# Outline

- Deep Generative Models
  - Generative adversarial learning

- Paper presentation:
  - Devanshi Garg, Shrenik Jain: "RHO-1: Not All Tokens Are What You Need"

# Generative modeling

- 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:



CC-LAPGAN: Dog

2009                    2015                              2018

# Generative modeling

- 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:



2014

2015

2016

2017

2018

# Generative modeling



Midjourney, 2025

# Generative modeling
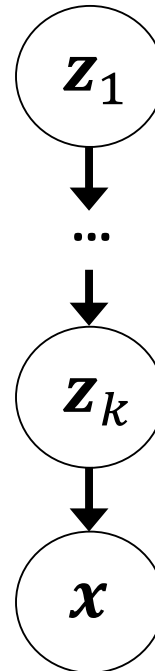
- In ge                                                      bution, such
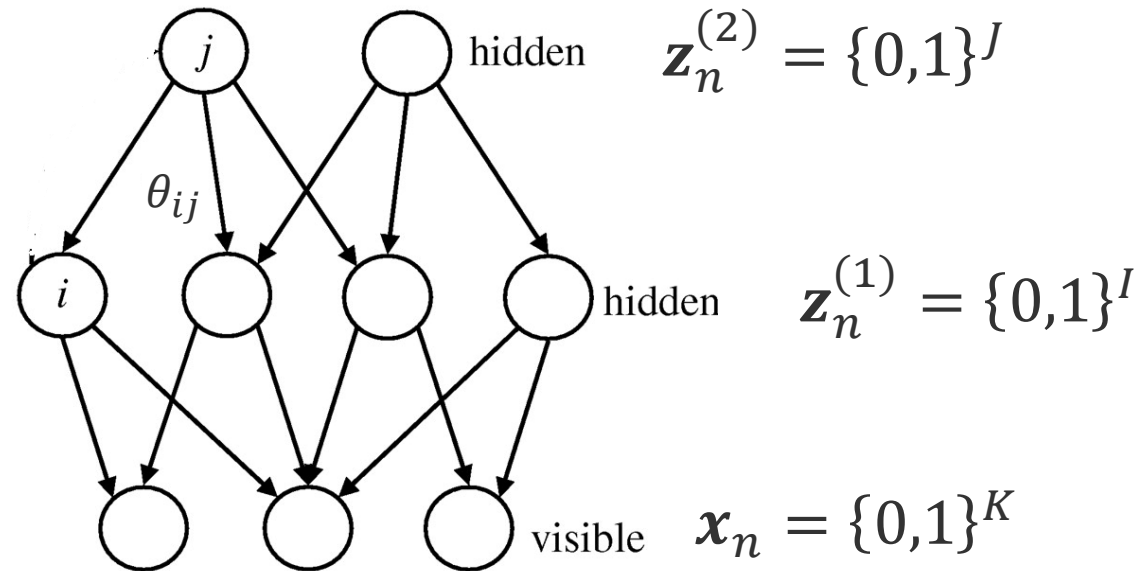
- One

- This



Google Veo2, 12/2024

# Deep generative models

- Define probabilistic distributions over a set of variables
- "Deep" means multiple layers of hidden variables!

$$z_1$$

$$\cdots$$

$$z_k$$

$$x$$

# Early forms of deep generative models

- Hierarchical Bayesian models
  - Sigmoid brief nets [Neal 1992]



$$\mathbf{z}_n^{(2)} = \{0,1\}^J$$

$$\mathbf{z}_n^{(1)} = \{0,1\}^I$$
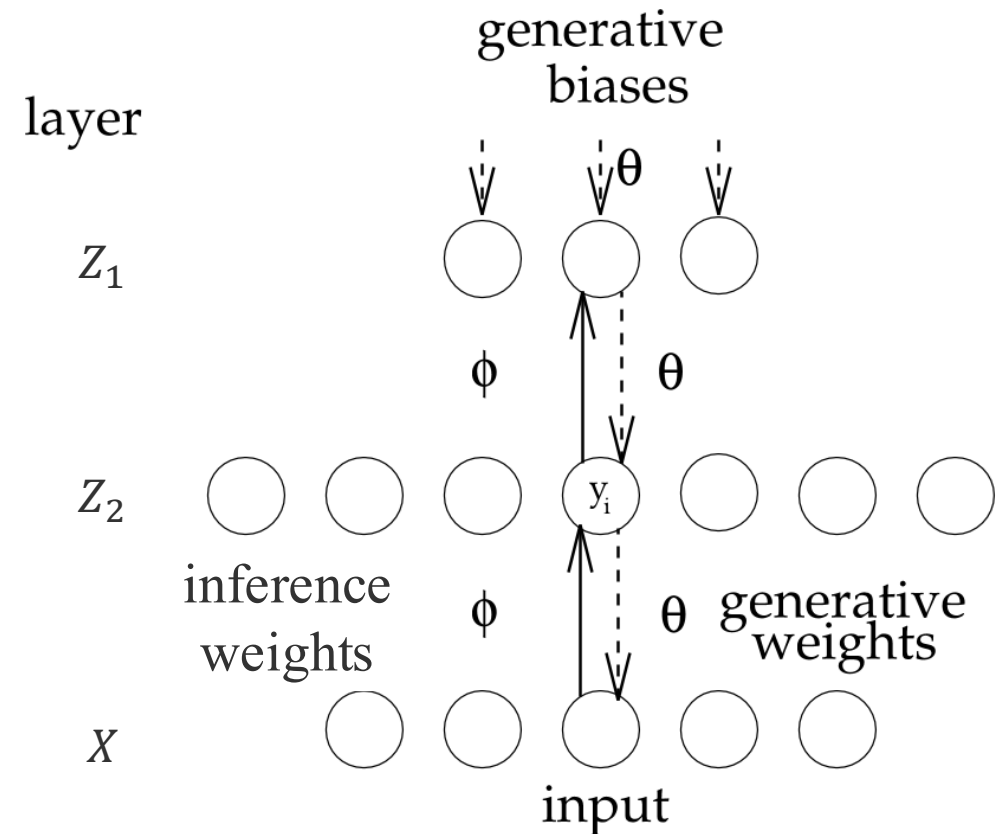
$$\mathbf{x}_n = \{0,1\}^K$$

$$p\left(x_{kn} = 1 \middle| \boldsymbol{\theta}_k, \mathbf{z}_n^{(1)}\right) = \sigma\left(\boldsymbol{\theta}_k^T \mathbf{z}_n^{(1)}\right)$$

$$p\left(z_{in}^{(1)} = 1 \middle| \boldsymbol{\theta}_i, \mathbf{z}_n^{(2)}\right) = \sigma\left(\boldsymbol{\theta}_i^T \mathbf{z}_n^{(2)}\right)$$

# Early forms of deep generative models

- Hierarchical Bayesian models
  - Sigmoid brief nets [Neal 1992]

- Neural network models
  - Helmholtz machines [Dayan et al.,1995]



[Dayan et al. 1995]

# Early forms of deep generative models

- Hierarchical Bayesian models
  - Sigmoid brief nets [Neal 1992]

- Neural network models
  - Helmholtz machines [Dayan et al.,1995]
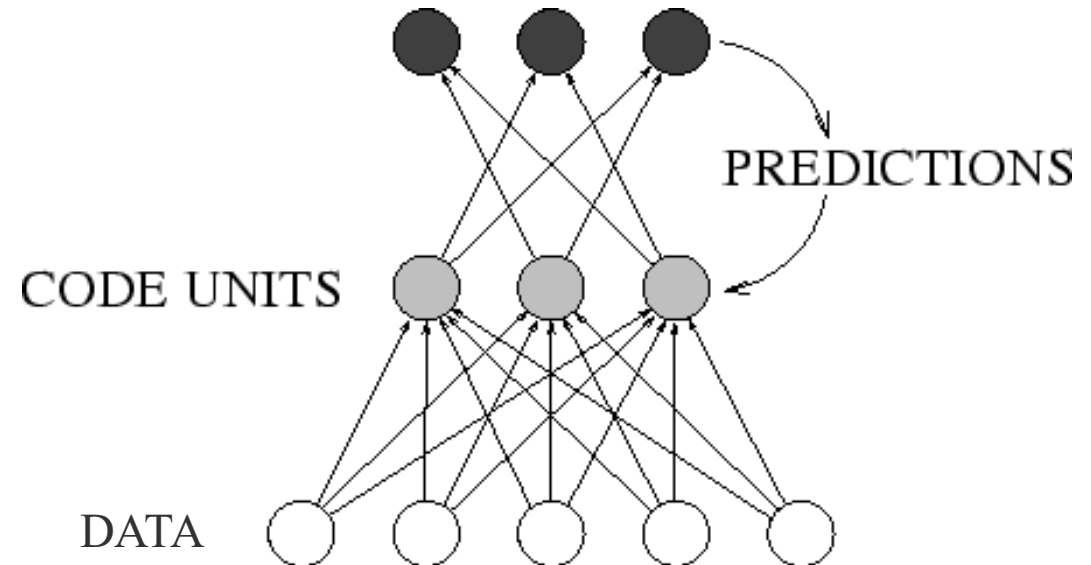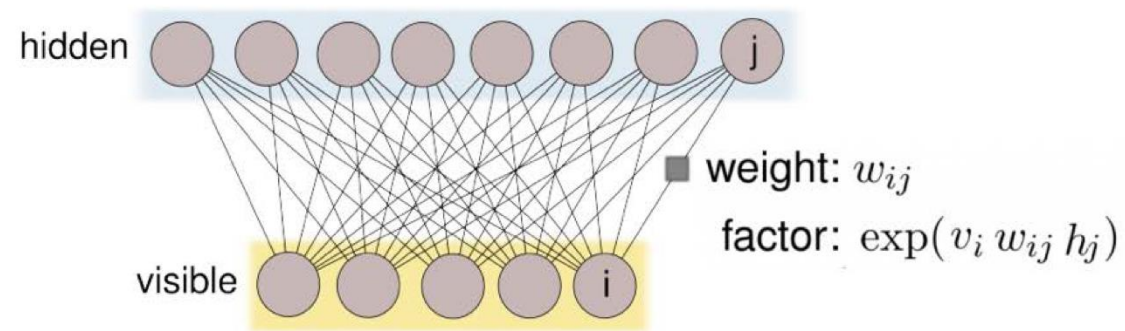  - Predictability minimization [Schmidhuber 1995]

Figure courtesy: Schmidhuber 1996

# Resurgence of deep generative models

- Restricted Boltzmann machines (RBMs) [Smolensky, 1986]
  - Building blocks of deep probabilistic models



weight: $w_{ij}$

factor: $\exp(v_i \, w_{ij} \, h_j)$

# Resurgence of deep generative models
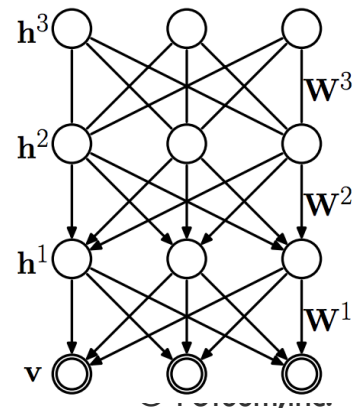
- Restricted Boltzmann machines (RBMs) [Smolensky, 1986]
  - Building blocks of deep probabilistic models

- Deep belief networks (DBNs) [Hinton et al., 2006]
  - Hybrid graphical model
  - Inference in DBNs is problematic due to explaining away

- Deep Boltzmann Machines (DBMs) [Salakhutdinov & Hinton, 2009]
  - Undirected model



Deep Belief Network        Deep Boltzmann Machine

# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]
  / Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]



$q_\phi(z|x)$
inference model

$p_\theta(x|z)$
generative model

Figure courtesy: Kingma & Welling, 2014

# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]

  / Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]

- Generative adversarial networks (GANs) [Goodfellow et al,. 2014]



code     data/gen

$G_\theta$ : generative model  **?**
$D_\phi$ : discriminator

# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]

  / Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]

- Generative adversarial networks (GANs) [Goodfellow et al,. 2014]

- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]

# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]
  / Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
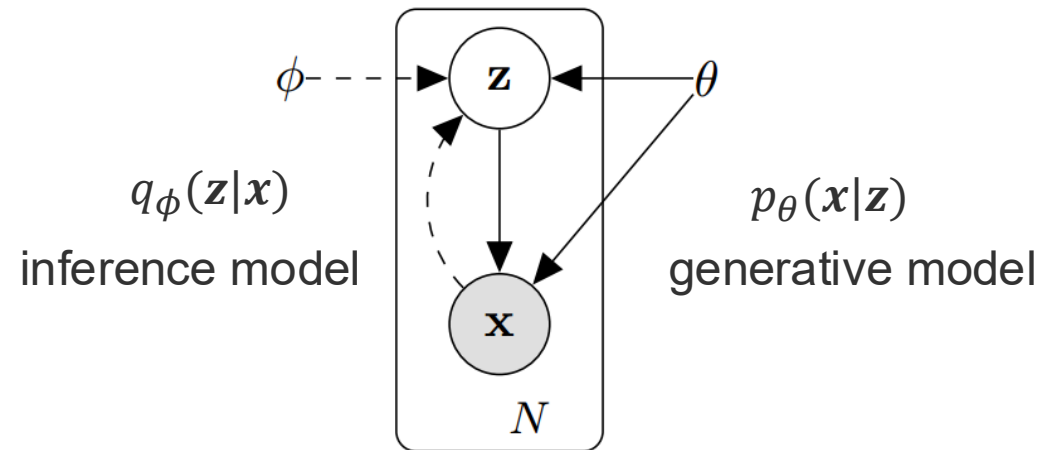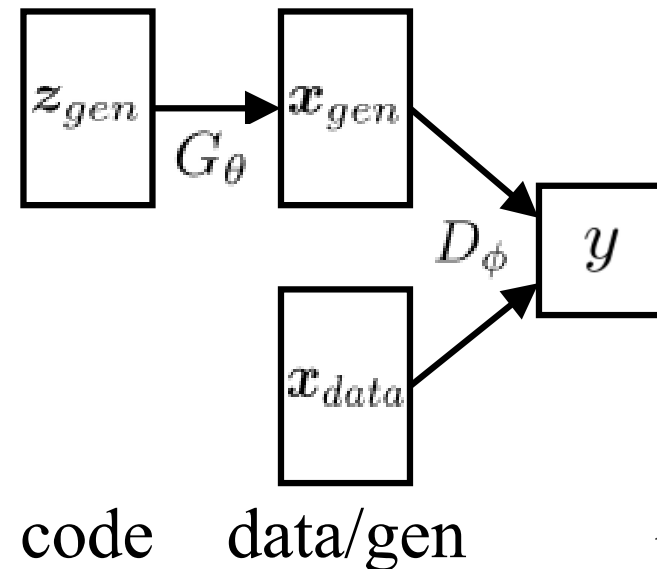- Generative adversarial networks (GANs) [Goodfellow et al,. 2014]
- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]
- Autoregressive neural networks

# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]
  / Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
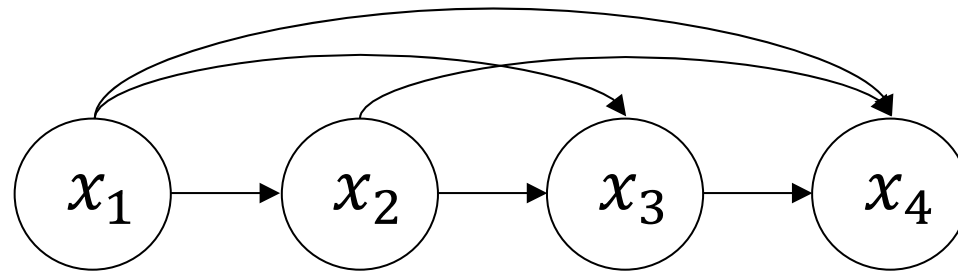- Generative adversarial networks (GANs) [Goodfellow et al,. 2014]
- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]
- Autoregressive neural networks
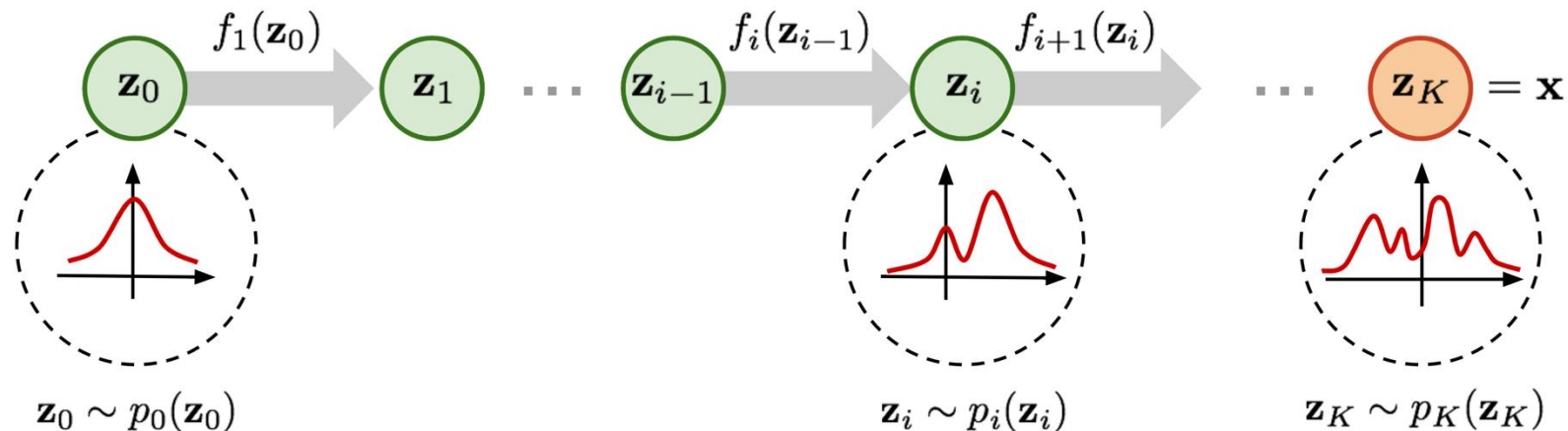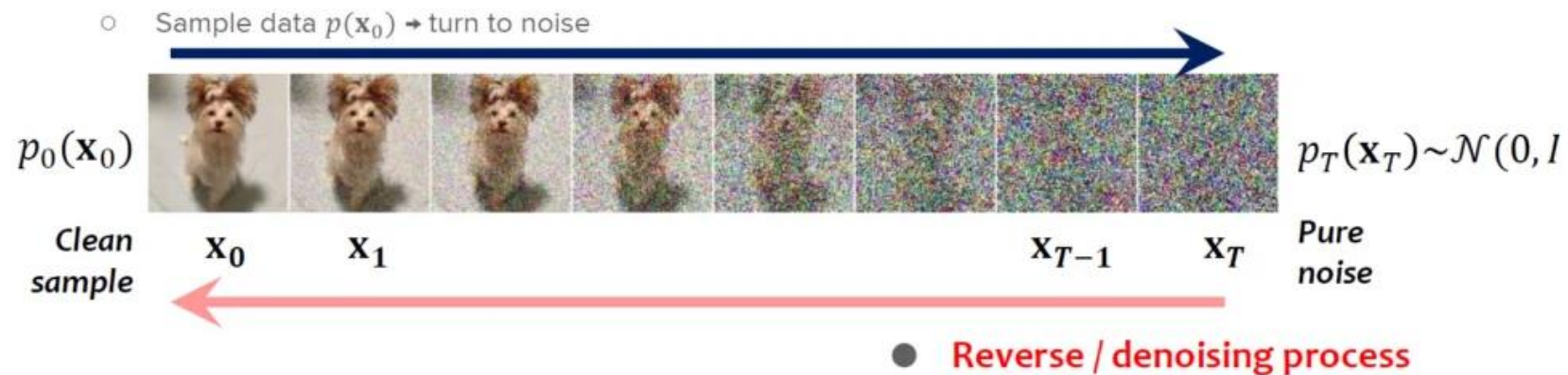- Reversible architectures (flow models)

# Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]

  / Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]

- Generative adversarial networks (GANs) [Goodfellow et al,. 2014]

- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]

- Autoregressive neural networks

- Reversible architectures (flow models)

- Diffusion models



- Sample data $p(x_0)$ ➔ turn to noise

$p_0(\mathbf{x}_0)$       $p_T(\mathbf{x}_T) \sim \mathcal{N}(0, I)$

Clean sample    $\mathbf{x}_0$    $\mathbf{x}_1$       $\mathbf{x}_{T-1}$    $\mathbf{x}_T$    Pure noise

- **Reverse / denoising process**

- Sample noise $p_T(\mathbf{x}_T)$ ➔ turn into data

# Generative Adversarial Networks

# Implicit Generative Models

- **Implicit generative models** implicitly define a probability distribution

- Start by sampling the code vector **z** from a fixed, simple distribution (e.g. spherical Gaussian)

- The generator network computes a differentiable function $G$ mapping **z** to an **x** in data space

sample    | $\mathbf{x} = G(\mathbf{z})$ |

- a stochastic process to simulate data $x$
- Intractable to evaluate likelihood

code vector    | **z** |

# Implicit Generative Models

A 1-dimensional example:

input
distribution

output
distribution

function
computed by
the network

# Implicit Generative Models

# Implicit Generative Models

- 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

# Generative Adversarial Nets (GANs)

- Generative model $x = G_\theta(z), \quad z \sim p(z)$
  - Maps noise variable $z$ to data space $x$
  - Defines an implicit distribution over $x$: $p_{g_\theta}(x)$

- Discriminator $D_\phi(x)$
  - Output the probability that $x$ came from the data rather than the generator



real image

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

G
(generator)

fake image

D
(discriminator)

**1**(Real)

**0**(fake)

**1**(real)

— **Discriminator training**
— **Generator training**

# Generative Adversarial Nets (GANs)

- Learning
  - A minimax game between the generator and the discriminator
  - Train $D$ to maximize the probability of assigning the correct label to both training examples and generated samples
  - Train $G$ to fool the discriminator

$$\max_D \mathcal{L}_D = \mathbb{E}_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})} \left[\log D(\boldsymbol{x})\right] + \mathbb{E}_{\boldsymbol{x} \sim G(\boldsymbol{z}), \boldsymbol{z} \sim p(\boldsymbol{z})} \left[\log(1 - D(\boldsymbol{x}))\right]$$

$$\min_G \mathcal{L}_G = \mathbb{E}_{\boldsymbol{x} \sim G(\boldsymbol{z}), \boldsymbol{z} \sim p(\boldsymbol{z})} \left[\log(1 - D(\boldsymbol{x}))\right].$$

# Generative Adversarial Nets (GANs)

# Generative Adversarial Nets (GANs)

Updating the discriminator:



$D(\mathbf{x})$

update the discriminator
weights using backprop
on the classification objective

$\mathbf{x}$    OR    $\mathbf{x} = G(\mathbf{z})$

real-world
image

generator

$\mathbf{z}$    code vector

# Generative Adversarial Nets (GANs)

Updating the generator:



$D(\mathbf{x})$

backprop the derivatives, but don't modify the discriminator weights

flip the sign of the derivatives

$\mathbf{x} = G(\mathbf{z})$

update the generator weights using backprop

$\mathbf{z}$

# Generative Adversarial Nets (GANs)

Alternating training of the generator and discriminator:

# Optimality of GANs

- Objectives:

$$\max_D \mathcal{L}_D = \mathbb{E}_{\boldsymbol{x} \sim p_{data}(\boldsymbol{x})} \left[ \log D(\boldsymbol{x}) \right] + \mathbb{E}_{\boldsymbol{x} \sim G(\boldsymbol{z}), \boldsymbol{z} \sim p(\boldsymbol{z})} \left[ \log(1 - D(\boldsymbol{x})) \right]$$

$$\min_G \mathcal{L}_G = \mathbb{E}_{\boldsymbol{x} \sim G(\boldsymbol{z}), \boldsymbol{z} \sim p(\boldsymbol{z})} \left[ \log(1 - D(\boldsymbol{x})) \right].$$

- Global optimality: $p_g = p_{data}$

- Proof:

# Optimality of GANs

**Proposition 1.** *For G fixed, the optimal discriminator D is*

$$D_G^*(\boldsymbol{x}) = \frac{p_{data}(\boldsymbol{x})}{p_{data}(\boldsymbol{x}) + p_g(\boldsymbol{x})} \tag{2}$$

[Goodfellow et al., 2014]

# Optimality of GANs

**Proposition 1.** *For G fixed, the optimal discriminator D is*

$$D_G^*(\boldsymbol{x}) = \frac{p_{data}(\boldsymbol{x})}{p_{data}(\boldsymbol{x}) + p_g(\boldsymbol{x})} \tag{2}$$

*Proof.* The training criterion for the discriminator D, given any generator $G$, is to maximize the quantity $V(G, D)$

$$V(G, D) = \int_{\boldsymbol{x}} p_{\text{data}}(\boldsymbol{x}) \log(D(\boldsymbol{x})) dx + \int_{\boldsymbol{z}} p_{\boldsymbol{z}}(\boldsymbol{z}) \log(1 - D(g(\boldsymbol{z}))) dz$$

$$= \int_{\boldsymbol{x}} p_{\text{data}}(\boldsymbol{x}) \log(D(\boldsymbol{x})) + p_g(\boldsymbol{x}) \log(1 - D(\boldsymbol{x})) dx \tag{3}$$

For any $(a, b) \in \mathbb{R}^2 \setminus \{0, 0\}$, the function $y \to a \log(y) + b \log(1 - y)$ achieves its maximum in $[0, 1]$ at $\frac{a}{a+b}$.

[Goodfellow et al., 2014]

33

# Optimality of GANs

- The minimax game can now be reformulated as

$$C(G) = \max_{D} V(G, D)$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} [\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}} [\log(1 - D_G^*(G(\boldsymbol{z})))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} [\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{x} \sim p_g} [\log(1 - D_G^*(\boldsymbol{x}))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \left[ \log \frac{p_{\text{data}}(\boldsymbol{x})}{P_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})} \right] + \mathbb{E}_{\boldsymbol{x} \sim p_g} \left[ \log \frac{p_g(\boldsymbol{x})}{p_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})} \right]$$

[Goodfellow et al., 2014]

# Optimality of GANs

- The minimax game can now be reformulated as

$$C(G) = \max_D V(G, D)$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}[\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}}[\log(1 - D_G^*(G(\boldsymbol{z})))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}[\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{x} \sim p_g}[\log(1 - D_G^*(\boldsymbol{x}))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}\left[\log \frac{p_{\text{data}}(\boldsymbol{x})}{P_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}\right] + \mathbb{E}_{\boldsymbol{x} \sim p_g}\left[\log \frac{p_g(\boldsymbol{x})}{p_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}\right]$$

**Theorem 1.** *The global minimum of the virtual training criterion $C(G)$ is achieved if and only if $p_g = p_{data}$. At that point, $C(G)$ achieves the value $-\log 4$.*

# Optimality of GANs

- The minimax game can now be reformulated as

$$C(G) = \max_D V(G, D)$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}[\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}}[\log(1 - D_G^*(G(\boldsymbol{z})))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}[\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{x} \sim p_g}[\log(1 - D_G^*(\boldsymbol{x}))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}\left[\log \frac{p_{\text{data}}(\boldsymbol{x})}{P_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}\right] + \mathbb{E}_{\boldsymbol{x} \sim p_g}\left[\log \frac{p_g(\boldsymbol{x})}{p_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}\right]$$

**Theorem 1.** *The global minimum of the virtual training criterion $C(G)$ is achieved if and only if $p_g = p_{data}$. At that point, $C(G)$ achieves the value $-\log 4$.*

$$C(G) = -\log(4) + KL\left(p_{\text{data}} \left\| \frac{p_{\text{data}} + p_g}{2}\right.\right) + KL\left(p_g \left\| \frac{p_{\text{data}} + p_g}{2}\right.\right)$$

$$= -\log(4) + 2 \cdot JSD\left(p_{\text{data}} \| p_g\right) \quad \text{Jensen-Shannon Divergence}$$

[Goodfellow et al., 2014]

## Optimality of GANs

- The minimax game can now be reformulated as

$$C(G) = \max_D V(G, D)$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}[\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}}[\log(1 - D_G^*(G(\boldsymbol{z})))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}[\log D_G^*(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{x} \sim p_g}[\log(1 - D_G^*(\boldsymbol{x}))]$$

$$= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}}\left[\log \frac{p_{\text{data}}(\boldsymbol{x})}{P_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}\right] + \mathbb{E}_{\boldsymbol{x} \sim p_g}\left[\log \frac{p_g(\boldsymbol{x})}{p_{\text{data}}(\boldsymbol{x}) + p_g(\boldsymbol{x})}\right]$$

**Theorem 1.** *The global minimum of the virtual training criterion $C(G)$ is achieved if and only if $p_g = p_{data}$. At that point, $C(G)$ achieves the value $-\log 4$.*

$$C(G) = -\log(4) + KL\left(p_{\text{data}} \left\| \frac{p_{\text{data}} + p_g}{2}\right.\right) + KL\left(p_g \left\| \frac{p_{\text{data}} + p_g}{2}\right.\right)$$

$$= -\log(4) + 2 \cdot JSD\left(p_{\text{data}} \| p_g\right) \quad \text{Jensen-Shannon Divergence}$$

[Goodfellow et al., 2014]

# Wasserstein GAN (WGAN)

- If our data are on a low-dimensional manifold of a high dimensional space, the model's manifold and the true data manifold can have a negligible intersection in practice
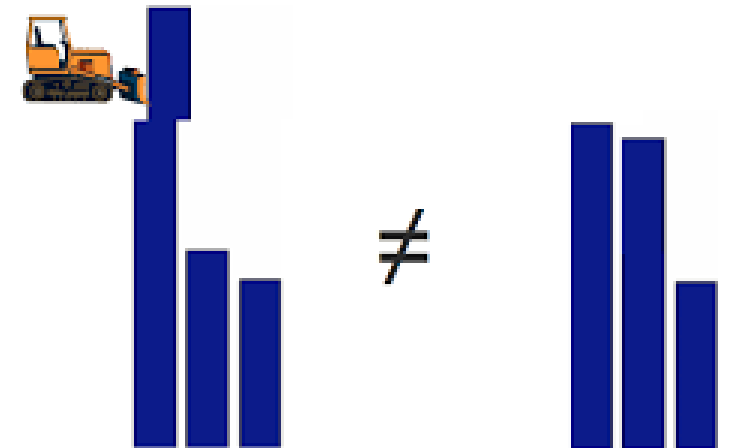
[Arjovsky et al., 2017]

# Wasserstein GAN (WGAN)

- If our data are on a low-dimensional manifold of a high dimensional space, the model's manifold and the true data manifold can have a negligible intersection in practice

- The loss function and gradients may not be continuous and well behaved

[Arjovsky et al., 2017]    Slide adapted from bhiksha

# Wasserstein GAN (WGAN)

- If our data are on a <span style="color:red">low-dimensional</span> manifold of a high dimensional space, the model's manifold and the true data manifold can have a <span style="color:red">negligible intersection in practice</span>

- The loss function and gradients may not be continuous and well behaved

- The <span style="color:red">Wasserstein Distance</span> is well defined
  - Earth Mover's Distance
  - Minimum transportation cost for making one pile
    of dirt in the shape of one probability distribution
    to the shape of the other distribution



[Arjovsky et al., 2017]     Slide adapted from bhiksha

# Wasserstein GAN (WGAN)

- Objective

$$W(p_{data}, p_g) = \frac{1}{K} \sup_{||D||_L \leq K} \mathrm{E}_{x \sim p_{data}}[D(x)] - \mathrm{E}_{x \sim p_g}[D(x)]$$

- $||D||_L \leq K$ : K- Lipschitz continuous
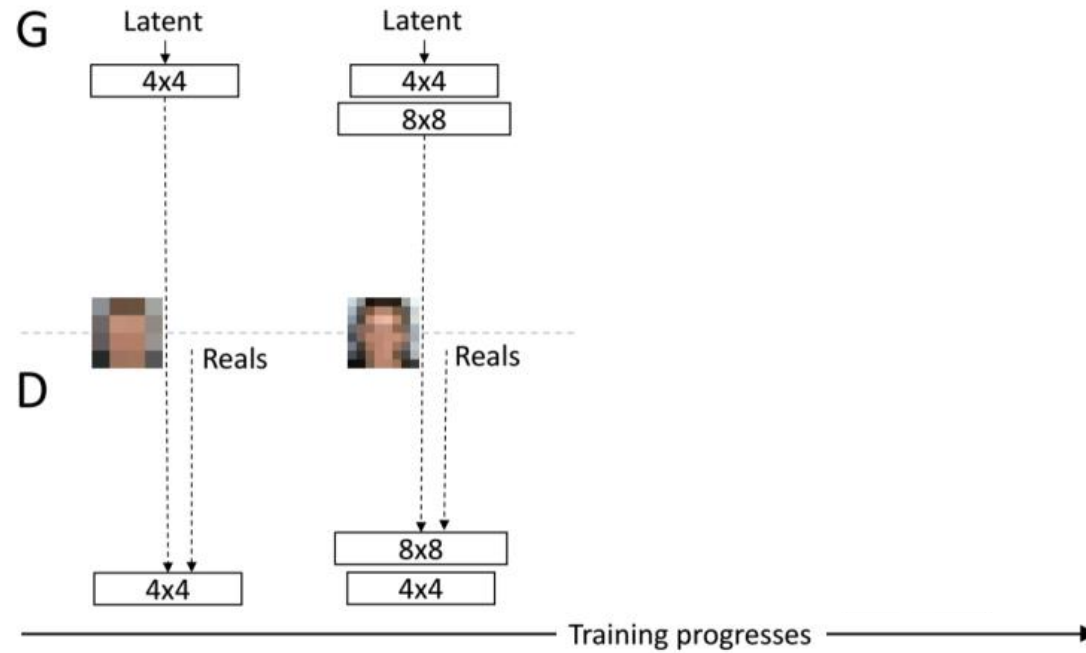- Use gradient-clipping to ensure $D$ has the Lipschitz continuity
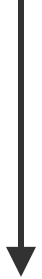
# Progressive GAN

Low resolution images

[Karras et al., 2018]

# Progressive GAN
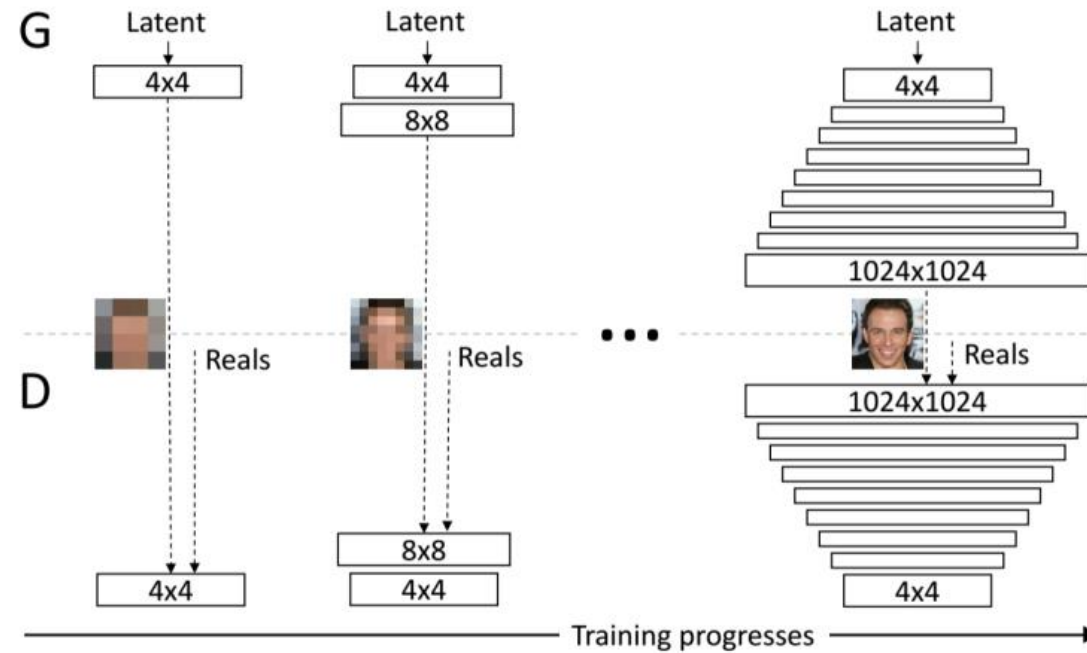
Low resolution images

add in
additional
layers



Training progresses

[Karras et al., 2018]

# Progressive GAN

Low resolution images

add in
additional
layers

High resolution images

[Karras et al., 2018]

# BigGAN

[Brock et al., 2018]

# BigGAN

- GANs benefit dramatically from <span style="color:red">scaling</span>

[Brock et al., 2018]

# BigGAN

- GANs benefit dramatically from <span style="color:red">scaling</span>

- 2x – 4x more parameters

- 8x larger batch size

- Simple architecture changes that improve scalability

[Brock et al., 2018]

# BigGAN

- GANs benefit dramatically from scaling

- 2x – 4x more parameters

- 8x larger batch size

- Simple architecture changes that improve scalability



[Brock et al., 2018]

# BigGAN

- GANs benefit dramatically from scaling
- 2x — 4x more parameters
- 8x
- Sim



[Brock et al., 2018]

# Key Takeaways

- Deep Generative Models: brief history

- GANs:
  - Implicit generative model
  - Minimax formulation
  - Wasserstein GAN

# Questions?