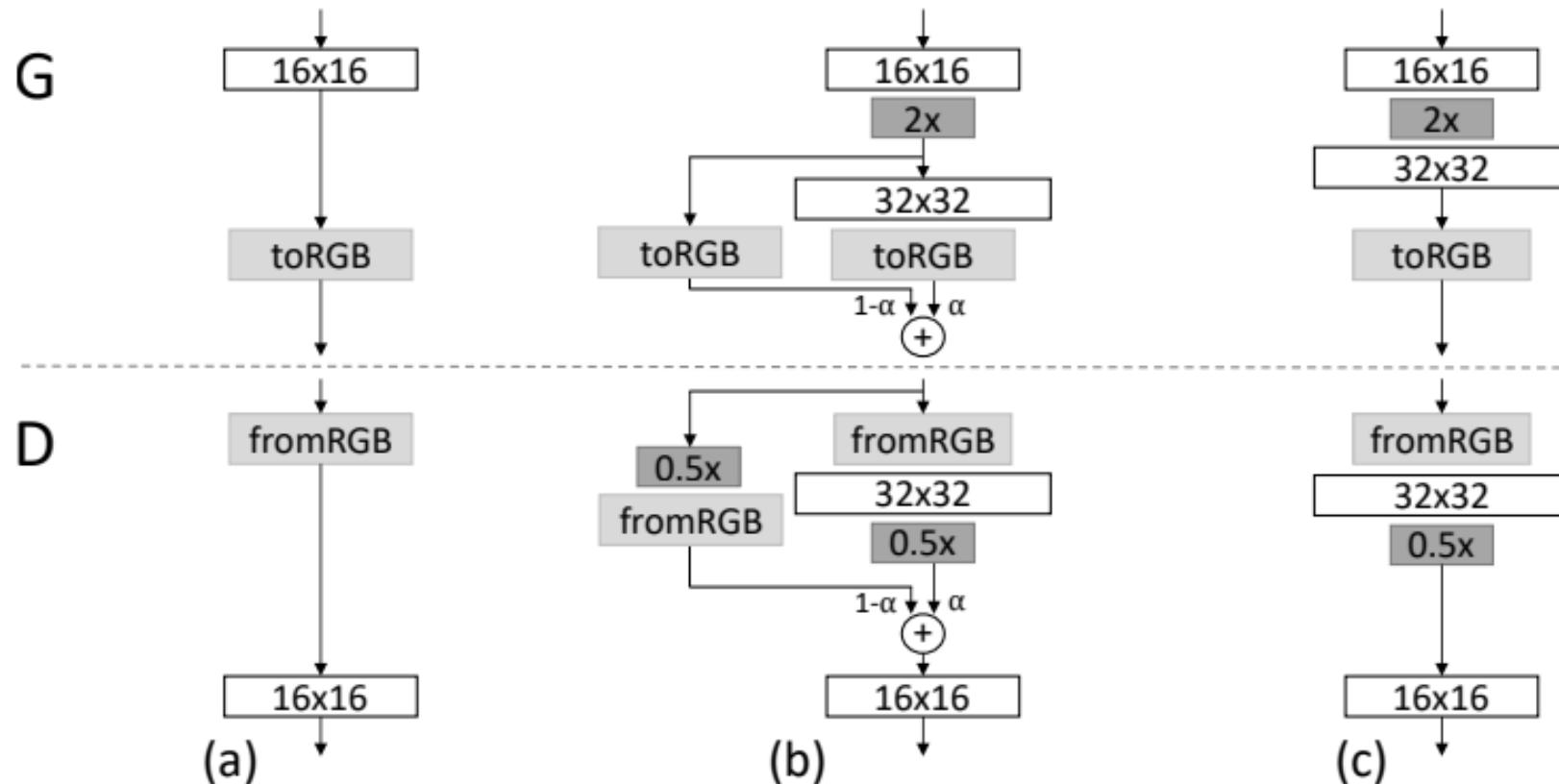
# Progressive Growing of GANs

Karras et al., Oct. 2017



# Progressive Growing of GANs

Karras et al., Oct. 2017



# Large Scale Gan Training

Brock et al., Sept. 2018



**Figure 1:** Class-conditional samples generated by our model.

This is a conditional GAN conditioned on the imagenet class label.

This generates 512 X 512 images without using progressive training.

# Issues

Jensen-Shannon Divergence

Vanishing Gradients

Unstable Training

Mode Collapse

Measuring Performance

## Discrimination: Jensen-Shannon Divergence

Assuming Universality:

$$\begin{aligned} & E_{(i,y) \sim (\text{Pop} \oplus Q)} - \ln P_{\Phi^*}(i|y) \\ = & E_{(i,y) \sim (\text{Pop} \oplus Q)} - \ln P(i|y) \\ = & \frac{1}{2} E_{y \sim \text{Pop}} - \ln \frac{\text{Pop}(y)}{\text{Pop}(y) + Q(y)} + \frac{1}{2} E_{y \sim Q} - \ln \frac{Q(y)}{\text{Pop}(y) + Q(y)} \\ = & (\ln 2) - \frac{1}{2} \left( KL \left( \text{Pop}, \frac{\text{Pop} + Q}{2} \right), KL \left( Q, \frac{\text{Pop} + Q}{2} \right) \right) \end{aligned}$$

## Discrimination: Jensen-Shannon Divergence

$$\begin{aligned} & E_{(i,y) \sim (\text{Pop} \oplus Q)} - \ln P_{\Psi^*}(i|y) \\ &= (\ln 2) - \frac{1}{2} \left( KL \left( \text{Pop}, \frac{\text{Pop} + Q}{2} \right), KL \left( Q, \frac{\text{Pop} + Q}{2} \right) \right) \\ &= (\ln 2) + JSD(\text{Pop}, Q) \end{aligned}$$

$$\Phi^* = \operatorname{argmin}_{\Phi} JSD(\text{Pop}, P_{\Phi})$$

$$0 \leq JSD(P, Q) \leq \ln 2$$

# Converting to Cross Entropy

Goodfellow, 2014

In Goodfellow's original paper he expressed a preference for cross entropy loss (the fundamental equation) over Jensen-Shannon loss.

$$\Phi^* = \operatorname{argmin}_{\Phi} JSD(\text{Pop}, P_{\Phi})$$

VS.

$$\Phi^* = \operatorname{argmin}_{\Phi} H(\text{Pop}, P_{\Phi})$$

He presented a modification to the GAN adversarial objective that yields cross-entropy loss rather than Jensen-Shanon loss.

# Converting to Cross Entropy

## Goodfellow, 2014

$$\Psi^*(\Phi) = \operatorname{argmin}_{\Psi} E_{(i,y) \sim (\text{Pop} \uplus P_\Phi)} - \ln P_\Psi(i|y)$$

Assume:  $P_{\Psi^*}(1|y) = \frac{\text{Pop}(y)}{\text{Pop}(y) + P_\Phi(y)}$

Define:  $f_{\Psi^*}(y) \doteq \frac{P_{\Psi^*}(1|y)}{P_{\Psi^*}(-1|y)}$

$$= \frac{\text{Pop}(y)}{P_\Phi(y)}$$

## Converting to Cross Entropy

$$\Phi^* = \underset{\Phi}{\operatorname{argmax}} E_{y \sim \text{Pop}} f_{\Psi^*}(y)$$

$$\begin{aligned} \nabla_{\Phi} E_{y \sim P_{\Phi}} f_{\Psi^*}(y) &= \nabla_{\Phi} \sum_y P_{\Phi}(y) f_{\Psi^*}(y) \\ &= \sum_y P_{\Phi}(y) f_{\Psi^*}(y) \nabla_{\Phi} \ln P_{\Phi}(y) \\ &= \sum_y \text{Pop}(y) \nabla_{\Phi} \ln P_{\Phi}(y) \\ &= E_{y \sim \text{Pop}} \nabla_{\Phi} \ln P_{\Phi}(y) \\ &= \nabla_{\Phi} E_{y \sim \text{Pop}} \ln P_{\Phi}(y) \end{aligned}$$

## Vanishing Gradients

The discriminator typically “wins”.

The log loss goes to zero (becomes exponentially small) and there is no gradient to guide the generator.

In this case the learning stops and the generator is blocked from minimizing  $\text{JSD}(\text{Pop}, P_\Phi)$ .

## A Heuristic Fix

We continue to use

$$\Psi^*(\Phi) = \operatorname{argmin}_{\Psi} E_{(y,s) \sim (\text{Pop} \uplus P_\Phi)} - \ln P_\Psi(s|y)$$

But switch the optimization for  $\Phi$  from

$$\Phi^* = \operatorname{argmax}_{\Phi} E_{y \sim P_\Phi} - \ln P_\Psi(-1|y)$$

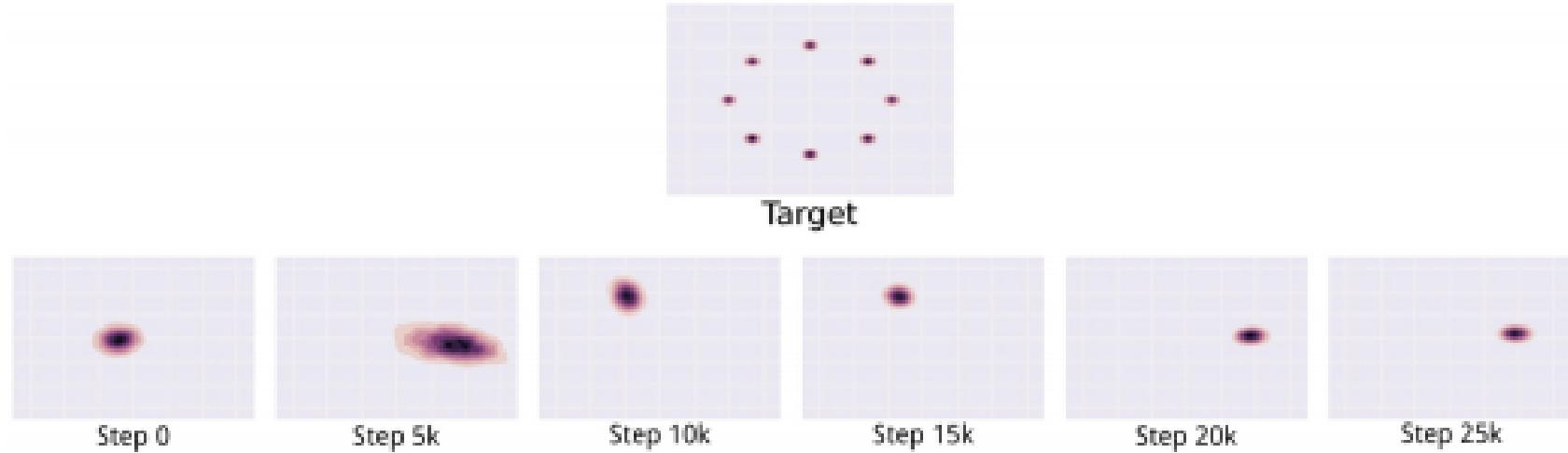
to

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{y \sim P_\Phi} - \ln P_\Psi(1|y)$$

It can be shown that  $-\ln P_\Psi(1|y)$  is essentially the margin of the binary classifier  $\Psi$ .

# Mode Collapse a.k.a Mode Dropping

The generator distribution drops portions of the population.



## Unstable Training

Joint SGD is not the same as nested max-min.

Consider

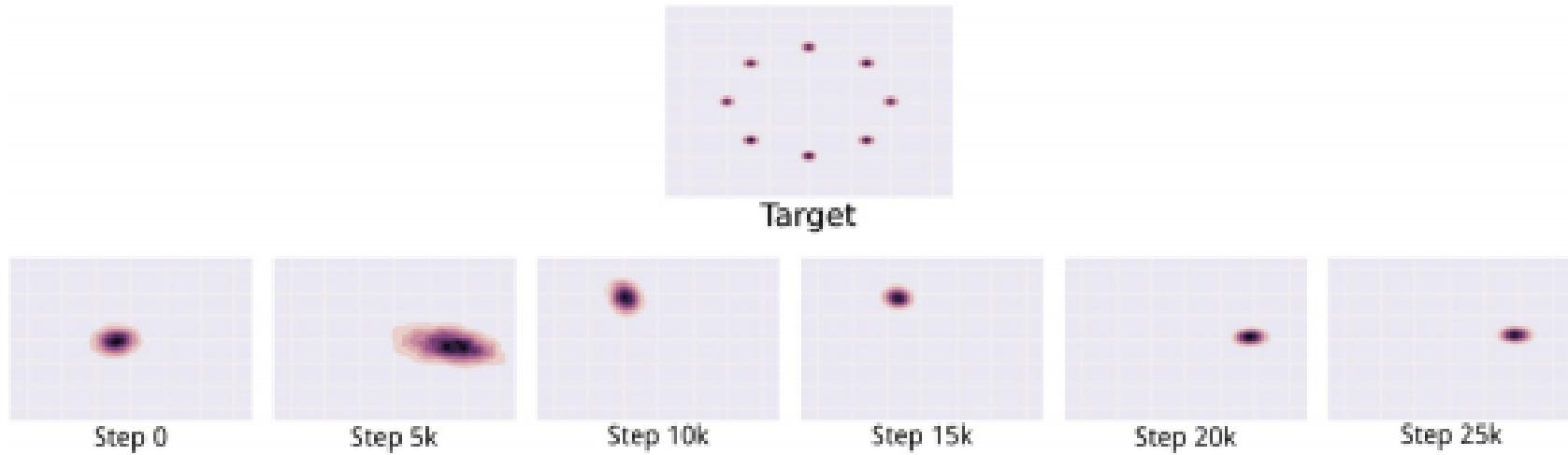
$$\max_x \min_y xy$$

A Nash equilibrium is  $x = y = 0$ .

Simultaneous gradient flow yields a circle.

## Mode Collapse a.k.a Mode Dropping

The generator distribution drops portions of the population.



# Pros and Cons of GAN Evaluation Measures

Borji, Oct 2018

We would like a rate-distortion metric on distribution models.

This has not yet been achieved for GANs.

Evaluation of GANs always involves, at least in part, subjective judgments of naturalness.

Sometimes automated metrics are also used.

The above paper discusses various proposed automated metrics of GAN performance. Current automated metrics are questionable.

## Contrastive GANS (TZ)

## The Discriminator also Models the Population

Fix the generator at a “noise distribution”  $Q$  where  $Q(y)$  is computable, and train a discriminator.

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{(i,y) \sim (\text{Pop} \uplus Q)} - \ln P_{\Psi}(i|y)$$

## The Discriminator also Models the Population

$$\Phi^* = \operatorname{argmin}_{\Phi} E_{(i,y) \sim (\text{Pop} \oplus Q)} - \ln P_{\Phi}(i|y)$$

Assuming Universality

$$P_{\Phi^*}(1|y) = P(1|y) = \frac{P(1 \text{ and } y)}{P(y)} = \frac{\frac{1}{2}\text{Pop}(y)}{\frac{1}{2}\text{Pop}(y) + \frac{1}{2}Q(y)}$$

$$P_{\Phi^*}(1|y)(\text{Pop}(y) + Q(y)) = \text{Pop}(y)$$

$$\frac{\text{Pop}(y)}{Q(y)} = \frac{P_{\Phi^*}(1|y)}{1 - P_{\Phi^*}(1|y)}$$

## Discrimination

The discrimination estimate of the population distribution is poor when the discrimination problem is easy — when the discrimination loss can be driven close to zero.

# Noise Contrastive Estimation

Gutmann and Hyvärinen, 2010

Consider a population distribution  $\text{Pop}$  and a fixed generating distribution  $Q$  where  $Q(y)$  is computable.

Define the distribution  $\text{Pop} \hookrightarrow Q^k$  to be the result of drawing one “positive” from  $\text{Pop}$  and  $k$  IID negatives from  $Q$ ; then inserting the positive at a random position among the negatives; and returning  $(i, y_1, \dots, y_{k+1})$  where  $i$  is the index of the positive.

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{(i, y_1, \dots, y_{k+1}) \sim (\text{Pop} \hookrightarrow Q^k)} -\ln P_\Psi(i | y_1, \dots, y_{k+1})$$

## Contrastive Estimation

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{(i, y_1, \dots, y_{k+1}) \sim (\text{Pop} \hookrightarrow Q^k)} -\ln P_\Psi(i | y_1, \dots, y_{k+1})$$

Note that  $\text{Pop} \hookrightarrow Q^1$  requires a choice between two  $y$ 's while  $\text{Pop} \uplus Q$  classifies a single  $y$  — these are different.

The contrastive task gets more difficult as  $k$  gets larger.

## Constrastive Estimation

$$\Psi^* = \operatorname{argmin}_{\Psi} E_{(i, y_1, \dots, y_{k+1}) \sim (\text{Pop} \hookrightarrow Q^k)} -\ln P_\Psi(i | y_1, \dots, y_{k+1})$$

$$P_\Psi(i | y_1, \dots, y_{k+1}) \doteq \underset{i}{\operatorname{softmax}} s_\Psi(y_i)$$

Theorem: Assuming universality:

$$s_{\Psi^*}(y) = \left( \ln \frac{\text{Pop}(y)}{Q(y)} \right) + C \quad \text{for arbitrary } C$$

$$\text{Pop}(y) = \underset{y}{\operatorname{softmax}} s_{\Psi^*}(y) - \ln Q(y) \quad Z \text{ must be estimated}$$

## Proof

$$\begin{aligned} P(i \text{ and } y_1, \dots, y_{k+1}) &= \frac{1}{k+1} \text{Pop}(y_i) \prod_{j \neq i} Q(y_j) \\ &= \alpha \frac{\text{Pop}(y_i)}{Q(y_i)}, \quad \alpha = \frac{1}{k+1} \prod_i Q(y_i) \end{aligned}$$

$$\begin{aligned} P(i \mid y_1, \dots, y_{k+1}) &= \frac{P(i \text{ and } y_1, \dots, y_{k+1})}{\sum_i P(i \text{ and } y_1, \dots, y_{k+1})} \\ &= \underset{i}{\text{softmax}} \left( \ln \frac{\text{Pop}(y_i)}{Q(y_i)} \right) + C \end{aligned}$$

## Constrastive Estimation

Like the discrimination estimate, the contrastive estimate of Pop is poor when the contrastive task is easy — when the contrastive loss can be driven near zero.

However, the contrastive task can be made more difficult by increasing  $k$ .

## Contrastive GANs (TZ)

Discriminative GAN:

$$\Phi^* = \operatorname{argmax}_{\Phi} \min_{\Psi} E_{(i,y) \sim (\text{Pop} \uplus P_\Phi)} - \ln P_\Psi(i|y)$$

Contrastive GAN:

$$\Phi^* = \operatorname{argmax}_{\Phi} \min_{\Psi} E_{(i,y_1, \dots, y_{k+1}) \sim (\text{Pop} \hookrightarrow P_\Phi^k)} - \ln P_\Psi(i|y_1, \dots, y_{k+1})$$

Assuming Universality:

$$P_{\Phi^*} = \text{Pop}$$

## Summary

GANs have not generally proved useful for representation learning for discriminative tasks such as image segmentation, speech recognition, or machine translation.

I predict that there will ultimately be better ways to model distributions (as in language modeling).

I predict that in a few years discriminators will be limited to enforcing perceptual naturalness in applications such as text to speech and image decompression.

**END**