

Computational Systems Biology Deep Learning in the Life Sciences

6.802 6.874 20.390 20.490 HST.506

David Gifford
Lecture 9
March 5, 2020

Generative Models



<http://mit6874.github.io>

Why generative models?

We can sample new examples from a generative models

- Generate new examples from model fit to training data
 - Sampling from input distribution
 - Optionally optimized with respect to a metric
- Reveals what models understand
 - What is the best example of a written digit?
 - What is the best example of a celebrity?
- Transform examples with respect to one or more metrics
 - Improve sentiment of text
 - Perform multi-objective optimization of antibodies

- In **generative modeling**, we'd like to train a network that models a distribution, such as a distribution over images.
- One way to judge the quality of the model is to sample from it.
- This field has seen rapid progress:

1 6 4 1 4 1 0
 7 2 8 8 4 9 4
 8 3 7 4 0 4 4
 3 7 2 1 7 7 7
 1 4 4 4 1 0 9
 3 0 5 9 5 2 7
 5 1 9 8 1 9 6

2009

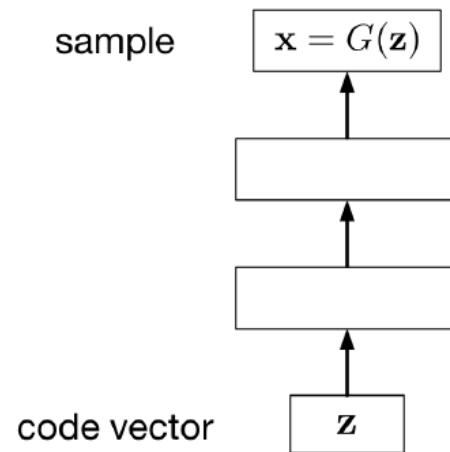


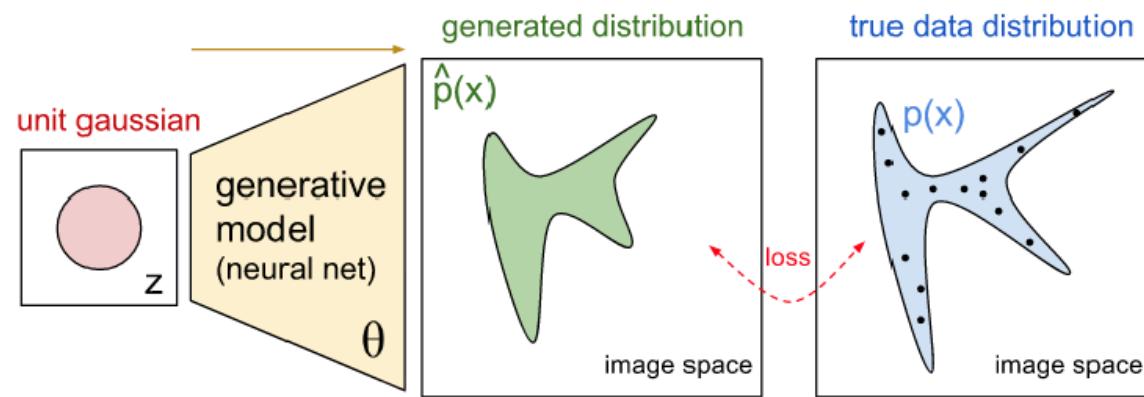
2015



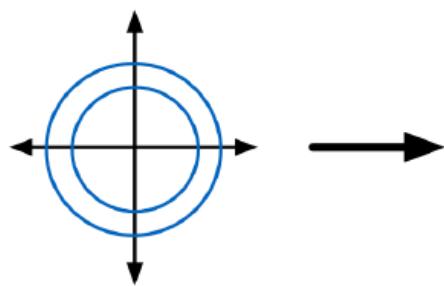
2018

- Implicit generative models implicitly define a probability distribution
- Start by sampling the code vector \mathbf{z} from a fixed, simple distribution (e.g. spherical Gaussian)
- The generator network computes a differentiable function G mapping \mathbf{z} to an \mathbf{x} in data space

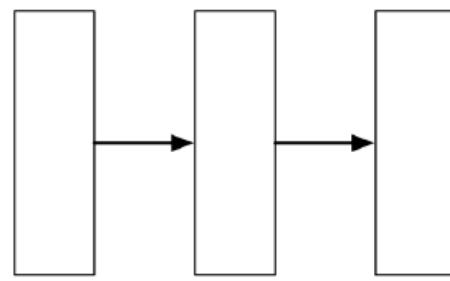




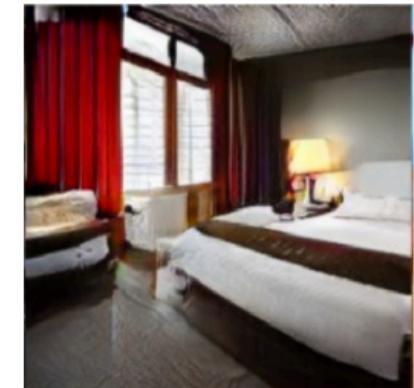
<https://blog.openai.com/generative-models/>



Each dimension of the code vector is sampled independently from a simple distribution, e.g. Gaussian or uniform.



This is fed to a (deterministic) generator network.



The network outputs an image.

This sort of architecture sounded preposterous to many of us, but amazingly, it works.

Three example generative models

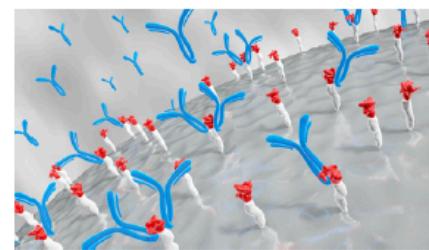
- Variational autoencoders
- Generative Adversarial Networks
- CycleGANs
- For each you should understand the loss function

Variational Autoencoders can provide improved examples

Improving outcomes for structured sequences¹²

- Discrete sequence data is commonplace (eg. text, proteins/genes)
sequence $x = (s_1, \dots, s_T) \in \mathcal{X}$ where each symbol $s_t \in \mathcal{S}$ (discrete vocabulary)
- Tiny fraction of \mathcal{X} represents sequences likely to occur naturally
(ie. those which appear *realistic*)
- Each sequence x is associated with outcome $y \in \mathbb{R}$

 [-] **DragonGodGrapha** 6 points 2 years ago
 $= y$
| This comment deserves more upvotes!
 $= x$
[permalink](#) [embed](#) [save](#) [parent](#) [give gold](#)



→ $y, x = \text{ASVKVSKC}$

¹² J. Mueller et al. Sequence to better sequence: Continuous revision of combinatorial structures. *ICML*, 2017.

Problem setup

- Dataset $\mathcal{D}_n = \{(x_i, y_i)\}_{i=1}^n \stackrel{iid}{\sim} p_{XY}$ of sequence-outcome pairs
- p_X = generative model of the *natural* sequences (unknown)

Goal: Given new sequence $x_0 \sim p_X$ (with unknown outcome),
quickly identify a revision x^* with superior expected outcome

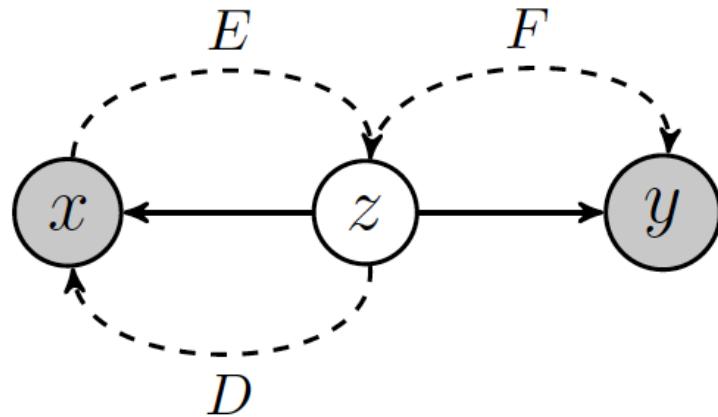
$$x^* = \operatorname{argmax}_{x \in \mathcal{C}_{x_0}} \mathbb{E}[Y \mid X = x]$$

$\mathcal{C}_{x_0} \subset \mathcal{X}$ = feasible set of sequences which appear natural
and are minor revisions of x_0

Desiderata for our revision procedure: $x_0 \rightarrow x^*$

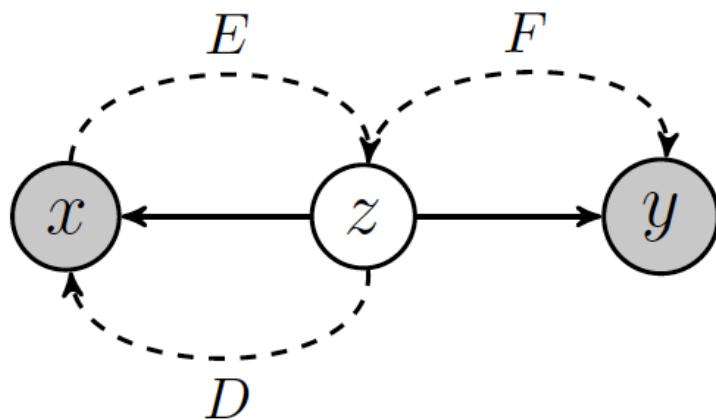
- Produces natural sequences
- Preserves intrinsic similarity
- Improves outcomes
- Computationally efficient

Probabilistic generative model



- Continuous latent factors $Z \in \mathbb{R}^d$ produce sequence $X +$ outcome Y
 - Prior: $p_Z = N(0, \mathbf{I})$
- Approximate inference maps $F \simeq y|z$, $E \simeq z|x$ & likelihood function $D \simeq x|z$ parameterized via three neural networks $\mathcal{F}, \mathcal{E}, \mathcal{D}$
- Parameters of $\mathcal{F}, \mathcal{E}, \mathcal{D}$ learned jointly using various objective functions that are added together

Probabilistic generative model



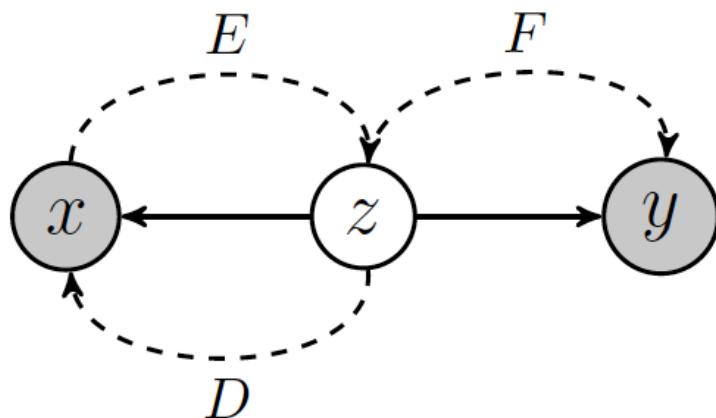
$$p_Z = N(0, \mathbf{I})$$

Why is this important?

Why does it make the task difficult?

- Continuous latent factors $Z \in \mathbb{R}^d$ produce sequence $X +$ outcome Y
 - Prior: $p_Z = N(0, \mathbf{I})$
- Approximate inference maps $F \simeq y|z$, $E \simeq z|x$ & likelihood function $D \simeq x|z$ parameterized via three neural networks $\mathcal{F}, \mathcal{E}, \mathcal{D}$
- Parameters of $\mathcal{F}, \mathcal{E}, \mathcal{D}$ learned jointly using various objective functions that are added together

Probabilistic generative model



$$p_Z = N(0, \mathbf{I})$$

Why is this important?

Find plausible revisions

Why does it make the task difficult?

$p(z | x)$ intractable

- Continuous latent factors $Z \in \mathbb{R}^d$ produce sequence $X +$ outcome Y
 - Prior: $p_Z = N(0, \mathbf{I})$
- Approximate inference maps $F \simeq y|z$, $E \simeq z|x$ & likelihood function $D \simeq x|z$ parameterized via three neural networks $\mathcal{F}, \mathcal{E}, \mathcal{D}$
- Parameters of $\mathcal{F}, \mathcal{E}, \mathcal{D}$ learned jointly using various objective functions that are added together

Kullback-Leibler Divergence

KL Divergence Definition

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

$$D_{KL}(P||Q) = - \sum_{x \in X} P(x) \log(Q(x)) + \sum_{x \in X} P(x) \log(P(x))$$

$$D_{KL}(P||Q) = H(P, Q) - H(P)$$

- Kullback-Leibler (KL) divergence measures of the difference between two probability distributions
- It is the extra information contained in P that is not in Q
- Thus it measures how far Q is from P

Learning the Encoder (E) and Decoder (D)

Variational inference is used to learn E and D

$$p(z|x) = \mathcal{N}(\mu_{z|x}, \text{diag}(\sigma_{z|x}^2))$$

We minimize the following two loss functions

$$\mathcal{L}_{rec}(x) = -E_{q_{E(z|x)}}[\log p_D(x|z)]$$

$$\mathcal{L}_{pri}(x) = KL(q_{E(z|x)}, p(z|x))$$

- We work to make $x' = D(E(x))$ lossless
- At the same time we shape observed sequences to a normal latent space
- More details on Variational Autoencoders - <https://arxiv.org/pdf/1312.6114.pdf>

Example of a bad latent representation

Suppose that for $x \in \mathcal{X}$:

$$E(x) = z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \in \mathbb{R}^d$$

$$\hat{y} = F(z) = F(z_1) \quad \text{and} \quad \hat{x} = D(z) = D(z_2)$$

Outcome prediction ignores z_2

Decoding step ignores z_1

- Changing $z_1 \implies$ change in predicted outcome \hat{y}
but *no* change in resulting sequence \hat{x} ☹

Preventing degenerate latent space encoding

Need to ensure E encodes z so both Y and D depend on the same components

We minimize \mathcal{L}_{inv} to ensure this dependence

$$\mathcal{L}_{inv} = E_{z \sim p_z}[F(z) - F(E(D(z)))]$$

- We can sample values of z to learn gradients to update E and F

Learning the objective function (F)

F maps a latent instance to its expected value Y

$$F(z) = E[Y|z]$$

We minimize \mathcal{L}_{mse} to learn F

$$\mathcal{L}_{mse}(x, y) = (y - F(E(x)))^2$$

- We will use gradients of F to optimize x

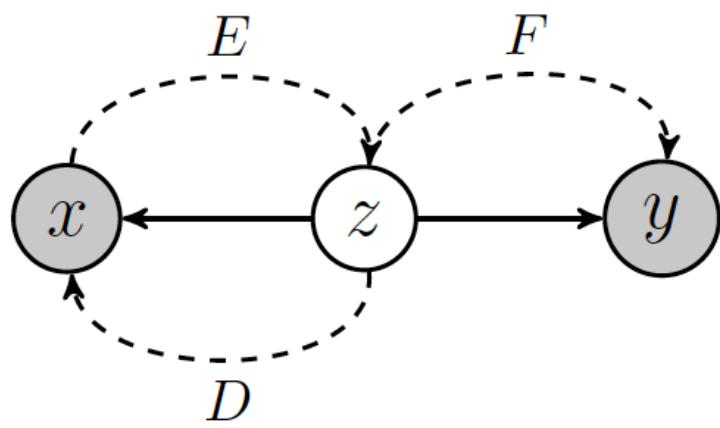
Overall VAE loss function

$$\mathcal{L}_{\text{rec}}(x) = -\mathbb{E}_{q_E(z|x)} [\log p_D(x \mid z)]$$

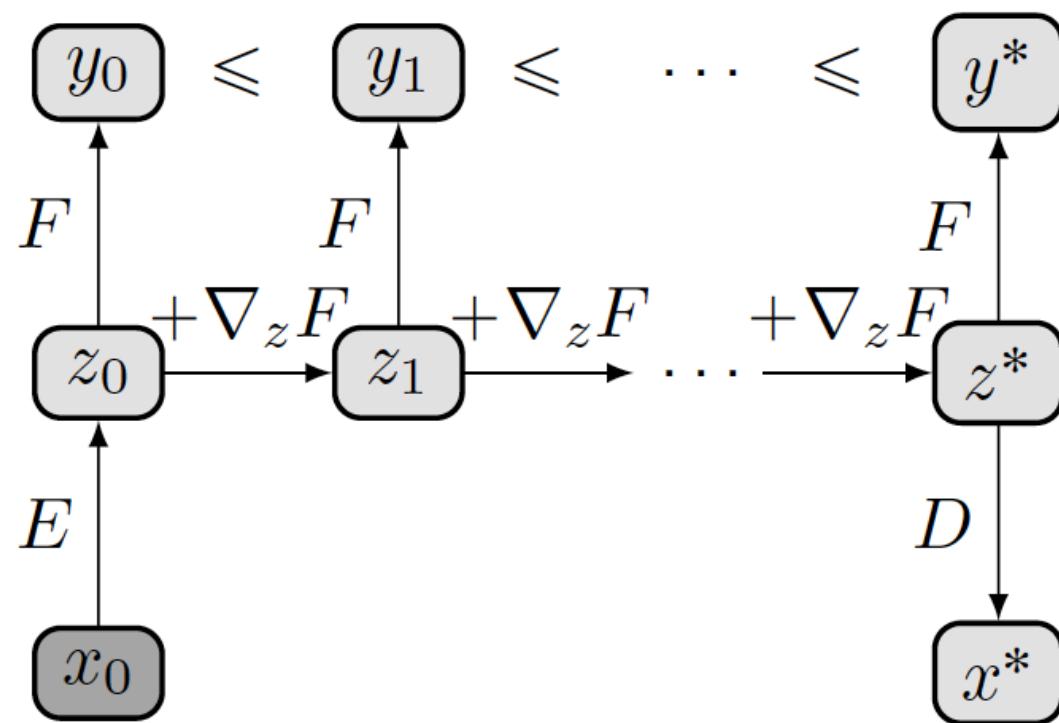
$$\mathcal{L}_{\text{pri}}(x) = \text{KL}(q_E(z \mid x) \parallel p_Z)$$

$$\mathcal{L}_{\text{mse}}(x, y) = [y - F(E(x))]^2$$

$$\mathcal{L}_{\text{inv}} = \mathbb{E}_{z \sim p_Z} [F(z) - F(E(D(z)))]^2$$



Revision framework



Proposing revisions

REVISE Algorithm

Input: sequence $x_0 \in \mathcal{X}$, constant $\alpha \in (0, |2\pi\Sigma_{z|x_0}|^{-\frac{1}{2}})$

Output: revised sequence $x^* \in \mathcal{X}$

- 1) Use \mathcal{E} to compute $q_E(z \mid x_0)$, $E(x_0) = \mathbb{E}_{q_E}[z \mid x_0]$
 - 2) Define $\mathcal{C}_{x_0} = \{z \in \mathbb{R}^d : q_E(z \mid x_0) \geq \alpha\}$ (ellipsoid)
 - 3) Find $z^* \approx \operatorname{argmax}_{z \in \mathcal{C}_{x_0}} F(z)$ (gradient ascent w/ log-barrier penalty)
 - 4) Return $x^* = D(z^*) \approx \operatorname{argmax}_{x \in \mathcal{X}} p_D(x \mid z^*)$ (greedy beam search)
-

Improving sentence positivity

- Data = 1M+ short sentences from BeerAdvocate reviews
- $y \in [0, 1]$: VADER sentiment compound score of each sentence¹⁴
- Apply methods to revise set of 1000 held-out sentences

¹⁴ C. Hutto and E. Gilbert. Vader: A parsimonious rule-based model for sentiment analysis of social media text.
ICWSM, 2014.

Improving sentence positivity

Model	$\Delta_Y(x^*)$	$L(x^*)$	$d(x^*, x_0)$
$\log \alpha = -10000$	0.52 ± 0.77	-8.8 ± 6.5	2.6 ± 3.3
$\log \alpha = -1$	0.31 ± 0.50	-7.6 ± 5.8	1.7 ± 2.6
$\lambda_{\text{inv}} = \lambda_{\text{pri}} = 0$	0.22 ± 1.03	-10.2 ± 7.0	3.3 ± 3.4
SEARCH	0.19 ± 0.56	-7.7 ± 4.2	3.0 ± 1.2

$\Delta_Y(x^*)$ = outcome improvement from revision (rescaled by variance of outcomes)

$L(x^*) = \hat{p}_X(x^*) - \hat{p}_X(x_0)$, \hat{p}_X = likelihood estimated via language model

$d(x^*, x_0)$ = Levenshtein (edit) distance

Improving sentence positivity

Model	Sentence	$\Delta_Y(x^*)$	$\Delta_L(x^*)$
x_0	this smells pretty bad.	-	-
$\log \alpha = -10000$	smells pretty delightful!	+2.8	-0.5
$\log \alpha = -1$	i liked this smells pretty.	+2.5	-2.8
$\lambda_{\text{inv}} = \lambda_{\text{pri}} = 0$	pretty this smells bad!	-0.2	-3.1
SEARCH	wow this smells pretty bad.	+1.9	-4.6
x_0	i like to support san diego beers.	-	-
$\log \alpha = -10000$	i love to support craft beers!	+0.5	+1.6
$\log \alpha = -1$	i like to support craft beers!	+0.1	+2.6
$\lambda_{\text{inv}} = \lambda_{\text{pri}} = 0$	i like to support you know.	0	+3.7
SEARCH	i like to super support san diego.	+0.7	-2.9
x_0	i'm not sure how old the bottle is.	-	-
$\log \alpha = -10000$	i definitely enjoy how old is the bottle is.	+3.0	-3.6
$\log \alpha = -1$	i'm sure not sure how old the bottle is.	+2.5	-6.8
$\lambda_{\text{inv}} = \lambda_{\text{pri}} = 0$	i'm sure better is the highlights when cheers.	+3.3	-9.2
SEARCH	i 'm not sure how the bottle is love.	+2.3	-3.3

Revising modern text in the language of Shakespeare

- Dataset of ~100K short sentences
- Each is either from Shakespeare with label $y = 0.9$ or a more contemporary source (from NLTK) with label $y = 0.1$
- Given new sentence, revise so that author is increasingly expected to be Shakespeare rather than contemporary source

Revising modern text in the language of Shakespeare

# Steps	Decoded Sentence
x_0	where are you, henry??
100	where are you, henry??
1000	where are you, royal??
5000	where art thou now?
10000	which cannot come, you of thee?
x^*	where art thou, keeper??
x_0	somewhere, somebody is bound to love us.
100	somewhere, somebody is bound to love us.
1000	courage, honey, somebody is bound to love us!
5000	courage man; 'tis love that is lost to us.
10000	thou, within courage to brush and such us brush.
x^*	courage man; somebody is bound to love us.
x_0	you are both the same size.
100	you are both the same.
1000	you are both wretched.
5000	you are both the king.
10000	you are both these are very.
x^*	you are both wretched men.

Desiderata for our revision procedure

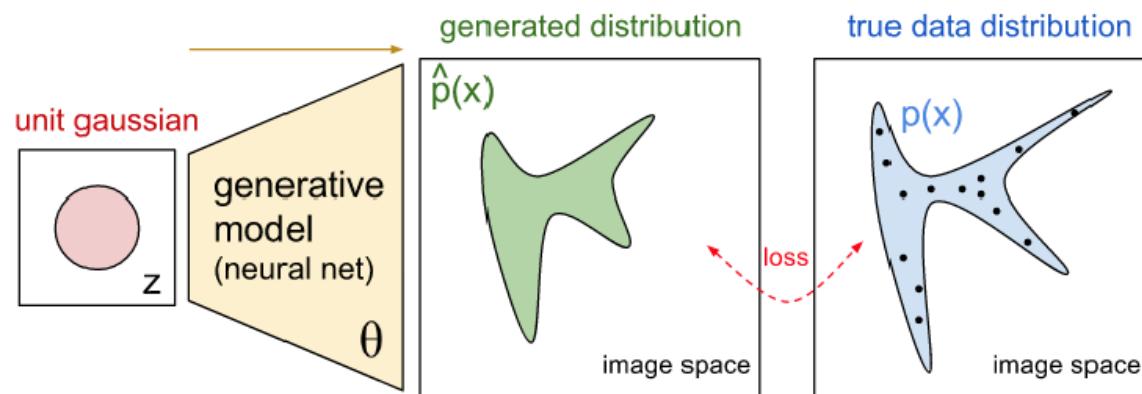
- Improves outcomes ✓
- Produces natural sequences ✓
- Computationally efficient ✓
- Preserves intrinsic similarity ✗

To model similarity, can leverage siamese RNN¹⁵ that embeds sentences into vector space where L_1 metric reflects human-labeled similarity

¹⁵ J. Mueller and A. Thyagarajan. Siamese Recurrent Architectures for Learning Sentence Similarity. AAAI, 2016.

Generative Adversarial Networks

We wish to learn a generative model that matches the true data distribution

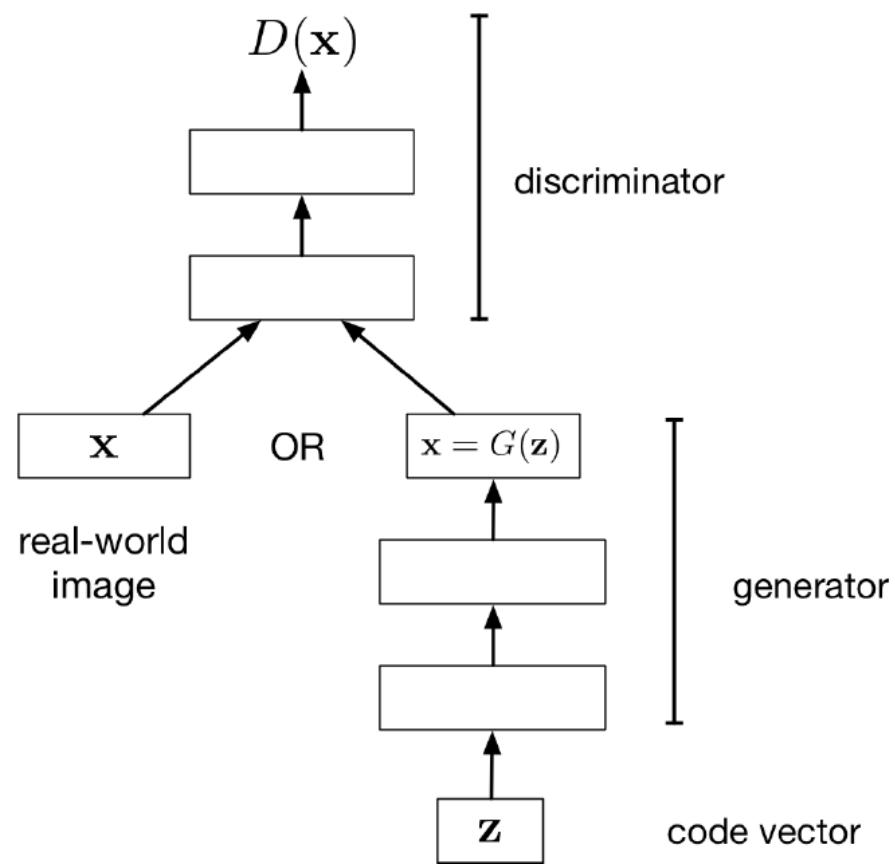


<https://blog.openai.com/generative-models/>

The Generative Adversarial Network (GAN) Game

- The advantage of implicit generative models: if you have some criterion for evaluating the quality of samples, then you can compute its gradient with respect to the network parameters, and update the network's parameters to make the sample a little better
- The idea behind **Generative Adversarial Networks (GANs)**: train two different networks
 - The **generator network** tries to produce realistic-looking samples
 - The **discriminator network** tries to figure out whether an image came from the training set or the generator network
- The generator network tries to fool the discriminator network

$D(x)$ is probability x is from the real-world



Learning the Discriminator (D) and Generator (G)

D is learned by a cross-entropy loss of scoring real vs. fake images

$$\mathcal{L}_D = E_{x \sim \text{Real}}[-\log D(x)] + E_z[-\log(1 - D(G(z)))]$$

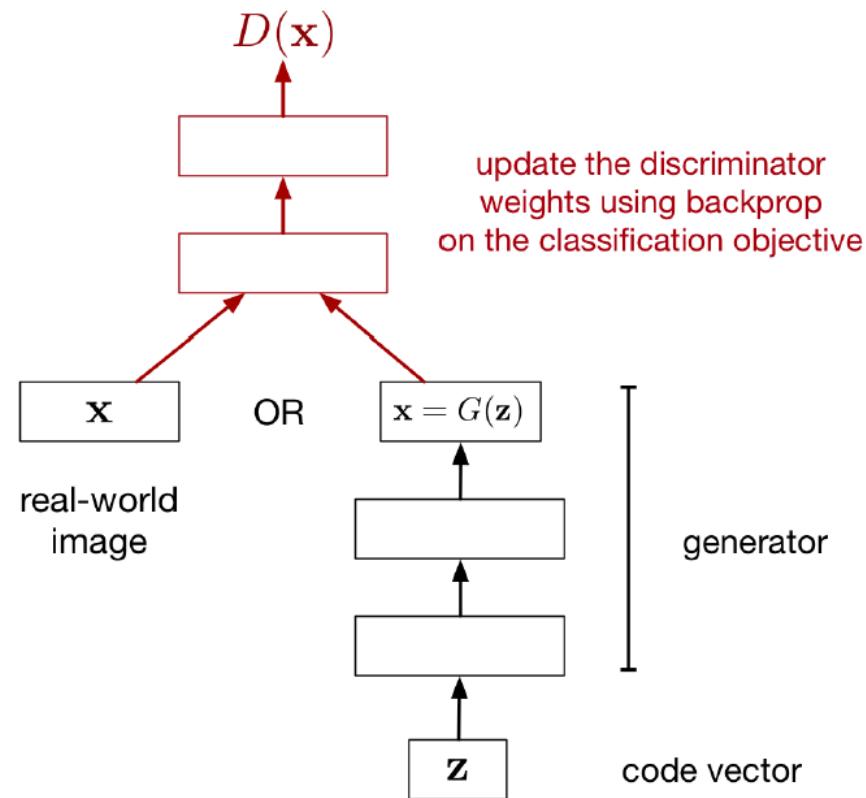
G is learned by trying to fool the discriminator

$$\mathcal{L}_G = E_z[\log(D(G(z)))]$$

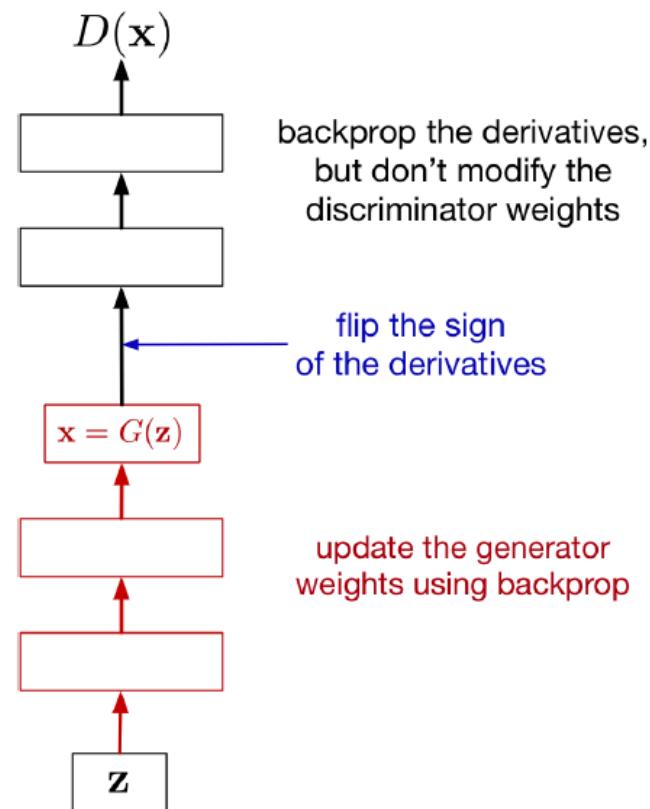
- This is a minimax formulation of a zero-sum game between D and G
- Either G wins (good outcome) or D wins (fakes are not very good!)
- Note that distribution of real examples is learned by D, and then communicated indirectly to G (it never sees a real example)

Update discriminator to maximize $D(\text{real})$ and minimize $D(G(z))$

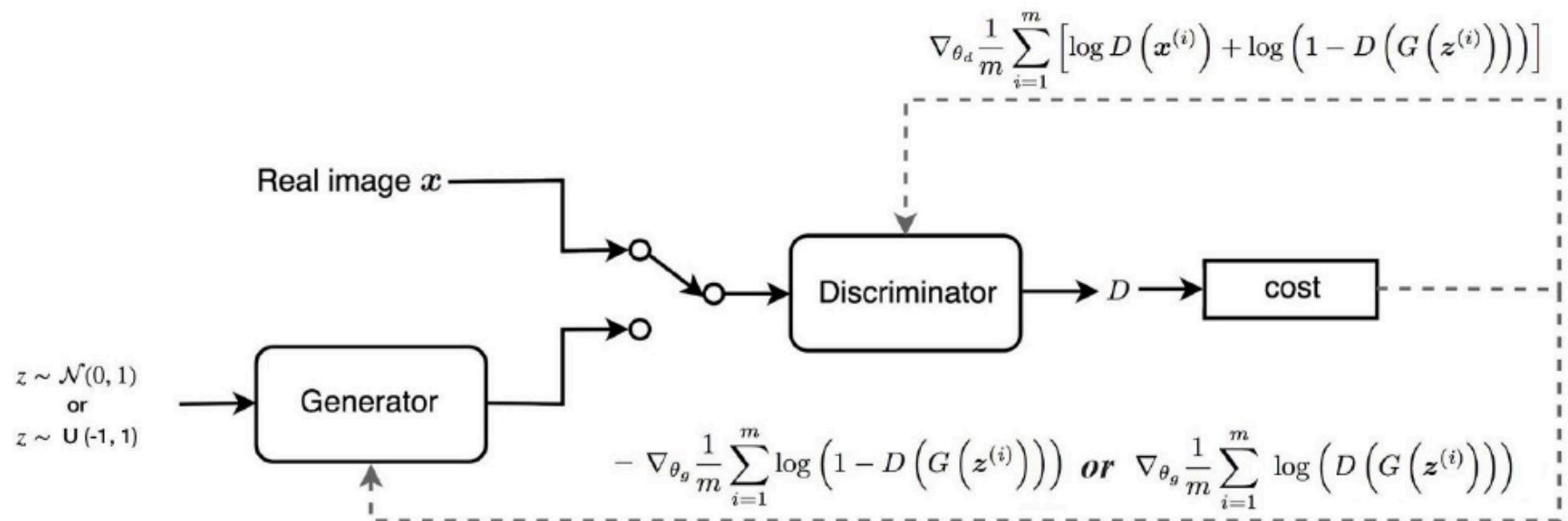
Updating the discriminator:



Update generator to maximize $D(G(z))$



Summary of GAN objective functions



GANs have become a bit of a fad

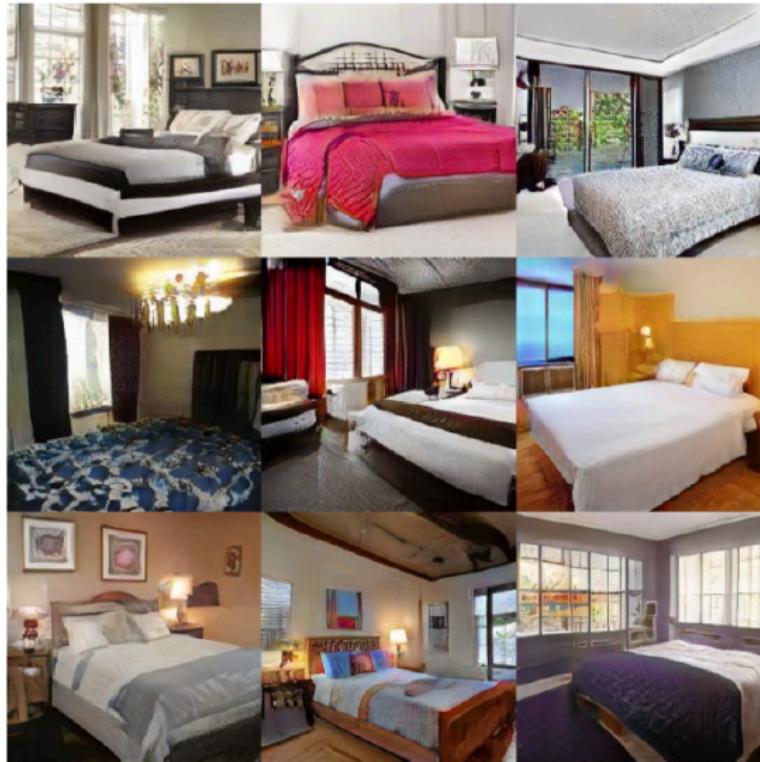
- Since GANs were introduced in 2014, there have been hundreds of papers introducing various architectures and training methods.
- Most modern architectures are based on the Deep Convolutional GAN (DC-GAN), where the generator and discriminator are both conv nets.
- GAN Zoo: <https://github.com/hindupuravinash/the-gan-zoo>
 - Good source of horrible puns (VEEGAN, Checkhov GAN, etc.)

Celebrities:



Karras et al., 2017. Progressive growing of GANs for improved quality, stability, and variation

Bedrooms:



Karras et al., 2017. Progressive growing of GANs for improved quality, stability, and variation

Objects:



Karras et al., 2017. Progressive growing of GANs for improved quality, stability, and variation

GANs can fail

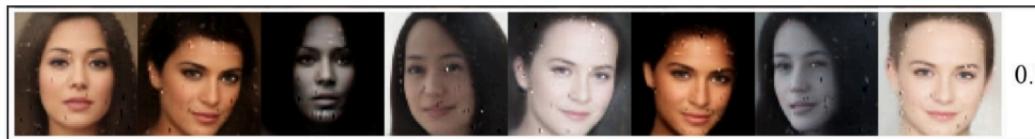
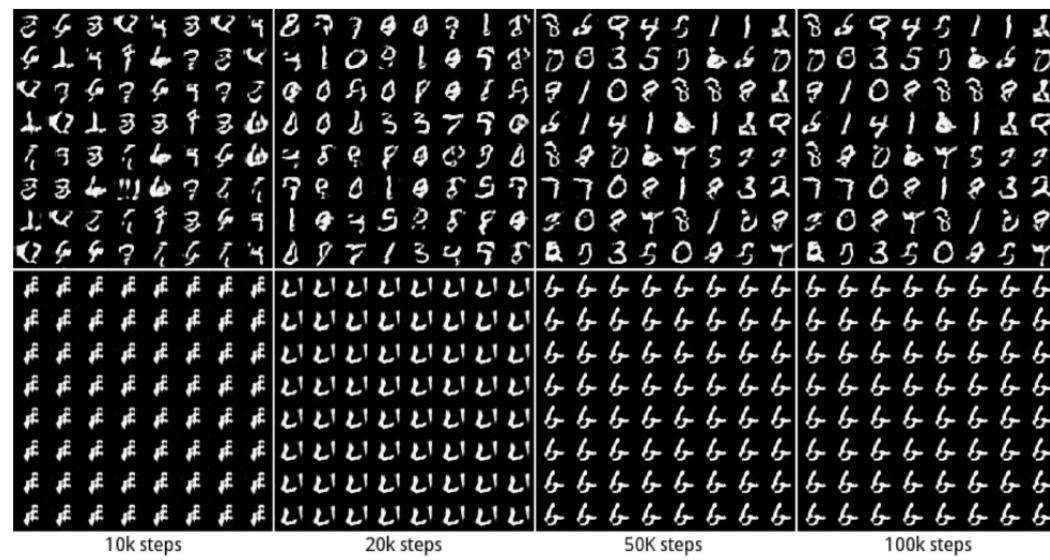
- GANs revolutionized generative modeling by producing crisp, high-resolution images.
- The catch: we don't know how well they're modeling the distribution.
 - Can't measure the log-likelihood they assign to held-out data.
 - Could they be memorizing training examples? (E.g., maybe they sometimes produce photos of real celebrities?)
 - We have no way to tell if they are dropping important modes from the distribution.
 - See Wu et al., "On the quantitative analysis of decoder-based generative models" for partial answers to these questions.

GAN Problems

- Non-convergence – model parameters oscillate and never converge
- Mode collapse – limited variety of samples from generator
- Diminished gradient – Discriminator is too successful and generator learns nothing
- Overfitting – imbalance between generator and discriminator
- Hyperparameter sensitivity – highly sensitive

https://medium.com/@jonathan_hui/gan-why-it-is-so-hard-to-train-generative-advisory-networks-819a86b3750b

Mode collapse shown in second row (all 6) and in images



<https://arxiv.org/pdf/1611.02163.pdf>

<https://arxiv.org/pdf/1703.10717.pdf>

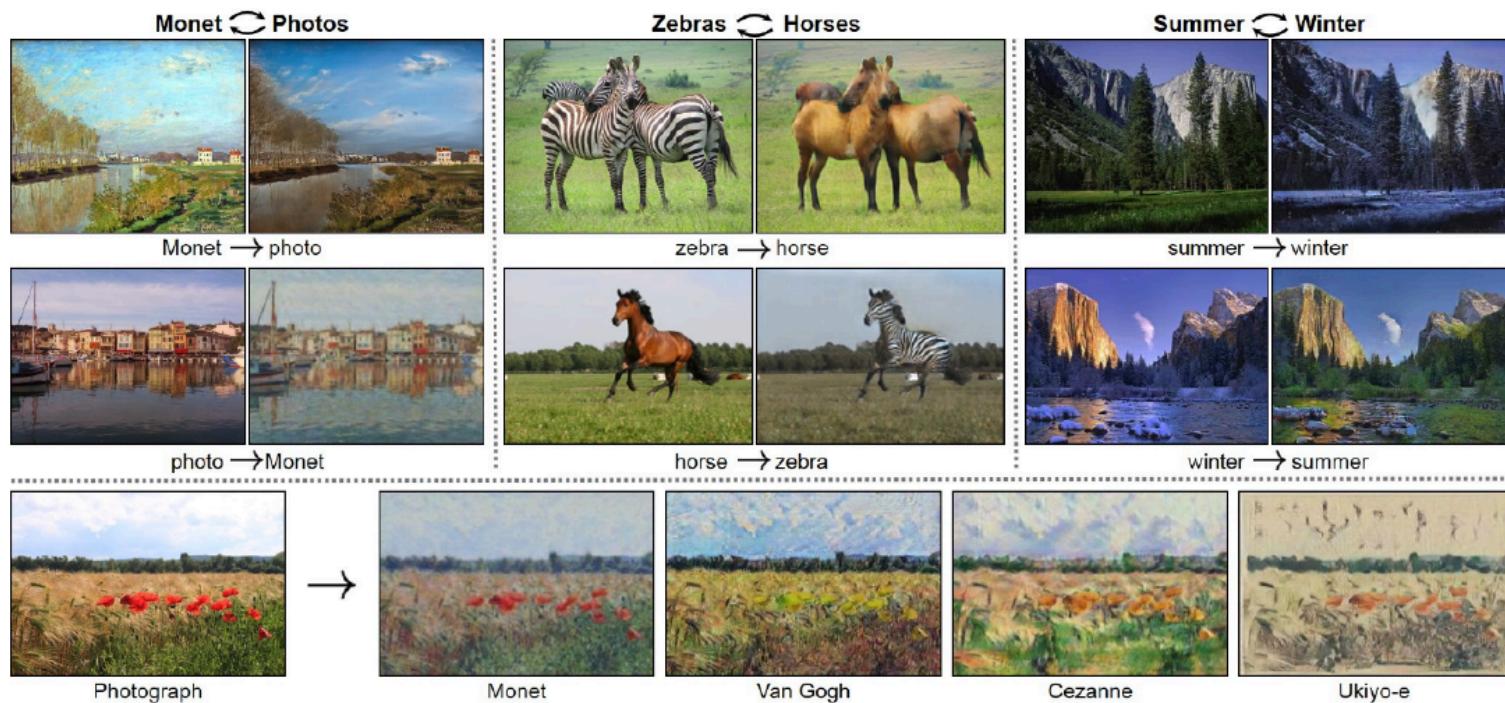
CycleGANs for style mapping

CycleGANs map between styles

- If we had paired data (same content in both styles), this would be a supervised learning problem. But this is hard to find.
- The CycleGAN architecture learns to do it from unpaired data.
 - Train two different generator nets to go from style 1 to style 2, and vice versa.
 - Make sure the generated samples of style 2 are indistinguishable from real images by a discriminator net.
 - Make sure the generators are **cycle-consistent**: mapping from style 1 to style 2 and back again should give you almost the original image.

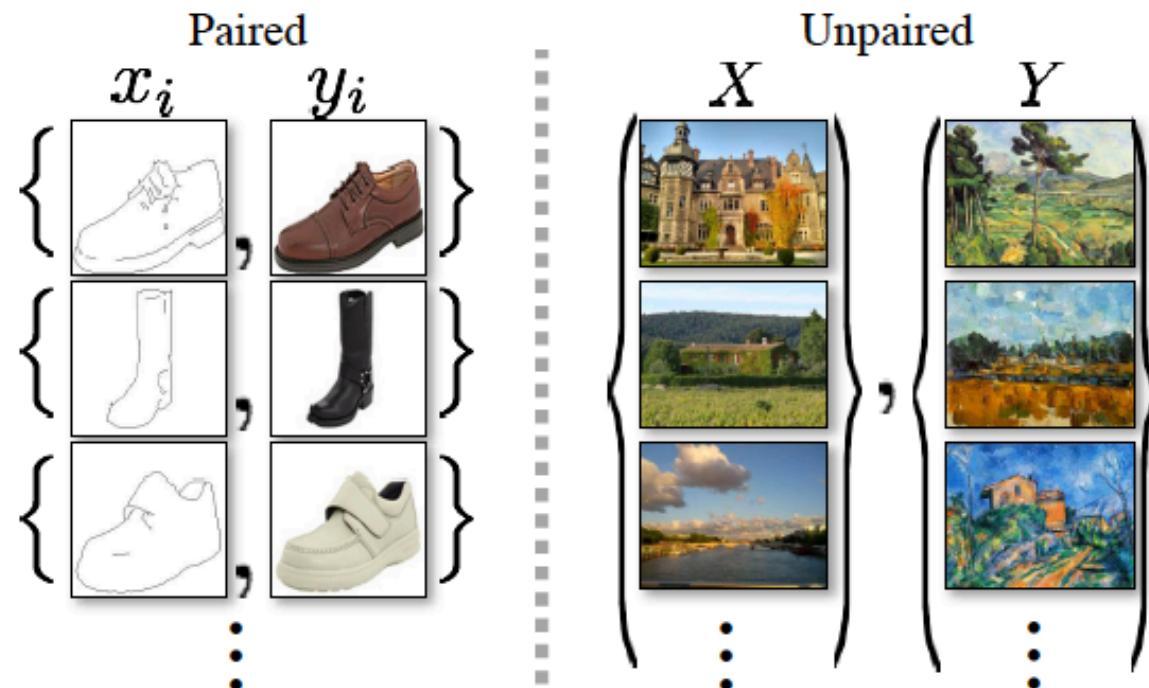
CycleGANs permit style transfer without matched training data

Style transfer problem: change the style of an image while preserving the content.



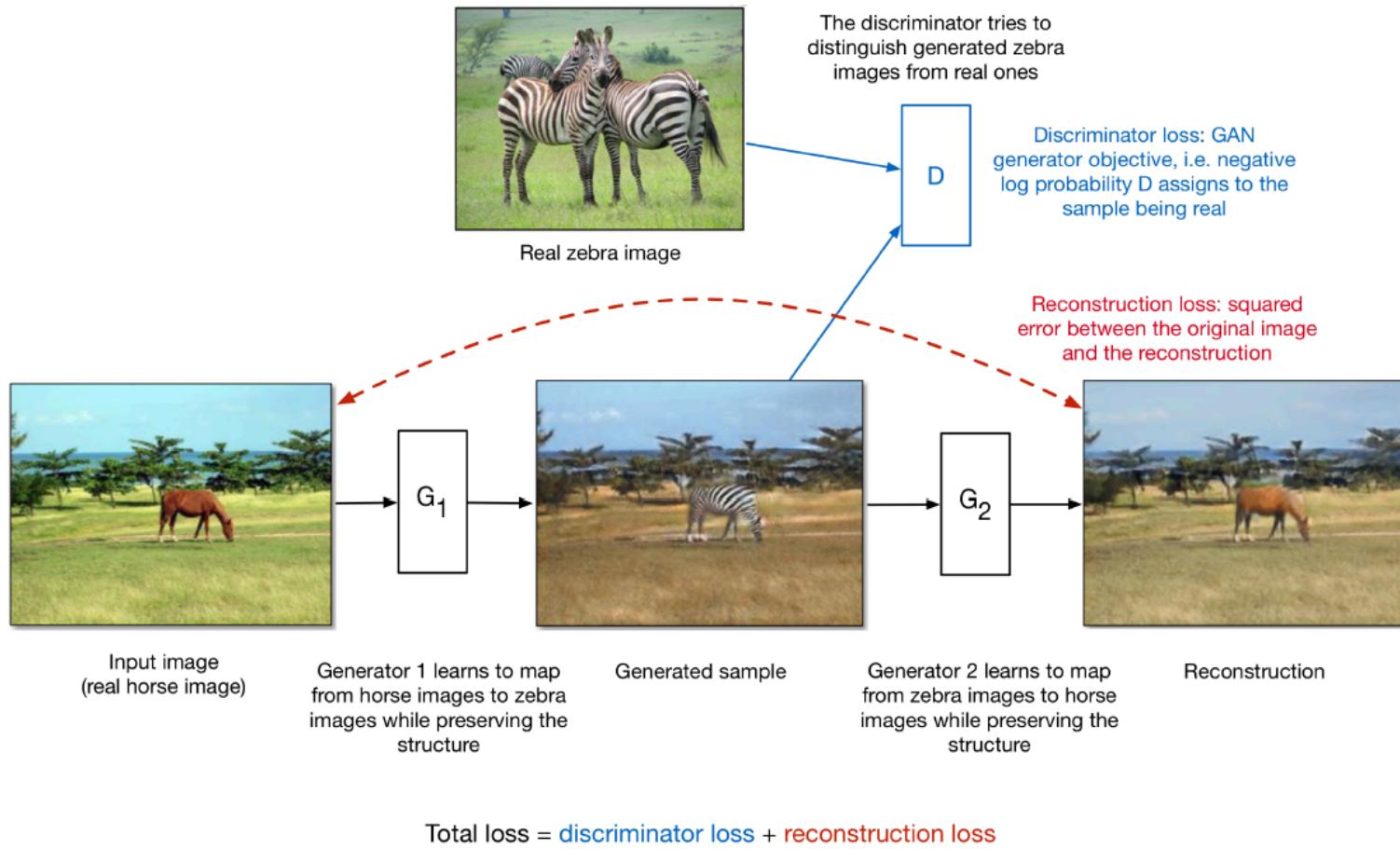
Data: Two unrelated collections of images, one for each style

CycleGANs permit style transfer without matched training data



<https://arxiv.org/pdf/1703.10593.pdf>

CycleGANs permit style transfer without matched training data



Cycle consistent loss functions for CycleGANs

The CycleGan loss consists of a GAN loss and a cycle consistency loss

$$\mathcal{L}_{GAN}(G, D_Y, X, Y) = E_{y \sim p_{data}(y)}[\log D_Y(y)] + E_{x \sim p_{data}(x)}[\log(1 - D_Y(G(x)))]$$

$$\mathcal{L}_{GAN}(F, D_X, X, Y) = E_{x \sim p_{data}(x)}[\log D_X(x)] + E_{y \sim p_{data}(y)}[\log(1 - D_X(F(y)))]$$

$$\mathcal{L}_{cyc}(G, F) = E_{x \sim p_{data}(x)}[||F(G(x)) - x||_1] + E_{y \sim p_{data}(y)}[||G(F(y)) - y||_1]$$

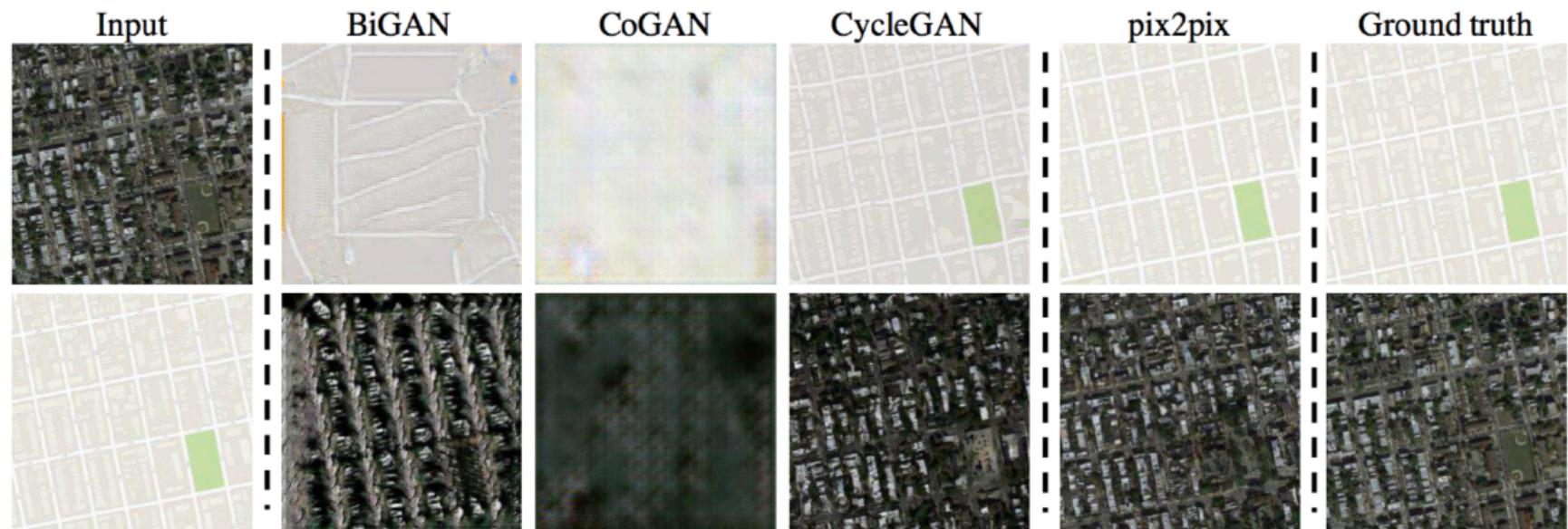
Overall loss: $\mathcal{L}(G, F, D_X, D_Y) = \mathcal{L}_{GAN}(G, D_Y, X, Y) + \mathcal{L}_{GAN}(F, D_X, X, Y) + \mathcal{L}_{cyc}(G, F)$

$$G^*, F^* = \arg \min_{G, F} \max_{D_X, D_Y} \mathcal{L}(G, F, D_X, D_Y)$$

- For CycleGans we do not need to use paired data
- Amazing that it works - See <https://arxiv.org/pdf/1703.10593.pdf>
- Once again, note that that distribution of real examples is learned by D_X and D_Y , and then communicated indirectly to F and G (they never sees a real example)

CycleGANs permit style transfer without matched training data

Style transfer between aerial photos and maps:



CycleGANs permit style transfer without matched training data

Style transfer between road scenes and semantic segmentations (labels of every pixel in an image by object category):

