

10-414/714 – Deep Learning Systems: Algorithms and Implementation

Generative Models

Fall 2025

Tianqi Chen (this time) and Tim Dettmers
Carnegie Mellon University

Outline

Generative adversarial training (GAN)

Diffusion models

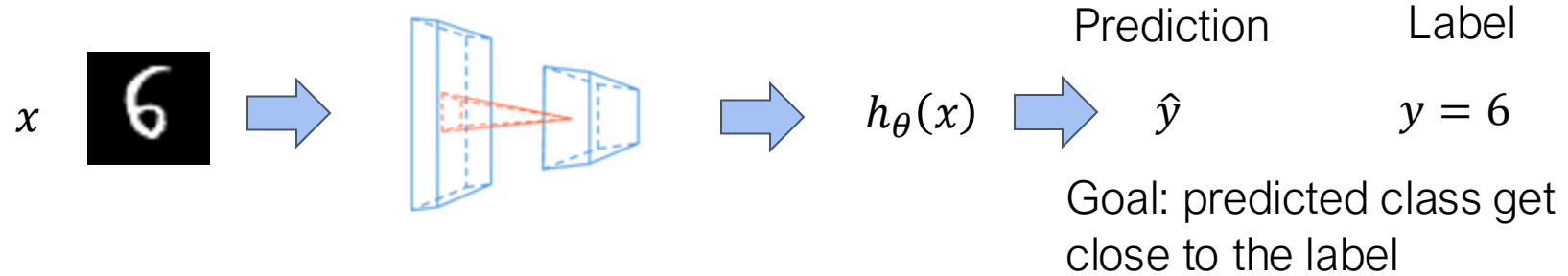
Outline

Generative adversarial training

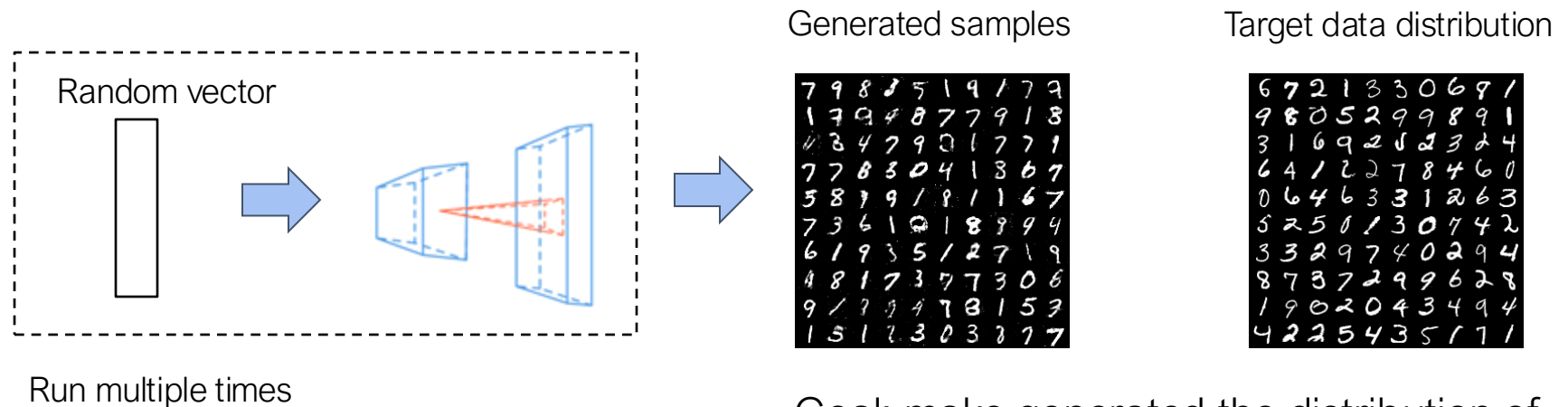
Diffusion models

From classifier to generator

Class probability predictor



Digits generator



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

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

Goal: make generated the distribution of generated samples “close” to target data distribution

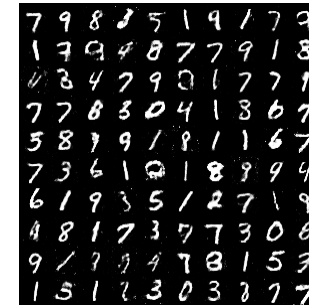
Define “distance” of distributions

Unlike supervised classification setting, the “goal” is less obvious

To build effective training mechanism, we need to define a “distance” between generated and real datasets and use that to drive the training.

What we really wanted, in text: make sure that the generated samples “looks real”.

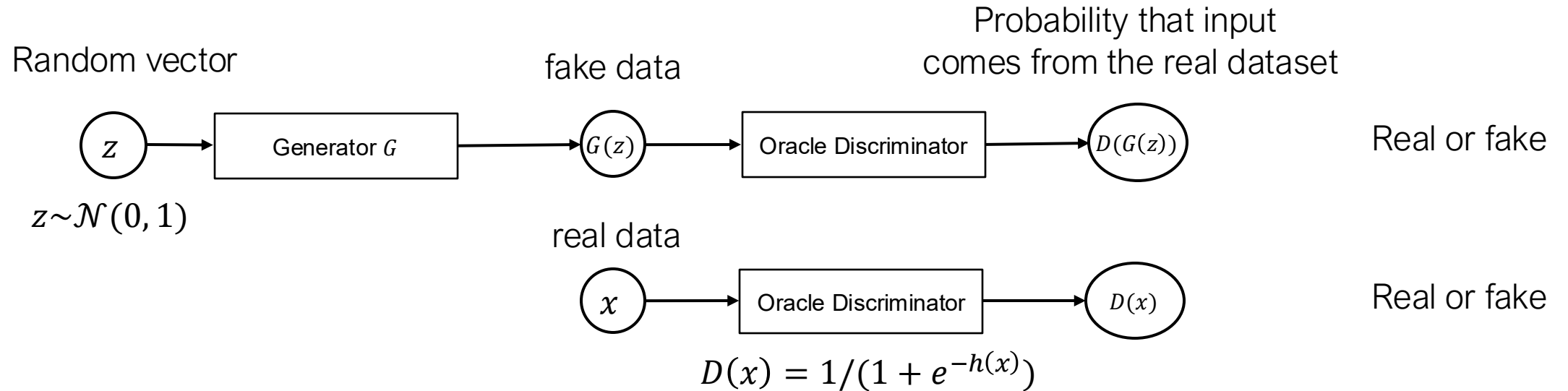
Generated samples



Target data distribution



Learn generator through an oracle discriminator

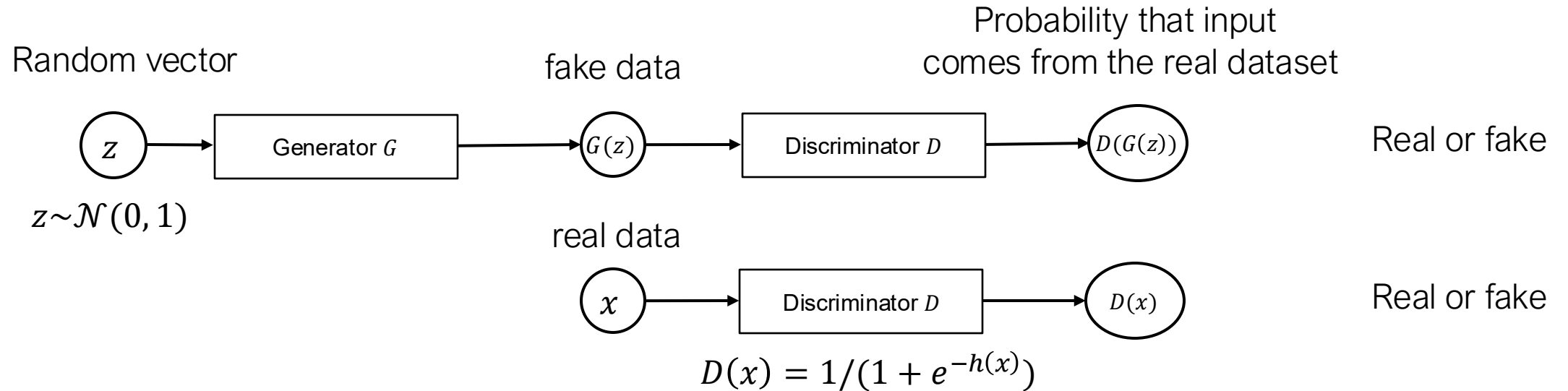


Assume that we have an oracle discriminator that can tell the difference between real and fake data. Then we need train the generator to “fool” the oracle discriminator.

We need to maximize the discriminator loss

$$\text{Generator objective: } \max_G \{-\mathbb{E}_{z \sim \text{Noise}} \log(1 - D(G(z)))\}$$

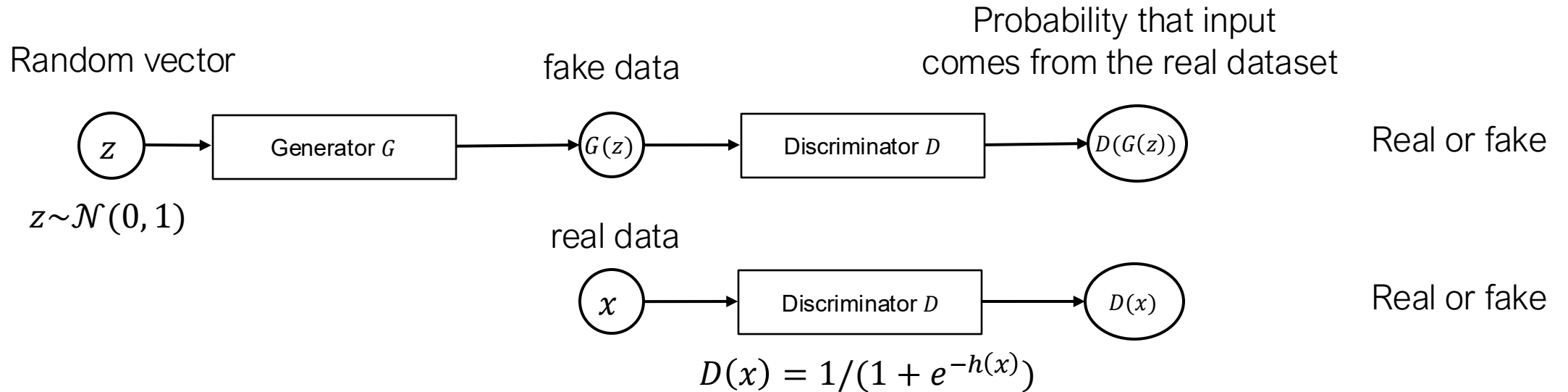
Learning the discriminator



We do not have an oracle discriminator, but we can learn it using the real and generated fake data.

Discriminator objective $\min_D \{-E_{x \sim \text{Data}} \log D(x) - E_{z \sim \text{Noise}} \log(1 - D(G(z)))\}$

Generative adversarial network

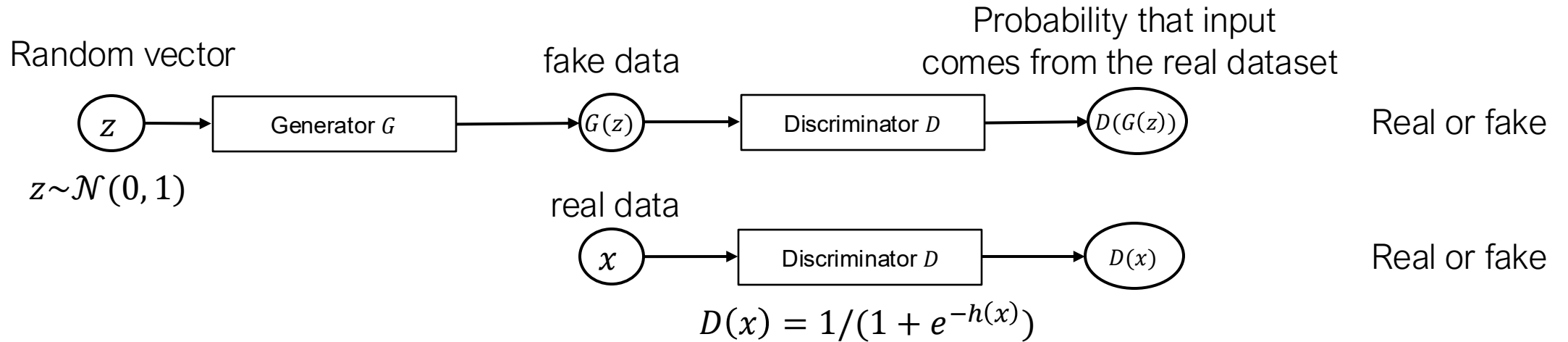


Putting it together, it becomes an “minimax” game between D and G

$$\min_D \max_G \{-E_{x \sim \text{Data}} \log D(x) - E_{z \sim \text{Noise}} \log(1 - D(G(z)))\}$$

In practice, we usually optimize G using $\min_G \{-E_{z \sim \text{Noise}} \log(D(G(z)))\}$, maximize the probability that discriminator predicts generated image is real

Generative adversarial training in practice



Iterative process

- Discriminator update
 - Sample minibatch of $D(G(z))$, get a minibatch of $D(x)$
 - Update D to minimize $\min_D \{-E_{x \sim \text{Data}} \log D(x) - E_{z \sim \text{Noise}} \log(1 - D(G(z)))\}$
- Generator update
 - Sample minibatch of $D(G(z))$
 - Update G to minimize $\min_G \{-E_{z \sim \text{Noise}} \log(D(G(z)))\}$, this can be done by feeding label=1 to to the model

Outline

Generative adversarial training

Diffusion models

Stochastic Differential Equation

Describes a stochastic movement of x_t in the space

Wiener process
(gaussian white noise, see
the discrete view)

Stochastic differential equation (SDE):

$$dx_t = f(x_t) dt + \sqrt{2D(x_t, t)} dW_t$$

Drift term,
pulls towards modes

Diffusion term,
inject noise

Discrete simulation of the process,
with small step size η_t

$$x_{t+1} \leftarrow x_t + \eta_t f(x_t) + \mathcal{N}(0, 2\eta_t D(x_t, t))$$

Forward Diffusion Process

Data distribution

$q(x_0)$



Close to standard normal distribution

$q(x_T)$

Gradually add gaussian noise

Forward SDE:
$$dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)} dW_t$$

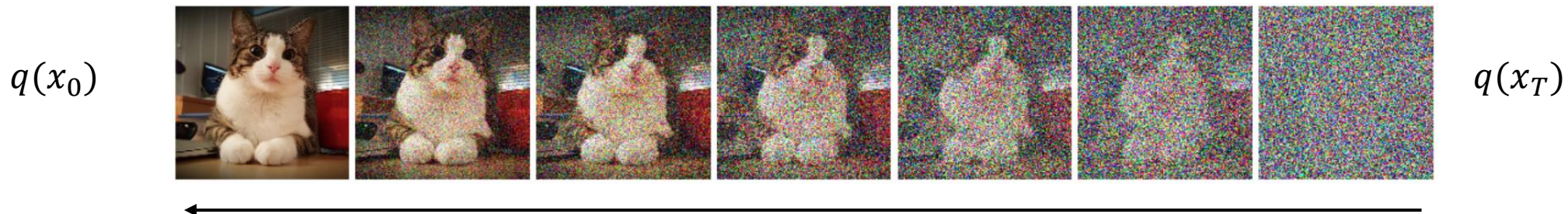
Discretized form:
$$x_{t+1} \leftarrow (1 - \frac{1}{2}\eta_t\beta(t))x_t + \mathcal{N}(0, \eta_t\beta(t))$$

Drifts toward 0

Add noise

Forward diffusion process takes image x_0 and generate white noise x_T

Reverse (Denoising) Process



Forward SDE: $dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)} dW_t$ It is easy to sample from $q(x_0, \dots, x_t \dots x_T | x_0)$

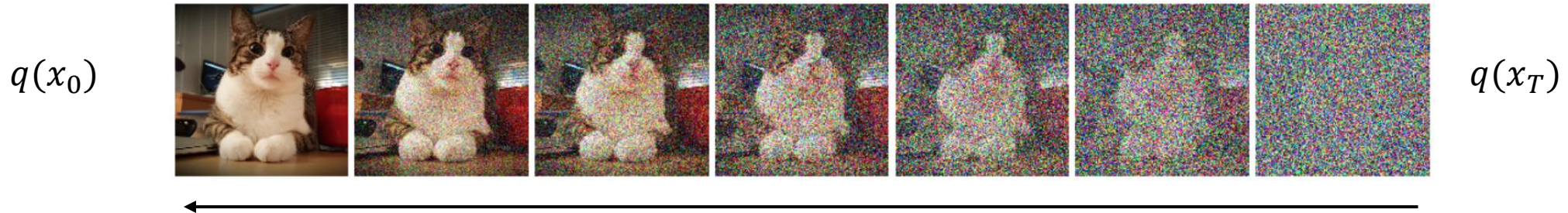
The reverse process $x_T, \dots, x_t, \dots x_0$ can be described by the following SDE

Reverse SDE: $dx_t = -\left[\frac{1}{2}\beta(t)x_t - \beta(t)\nabla_{x_t} \log q(x_t)\right] dt + \sqrt{\beta(t)} dW_t$

Score function

If we simulate this reverse process, we can generate data distribution from noise x_T

Score Matching the Reverse Process



Forward SDE: $dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)} dW_t$

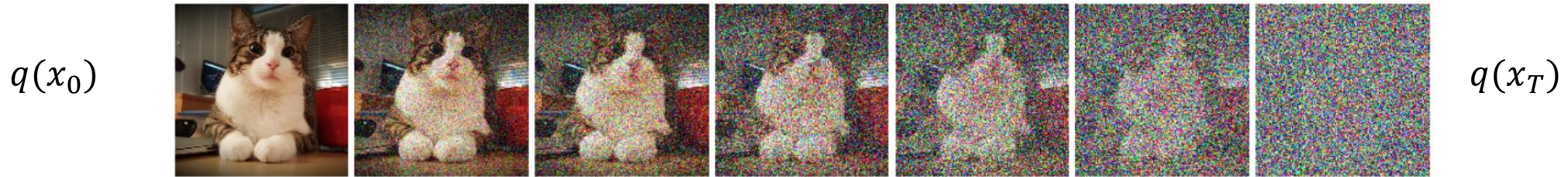
Reverse SDE: $dx_t = -\left[\frac{1}{2}\beta(t)x_t - \beta(t)\nabla_{x_t} \log q(x_t)\right] dt + \sqrt{\beta(t)} dW_t$

Approximate with
neural network $s_\theta(x_t, t)$

$$\min_{\theta} E_{t \sim U(0, T), x_t \sim q(x_t)} \|s_\theta(x_t, t) - \nabla_{x_t} \log q(x_t)\|^2$$

Issue: we don't have close-form formula for $\nabla_{x_t} \log q(x_t)$

Score Matching the Reverse Process



Forward SDE: $dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)} dW_t$

can derive

$$q(x_t | x_0) = \mathcal{N}(\gamma_t x_0, \sigma_t^2 \mathbf{I})$$

$$\gamma_t = e^{-\frac{1}{2} \int_0^t \beta(s) ds}, \sigma_t^2 = 1 - e^{-\int_0^t \beta(s) ds}$$

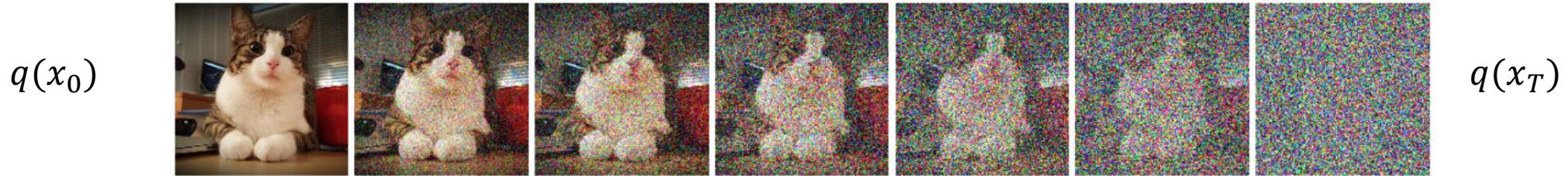
Modified matching objective: $\min_{\theta} \mathbb{E}_{t \sim U(0, T), x_0 \sim q(x_0), x_t \sim q(x_t | x_0)} \|s_{\theta}(x_t, t) - \nabla_{x_t} \log q(x_t | x_0)\|^2$

Sample x_t : $x_t = \gamma_t x_0 + \sigma_t \epsilon \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$

Simplify the score: $\nabla_{x_t} \log q(x_t | x_0) = -\nabla_{x_t} \frac{(x_t - \gamma_t x_0)^2}{2\sigma_t^2} = -\frac{x_t - \gamma_t x_0}{\sigma_t^2} = -\frac{\epsilon}{\sigma_t}$

Parameterize: $s_{\theta}(x_t, t) = -\frac{\epsilon_{\theta}(x_t, t)}{\sigma_t}$ \Rightarrow $\min_{\theta} \mathbb{E}_{t \sim U(0, T), x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(0, \mathbf{I})} \frac{1}{\sigma_t^2} \|\epsilon_{\theta}(x_t, t) - \epsilon\|^2$

Training Diffusion Model



Forward SDE:
$$dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)} dW_t$$

Sample $t \sim U(0, T)$

$$x_t = \gamma_t x_0 + \sigma_t \epsilon \quad \epsilon \sim \mathcal{N}(0, I)$$

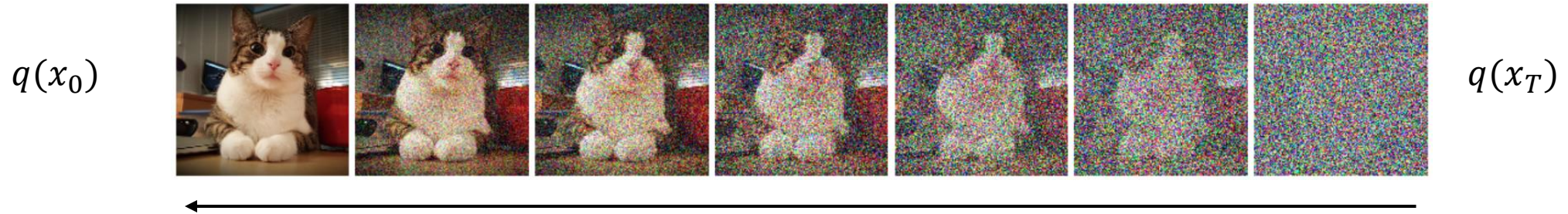
$$q(x_t | x_0) = \mathcal{N}(\gamma_t x_0, \sigma_t^2 I)$$
$$\gamma_t = e^{-\frac{1}{2} \int_0^t \beta(s) ds}, \sigma_t^2 = 1 - e^{-\int_0^t \beta(s) ds}$$

Update $\epsilon_\theta(x_t, t)$ to minimize $\frac{\lambda_t}{\sigma_t^2} \|\epsilon_\theta(x_t, t) - \epsilon\|^2$

Time step weighting, usually set $\lambda_t = \sigma_t^2$

Intuition: we are predicting the noise ϵ !

Generation Process in Diffusion Model



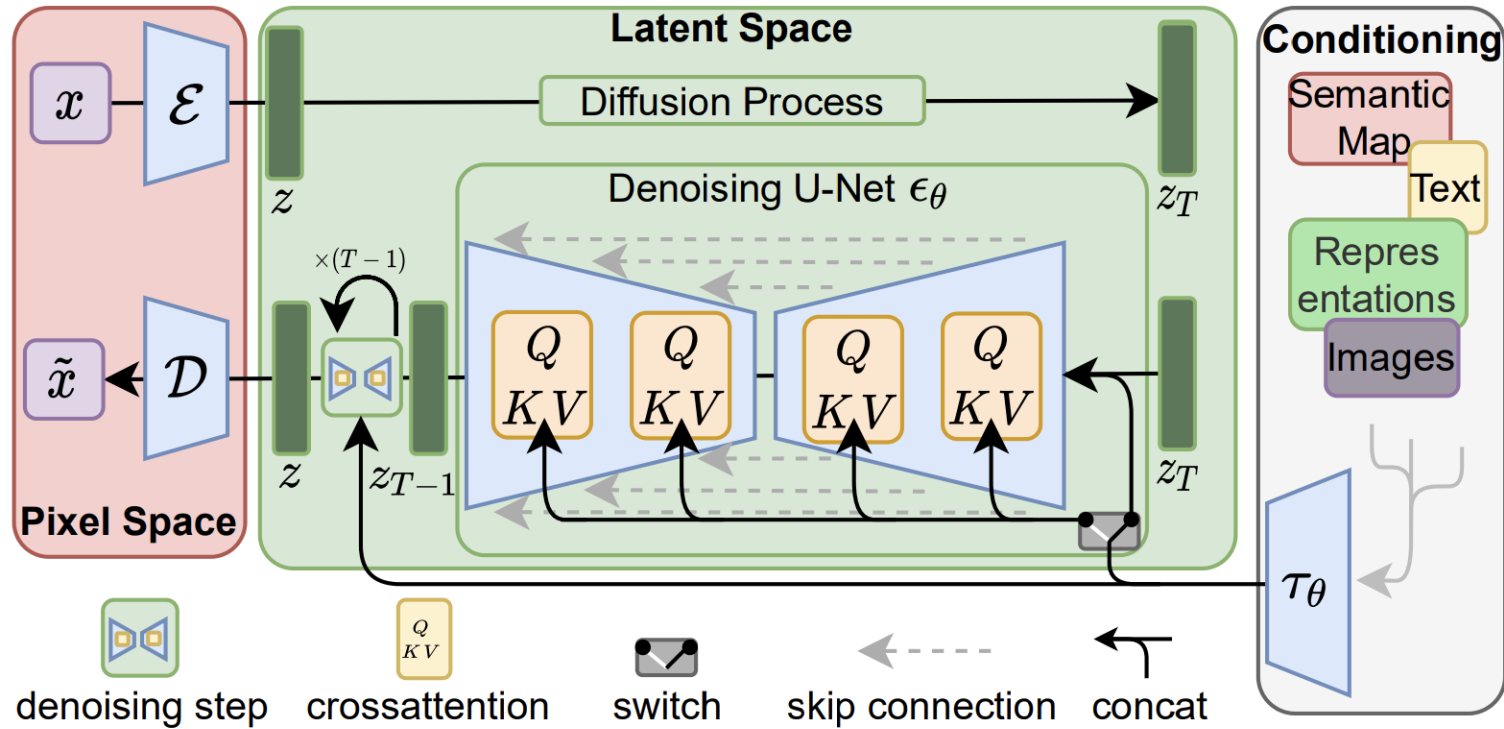
Forward SDE: $dx_t = -\frac{1}{2}\beta(t)x_t dt + \sqrt{\beta(t)} dW_t$

Generation process: $dx_t = -\left[\frac{1}{2}\beta(t)x_t - \frac{\beta(t)\epsilon_\theta(x_t, t)}{\sigma_t}\right] dt + \sqrt{\beta(t)} dW_t$

Intuition take a small step in reverse direction of predicted noise

Different ways to numerically run the generation process

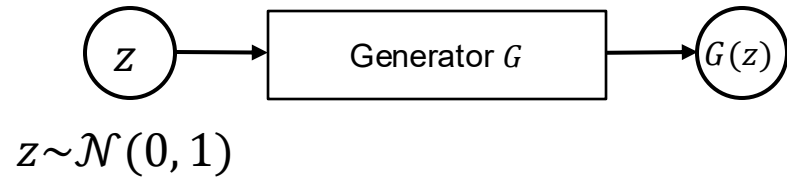
Latent Space Diffusion Models



Running diffusion process in latent space
Decode back to higher resolution pixel space

Comparing GAN and Diffusion Models

GAN



Generate output by single step G

Diffusion
Models



Learning iterative refinement
instead of single step generation

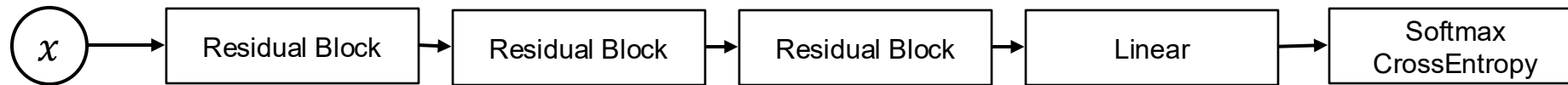
Outline

Generative adversarial training (GAN)

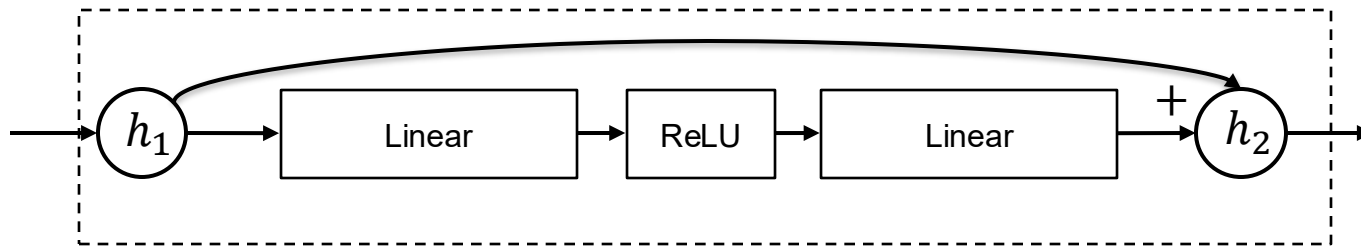
Diffusion models

Deep learning is modular in nature

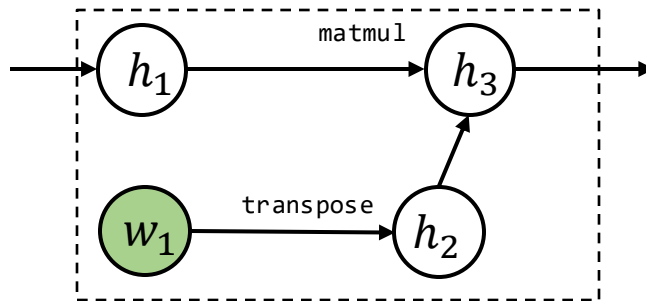
Multi-layer Residual Net



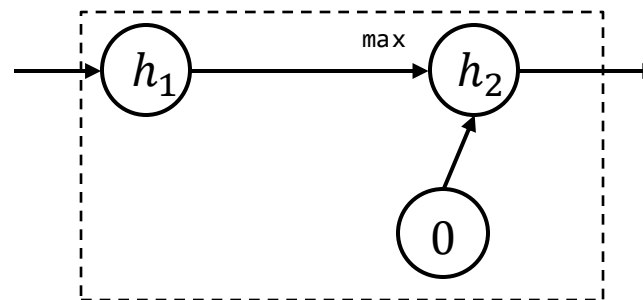
Residual block



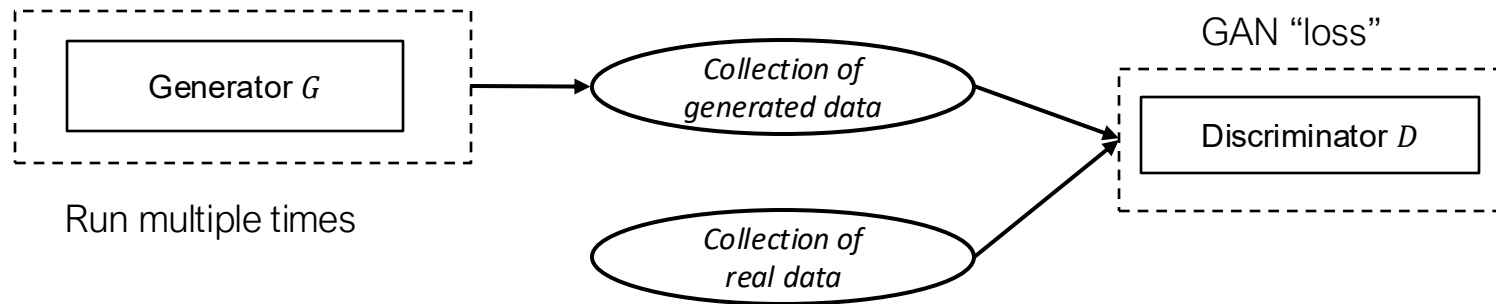
Linear



ReLU



Use GAN as a compositional module



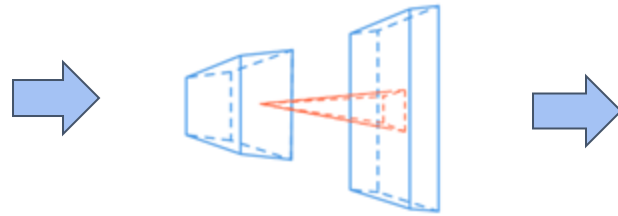
GAN is not exactly like a loss function, as it involves an iterative update recipe. But we can compose it with other neural network modules in a similar way like loss function.

Use GAN "loss" whenever we want a collection of data to "look like" another collection

DCGAN: Deep convolutional generative adversarial networks

Generator

z
 $z \sim \mathcal{N}(0, 1)$

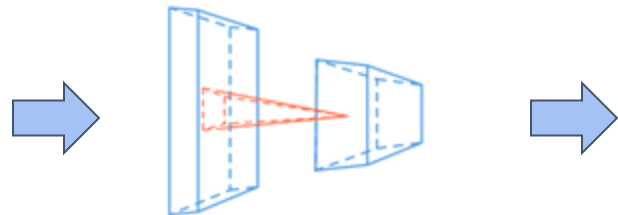


Fake image

Convolutional units with
Conv2dTranpose

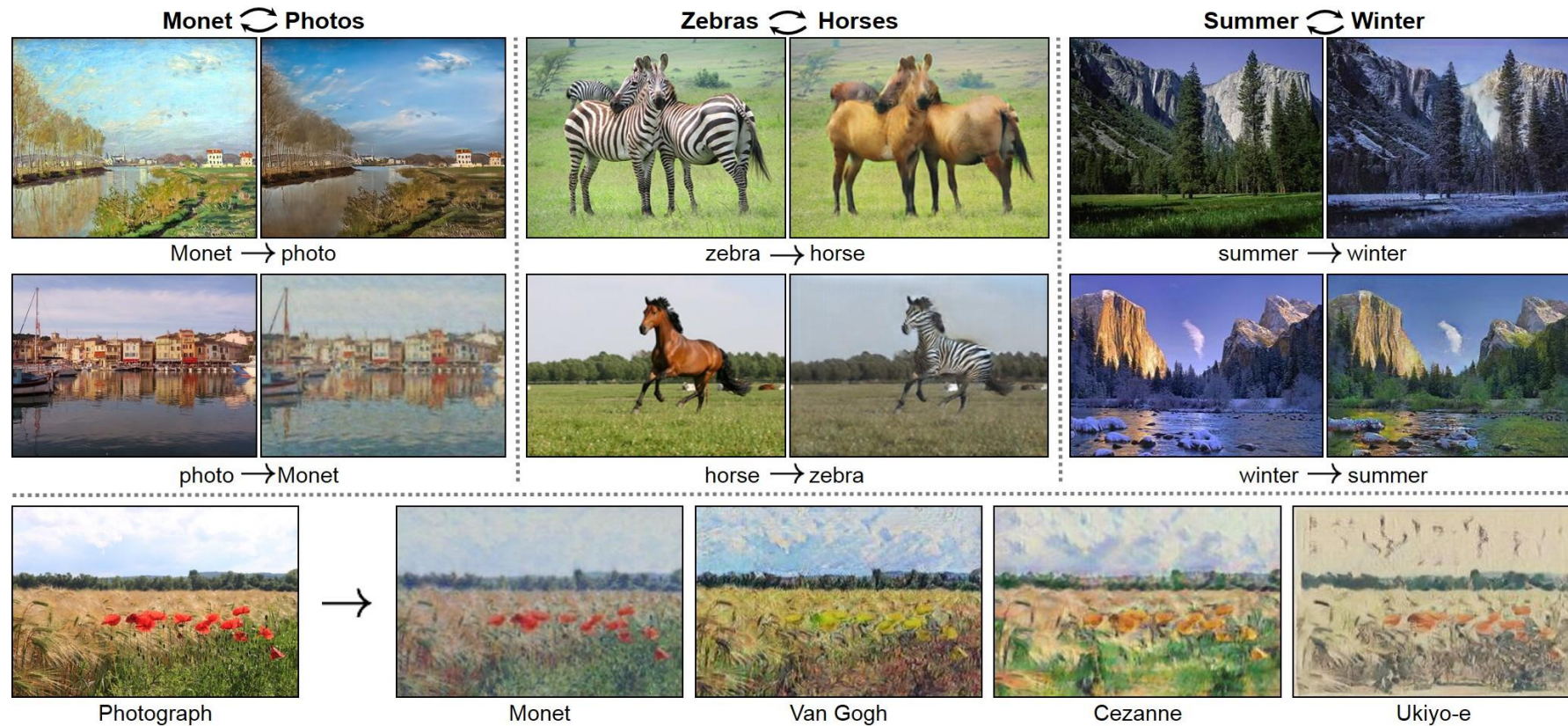
Discriminator

Fake image



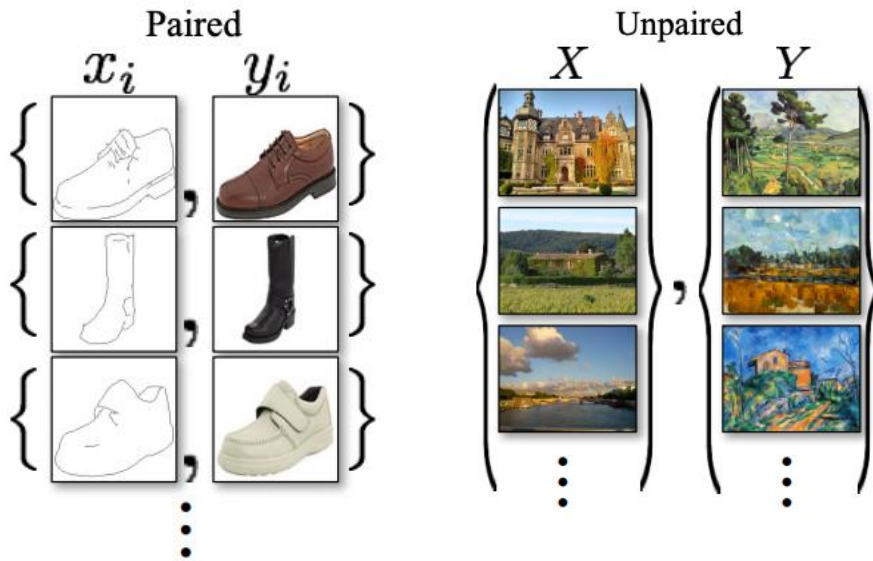
Prediction

CycleGAN: Image to image translation



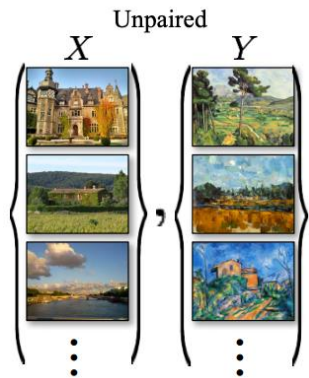
Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, ICCV 2017

CycleGAN: Structure

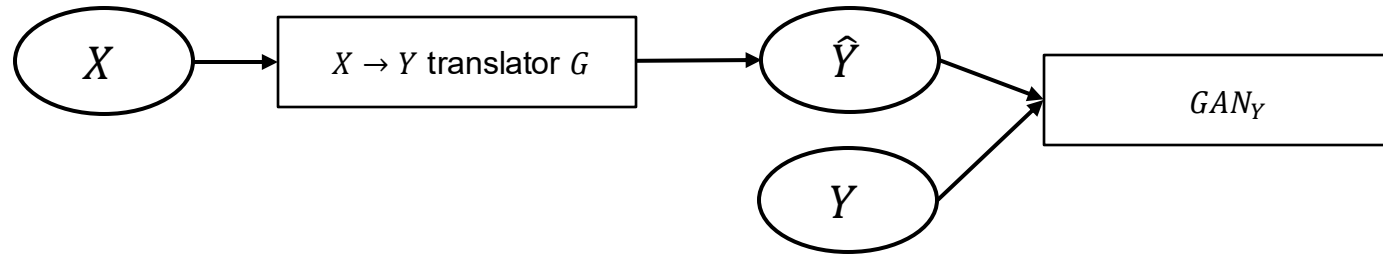


The goal of CycleGAN is to learn bi-directional translator between two **unpaired** collections of data

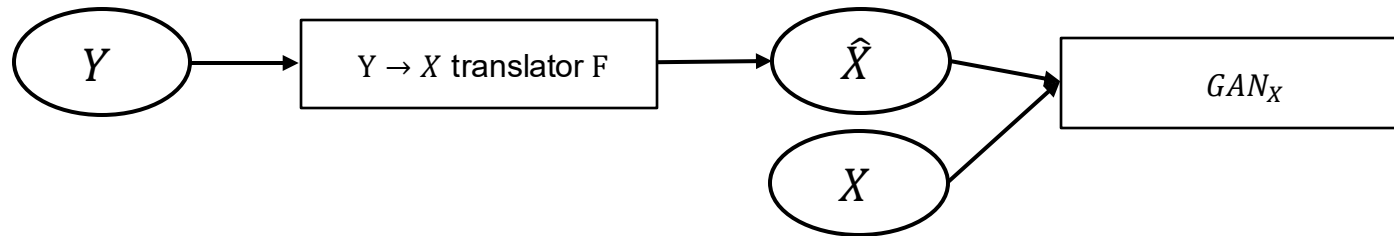
CycleGAN: The model structure



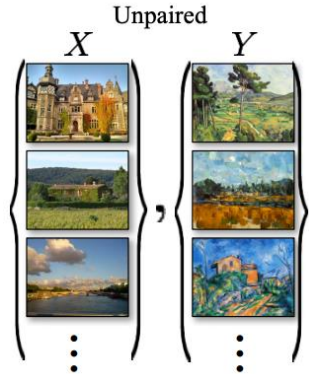
Collection of \hat{Y} should look like collection of Y



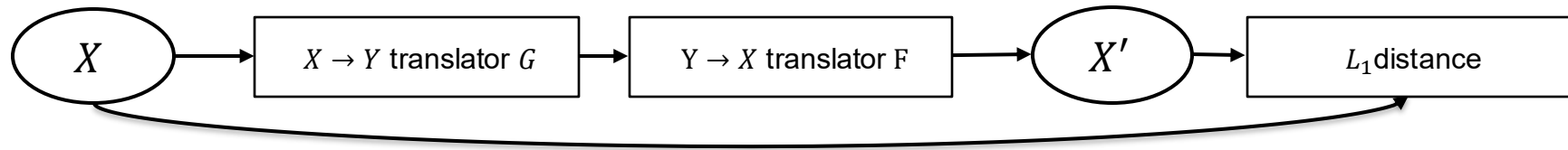
Collection of \hat{X} should look like collection of X



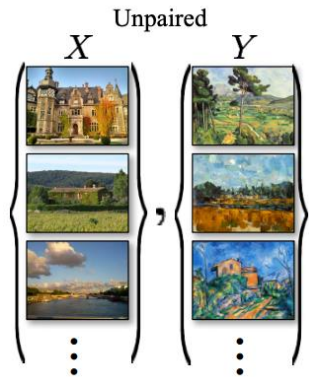
CycleGAN: The model structure



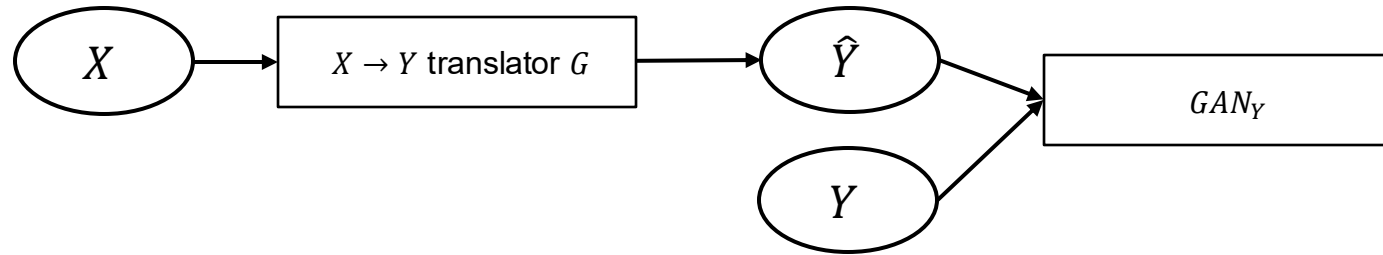
Cycle consistency: map forward and back should map back to the original image



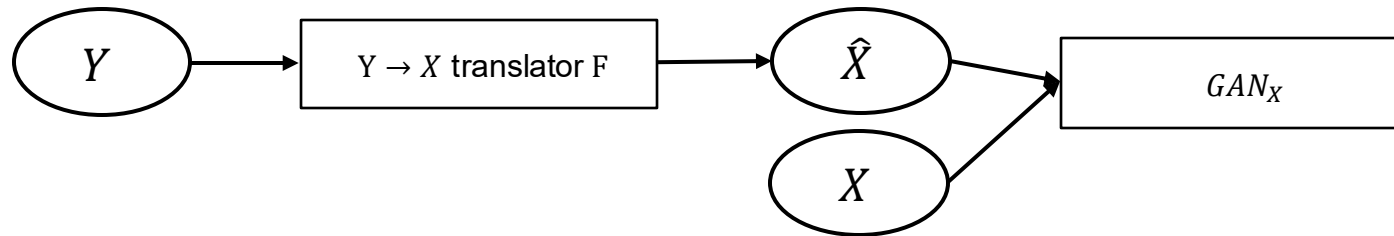
CycleGAN: Putting it together



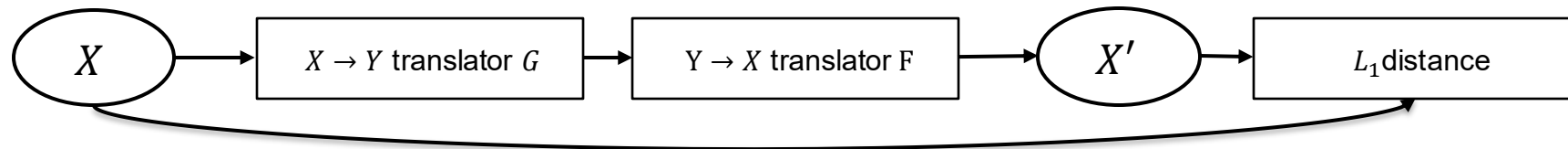
Collection of \hat{Y} should look like collection of Y



Collection of \hat{X} should look like collection of X



Cycle consistency: map forward and back should map back to the original image



Some results

