

Learning to Generate Johannes Vermeer's 'The music lesson'. Left: Tim Jenison's version (Tim's Vermeer, 2013) Right: Original (1662 – 1665)



Part 2 - Generative Adversarial Networks

Aykut Erdem

Computer Vision Lab, Hacettepe University



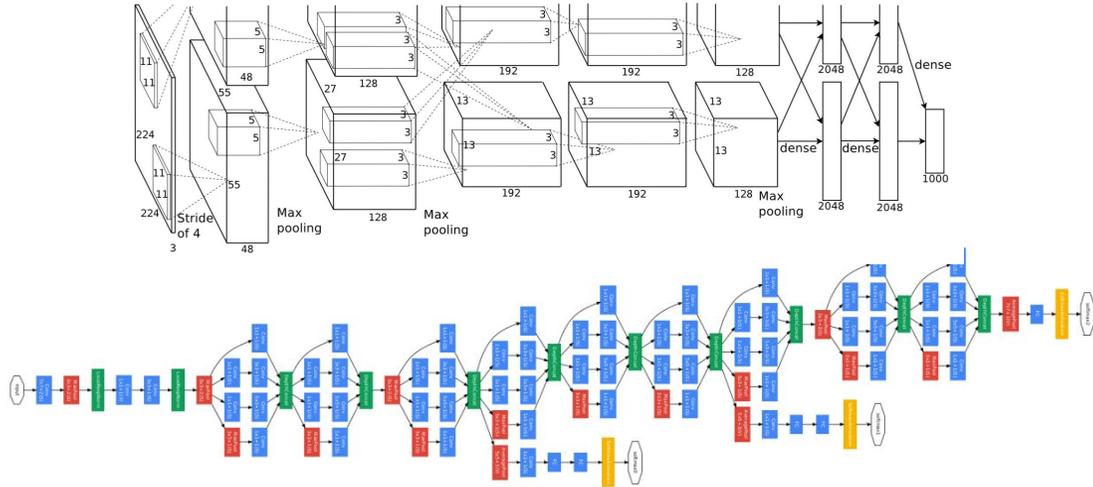
HACETTEPE
UNIVERSITY
COMPUTER
VISION LAB



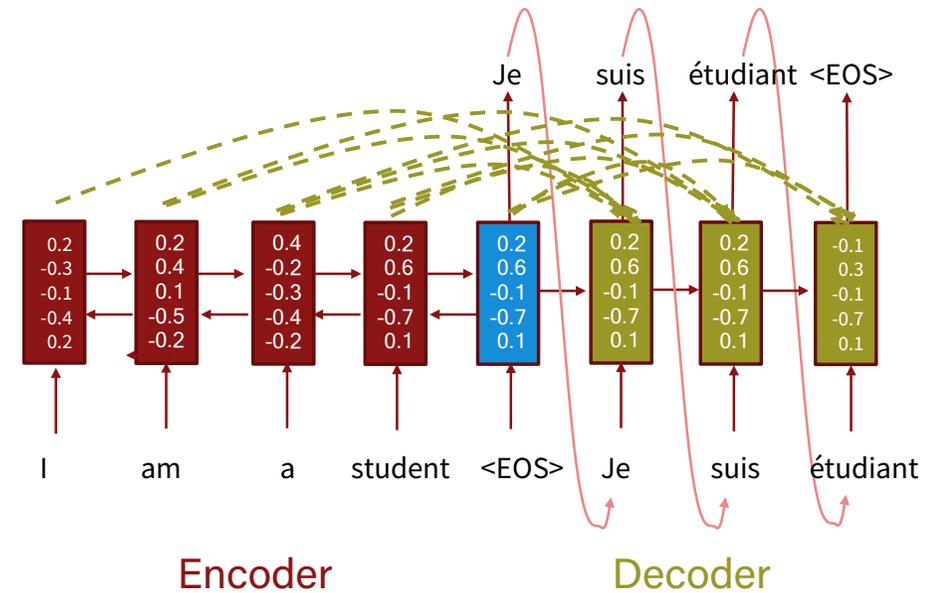
HACETTEPE
UNIVERSITY

Deep Supervised Learning: A Success Story

- Obtain lots of input-output examples
- Train a deep neural network



Deep CNN



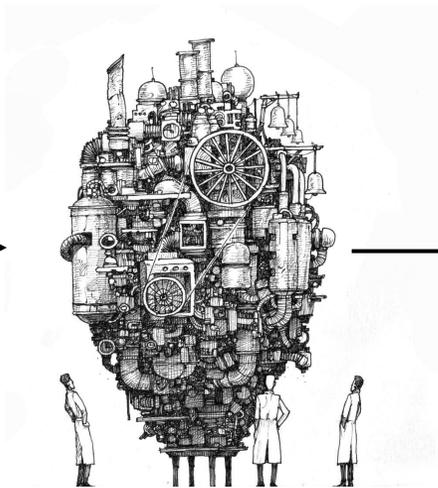
RNN with attention

- Achieve superior results

Discriminative vs. Generative Models

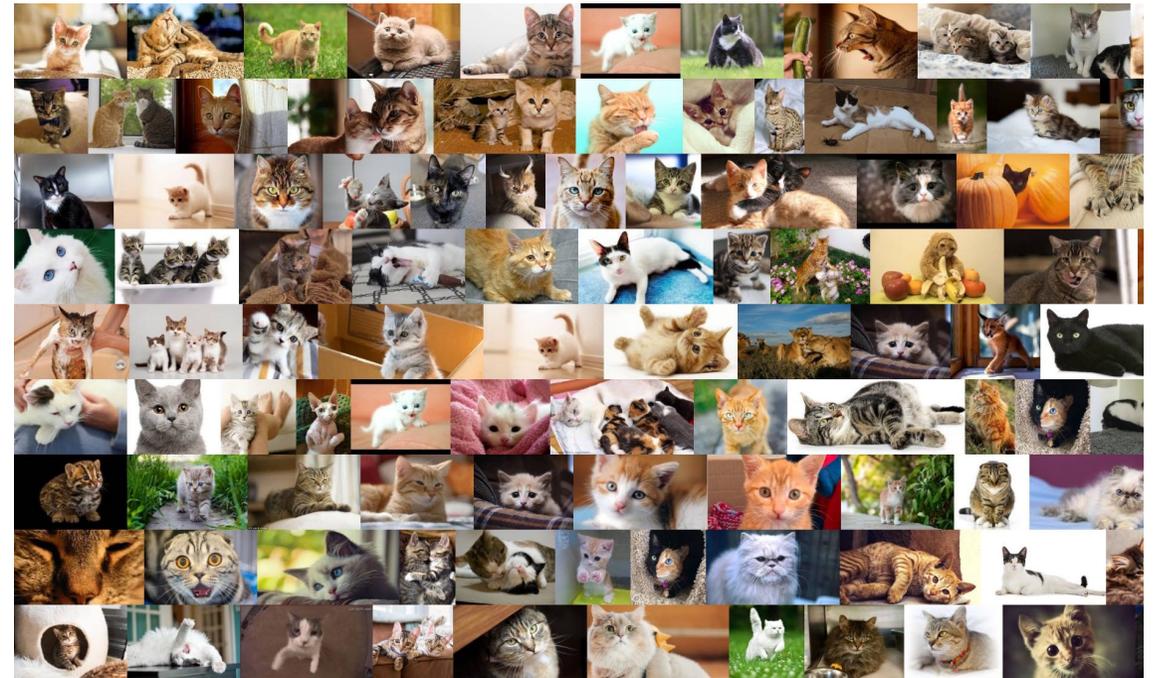
$$p(y|x)$$

$$p(x|y)$$



"Cat"

Discriminative models



Generative models

Why study deep generative models?

- Go beyond associating inputs to outputs
- Understand high-dimensional, complex probability distributions
- Discover the “true” structure of the data
 - Detect surprising events in the world (*anomaly detection*)
 - Missing Data (*semi-supervised learning*)
 - Generate models for planning (*model-based reinforcement learning*)

Why study Generative Adversarial Networks?

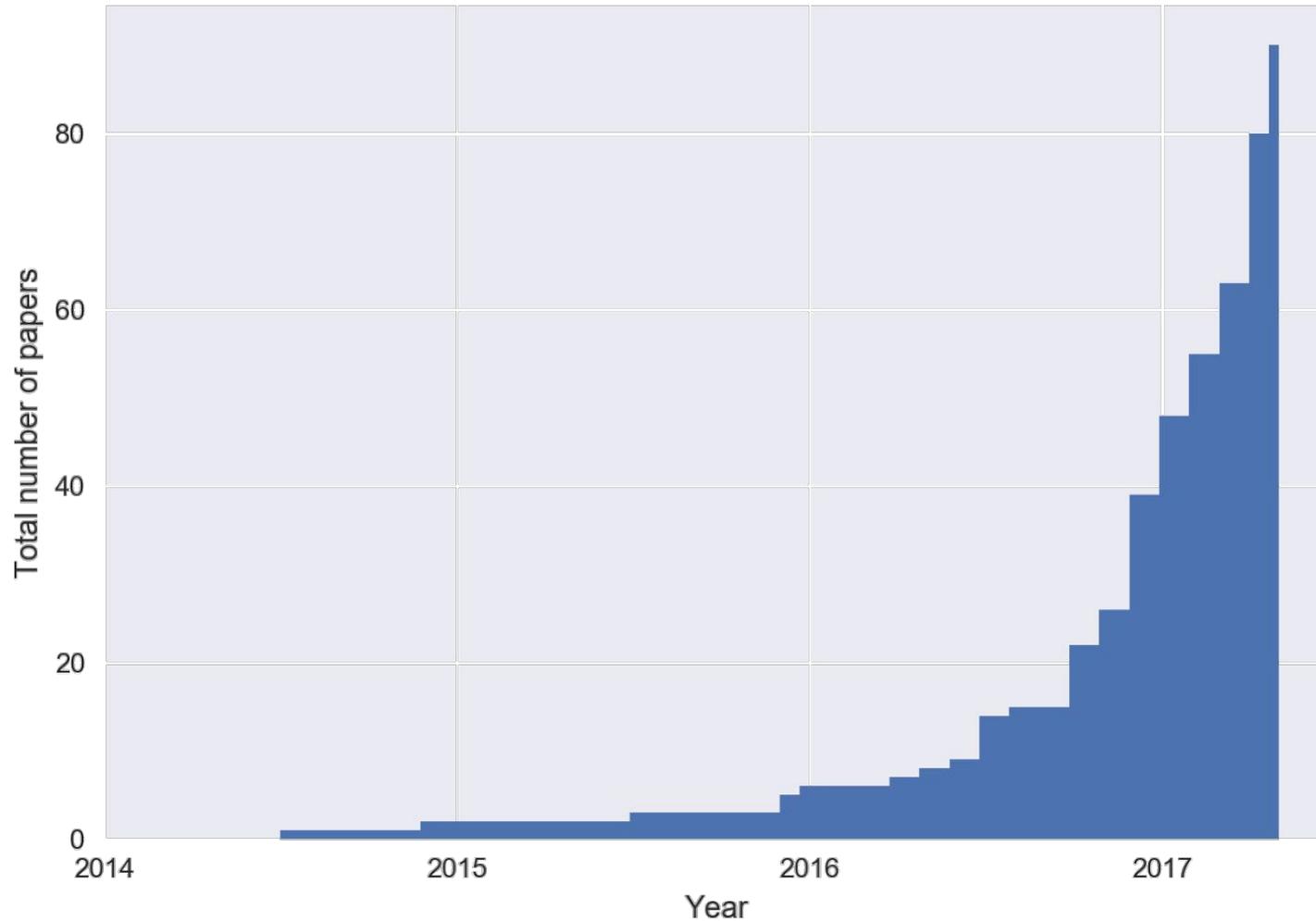


Q: What are some recent and potentially upcoming breakthroughs in deep learning?

A: The most important one, in my opinion, is adversarial training (also called GAN for Generative Adversarial Networks) ... This, and the variations that are now being proposed is **the most interesting idea in the last 10 years in ML**, in my opinion.

Progress in GANs

Cumulative number of GAN papers by year



Source: <https://deephunt.in/the-gan-zoo-79597dc8c347>

Ian Goodfellow Retweeted

Terry Taewoong Um @TerryUm_ML · Apr 6
I developed a GANN (Generative adversarial name-making networks), for you, @hardmaru @karpathy. The source code is available in powerpoint.

GANN

Generative Adversarial Name-making Networks

HooliGAN
CardiGAN
GANg

GANGNAM style transfer

G
Generate prefix or postfix of GAN

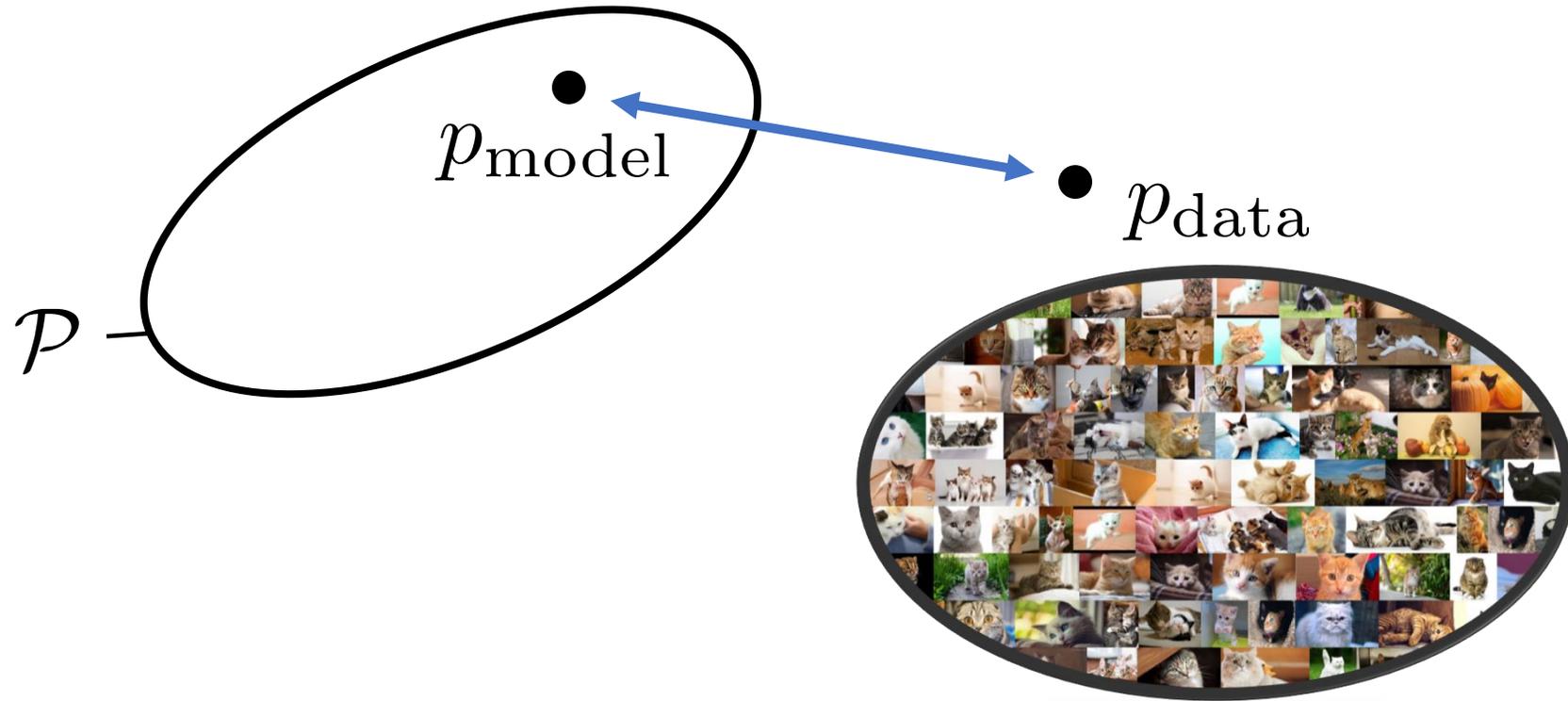
D
Determine if the name is cool or not

Character-level input

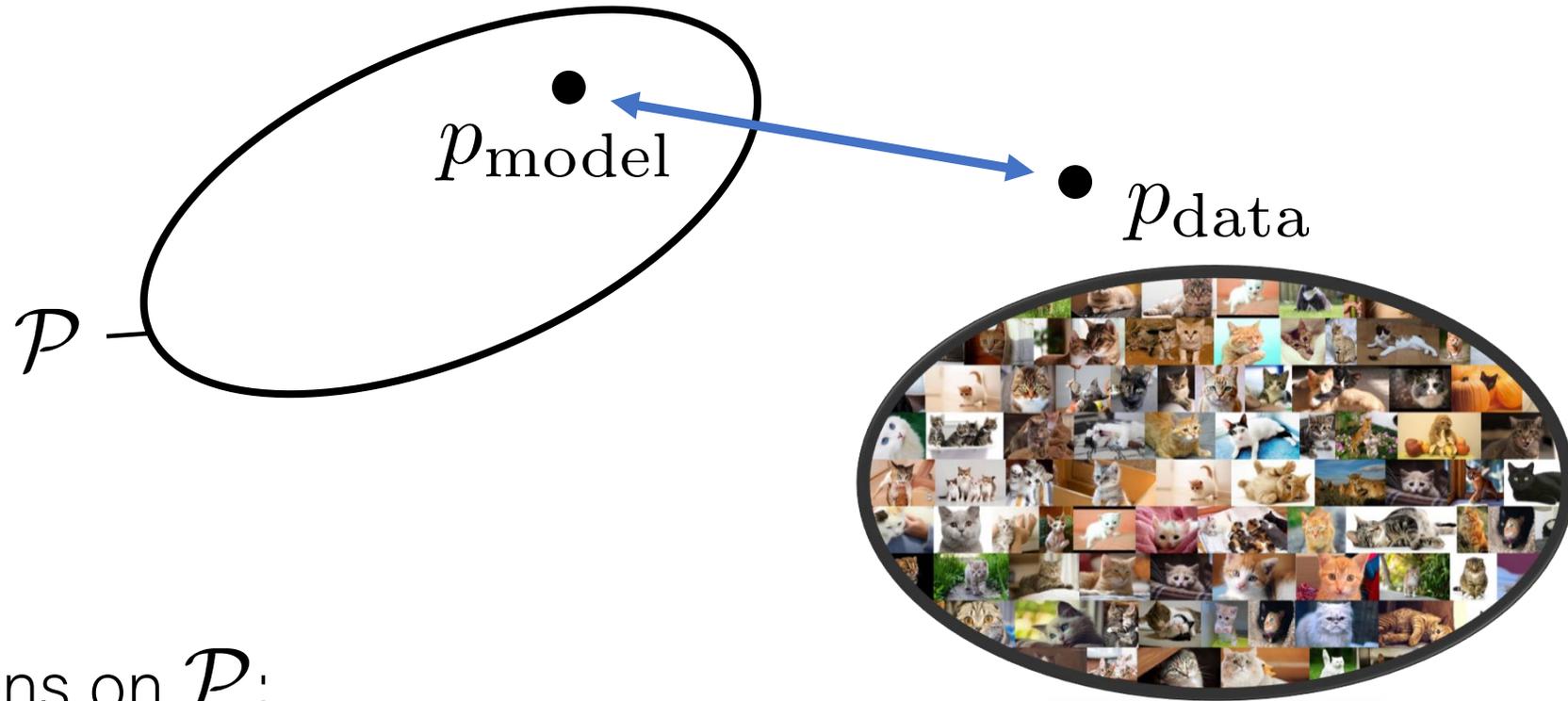
@TerryUm_ML

3 101 290

Generative Modeling



Generative Modeling

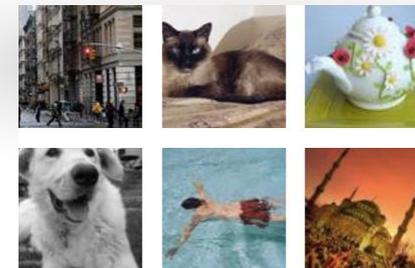


Assumptions on \mathcal{P} :

- tractable sampling

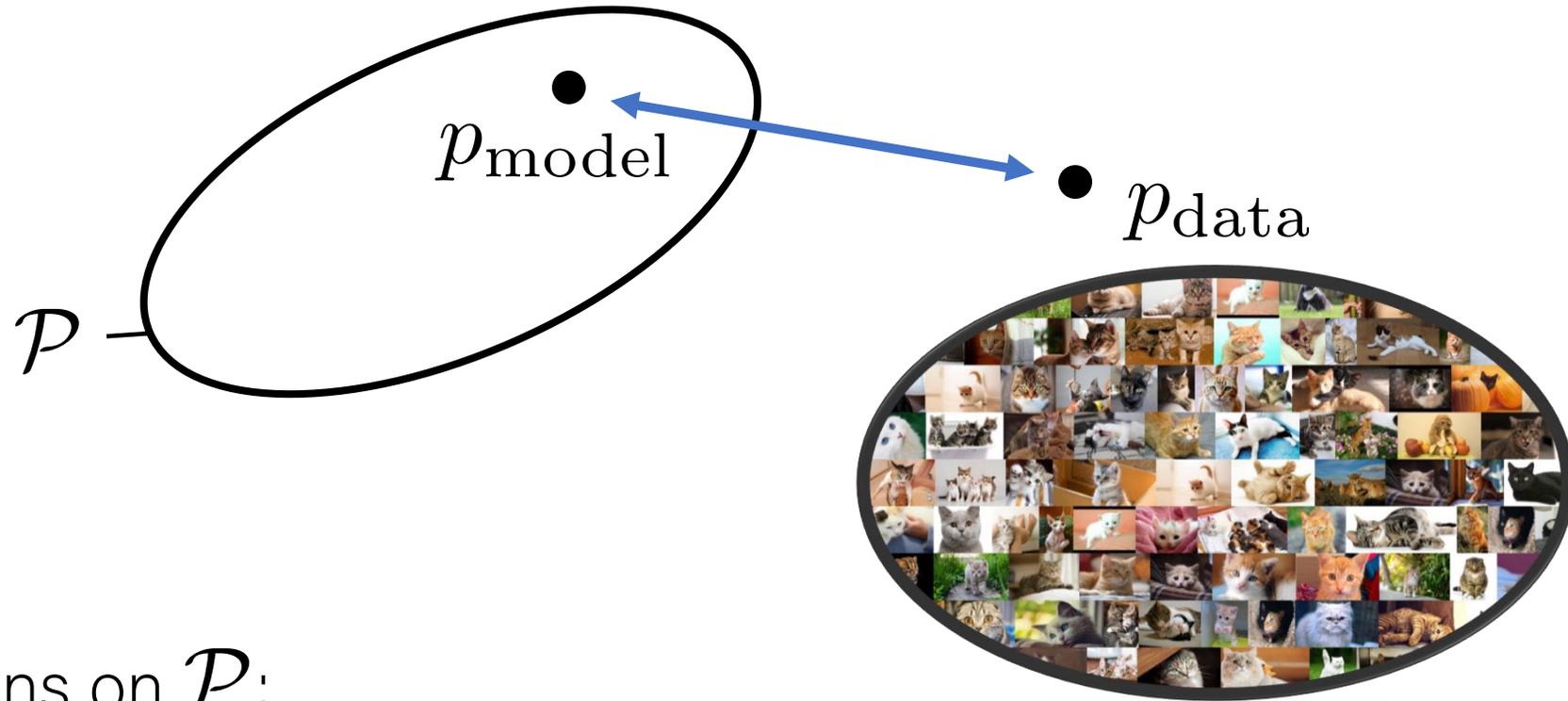


Training examples



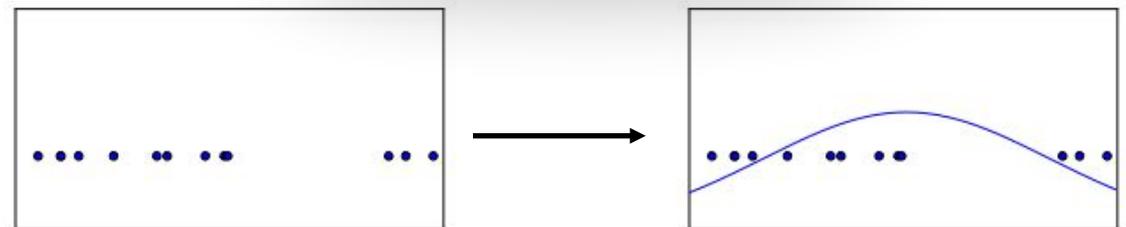
Model samples

Generative Modeling



Assumptions on \mathcal{P} :

- tractable sampling
- tractable likelihood function



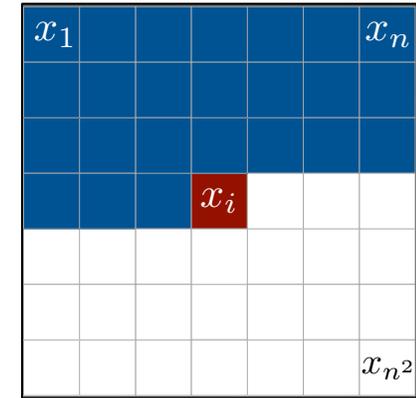
Three Broad Categories

- Autoregressive Models
- Variational Autoencoders
- Generative Adversarial Networks (GANs)

Autoregressive Models

- Explicitly model conditional probabilities:

$$p_{\text{model}}(\mathbf{x}) = p_{\text{model}}(x_1) \prod_{i=2}^n p_{\text{model}}(x_i \mid x_1, \dots, x_{i-1})$$



Each conditional can be a complicated neural net

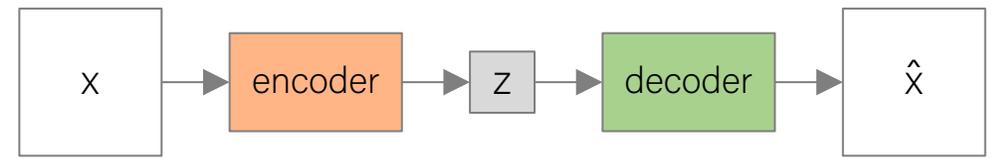
Disadvantages:

- Generation can be too costly
- Generation can not be controlled by a latent code

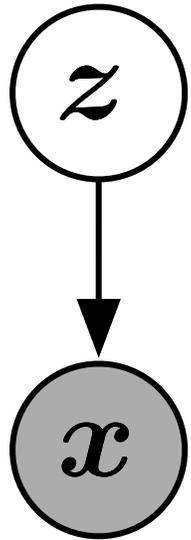


PixelCNN elephants
(van den Ord et al. 2016)

Variational Autoencoder



- Maximizes a variational lower bound on log-likelihood of \mathbf{x}



$$\begin{aligned} \log p(\mathbf{x}) &\geq \log p(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x})) \\ &= \mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{x}, \mathbf{z}) + H(q) \end{aligned}$$



Face samples for
Labeled Faces in the Wild (LFW)
(Alec Radford)

Disadvantages:

- Not asymptotically consistent unless q is perfect
- Tends to produce blurry samples

GANs

Generative
Adversarial
Networks

Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)



Noise
(random input)

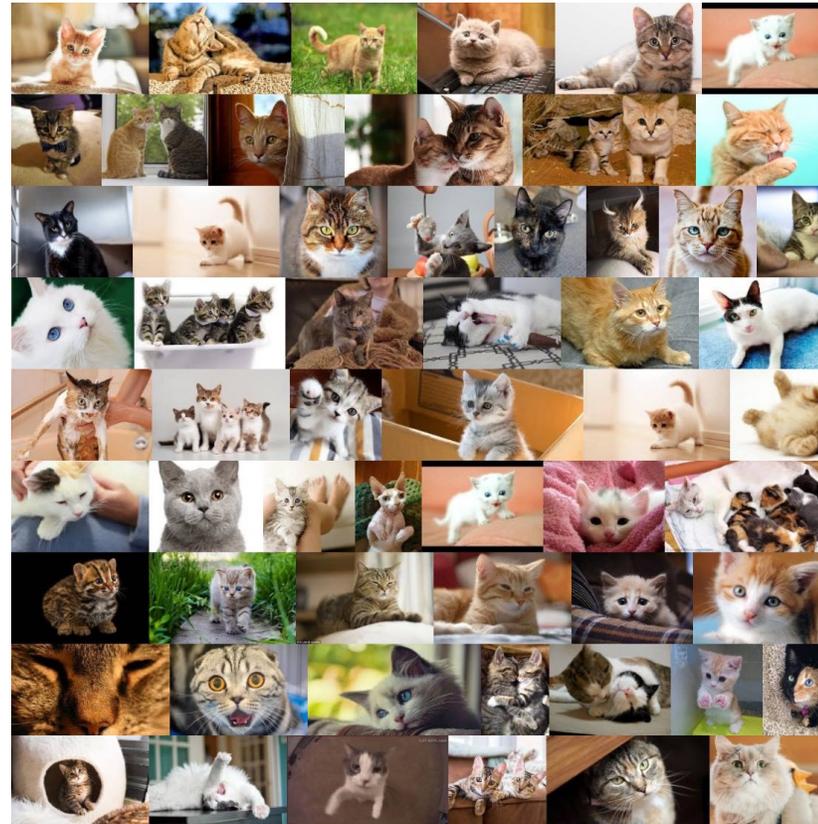


Generative
Model



$z \sim \text{Uniform}_{100}$

*think of this as
a transformation*



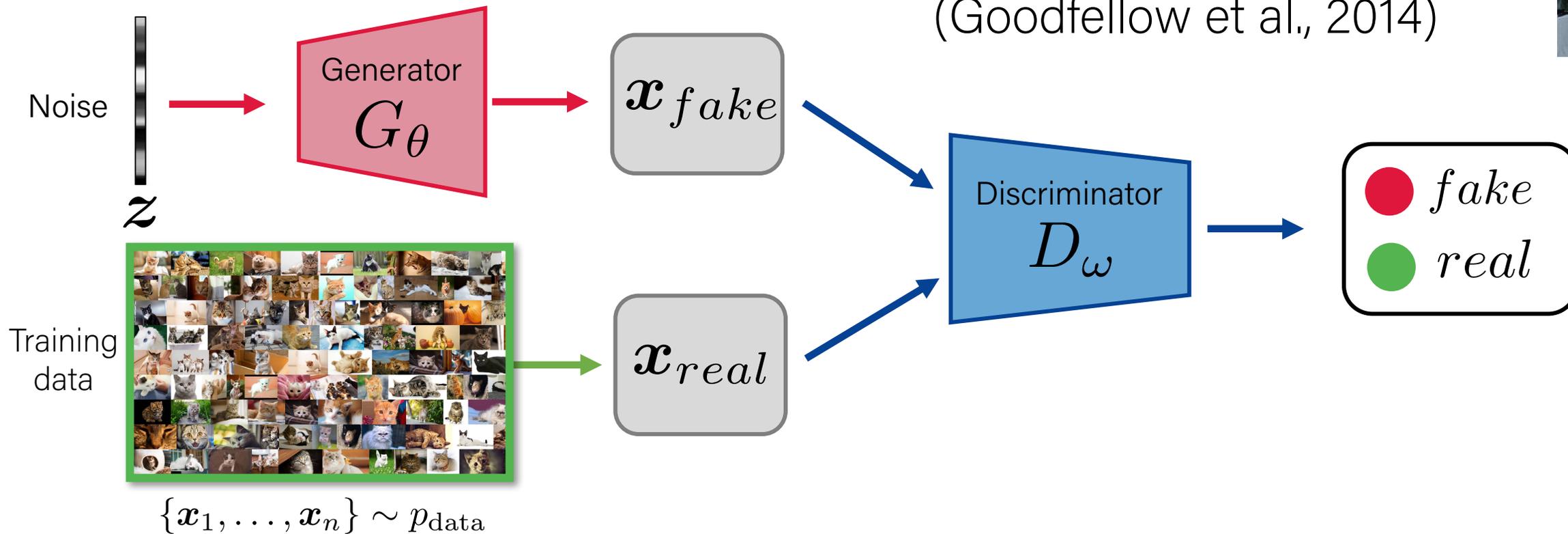
- A game-theoretic likelihood free model

Advantages:

- Uses a latent code
- No Markov chains needed
- Produces the best samples

Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)



- A game between a generator $G_\theta(z)$ and a discriminator $D_\omega(x)$
 - Generator tries to fool discriminator (i.e. generate realistic samples)
 - Discriminator tries to distinguish fake from real samples

Intuition behind GANs



D_ω : Discriminator (*Art Critic*)



x_{real}



x_{fake}



G_θ : Generator (*Forger*)

Intuition behind GAN Training



<https://www.youtube.com/watch?v=No26JKQKZNE>

GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\omega}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log (1 - D_{\omega}(G_{\theta}(\mathbf{z})))]$$

Real data

Noise vector used to
generate data

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

Cross-entropy
loss for binary
classification

Generator maximizes the log-probability
of the discriminator being mistaken

- Equilibrium of the game
- Minimizes the Jensen-Shannon divergence

GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\omega}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log (1 - D_{\omega}(G_{\theta}(\mathbf{z})))]$$

Real data

Noise vector used to

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\omega}(\mathbf{x})]$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log (1 - D_{\omega}(G_{\theta}(\mathbf{z})))]$$

**Important question is
"Does this converge??"**

Cross-entropy loss for binary classification

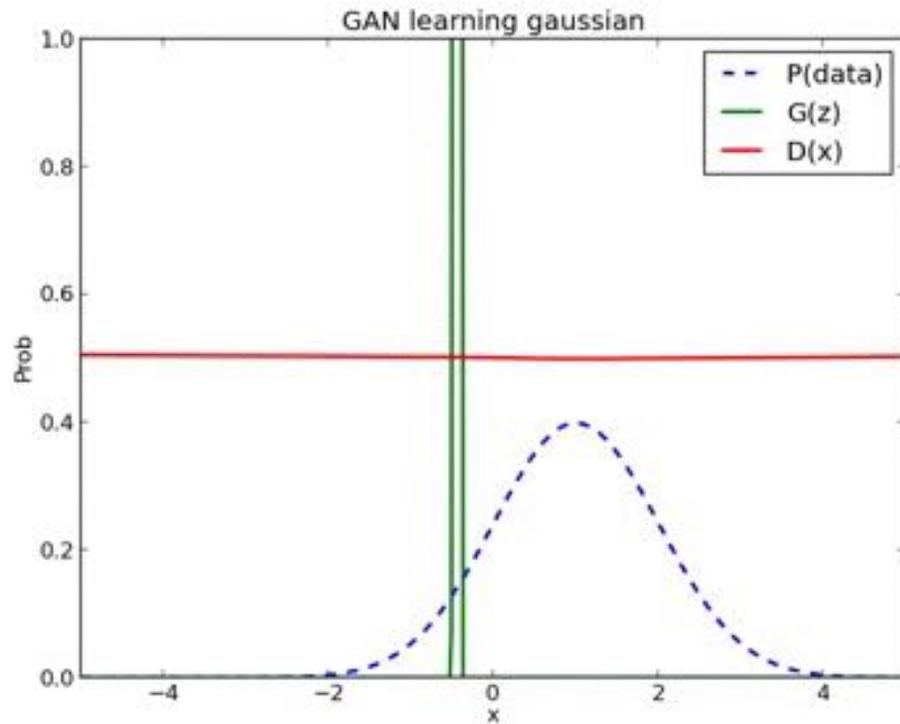
probability

of the discriminator being mistaken

- Equilibrium of the game
- Minimizes the Jensen-Shannon divergence

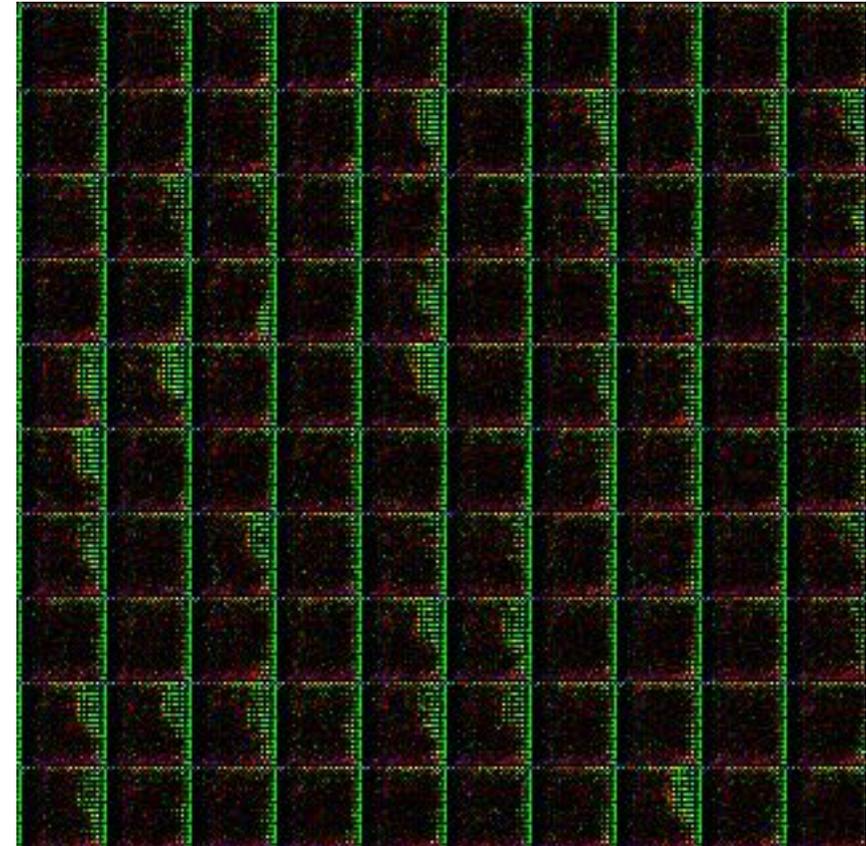
Training Procedure

(Goodfellow et al., 2014)



Source: Alec Radford

Generating 1D points



Source: OpenAI blog

Generating images

Results

(Goodfellow et al., 2014)

- The generator uses a mixture of rectifier linear activations and/or sigmoid activations
- The discriminator net used maxout activations.



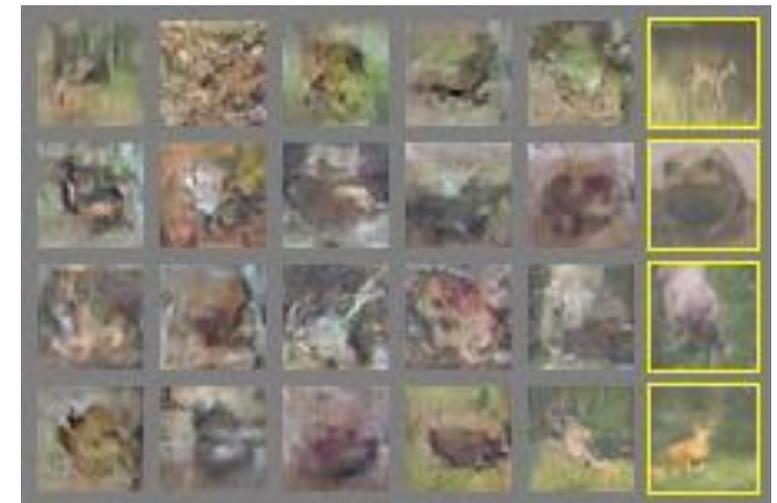
MNIST samples



TFD samples



CIFAR10 samples
(fully-connected model)

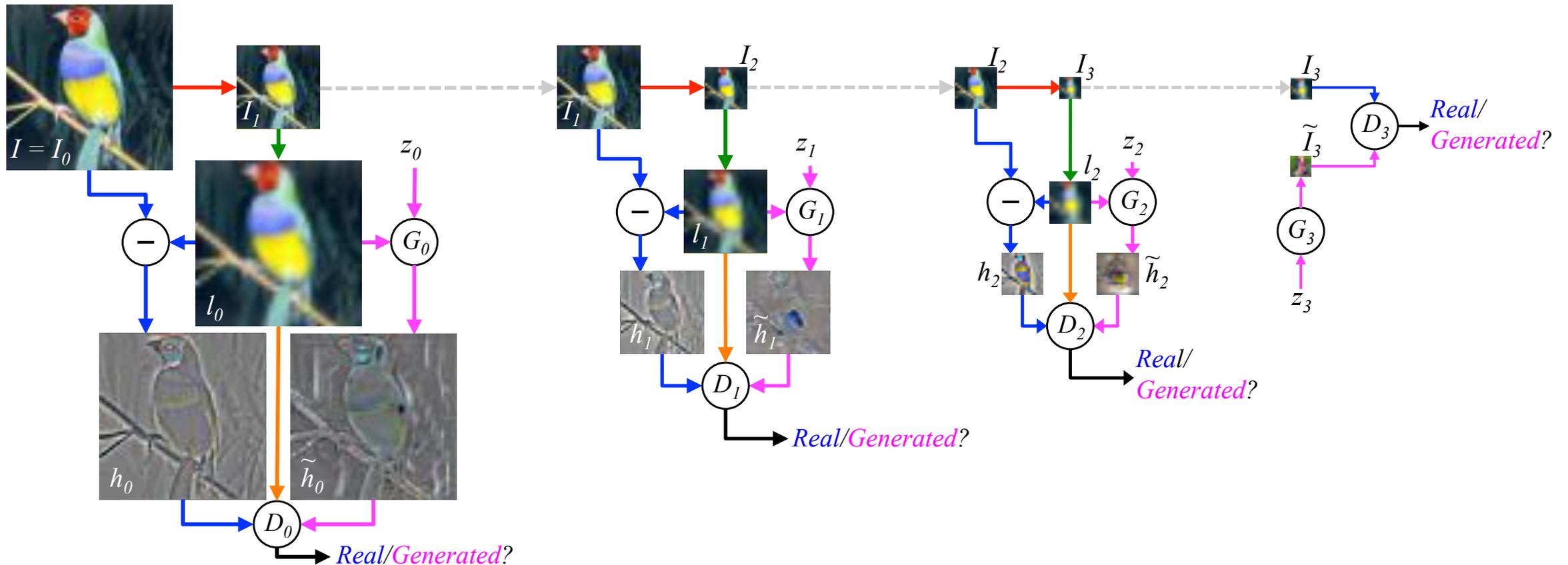


CIFAR10 samples
*(convolutional discriminator,
deconvolutional generator)*

Laplacian GANs (LAPGAN)

(Denton et al., 2015)

- **Idea:** Combine GAN with a multi-scale image representation (Laplacian pyramid)

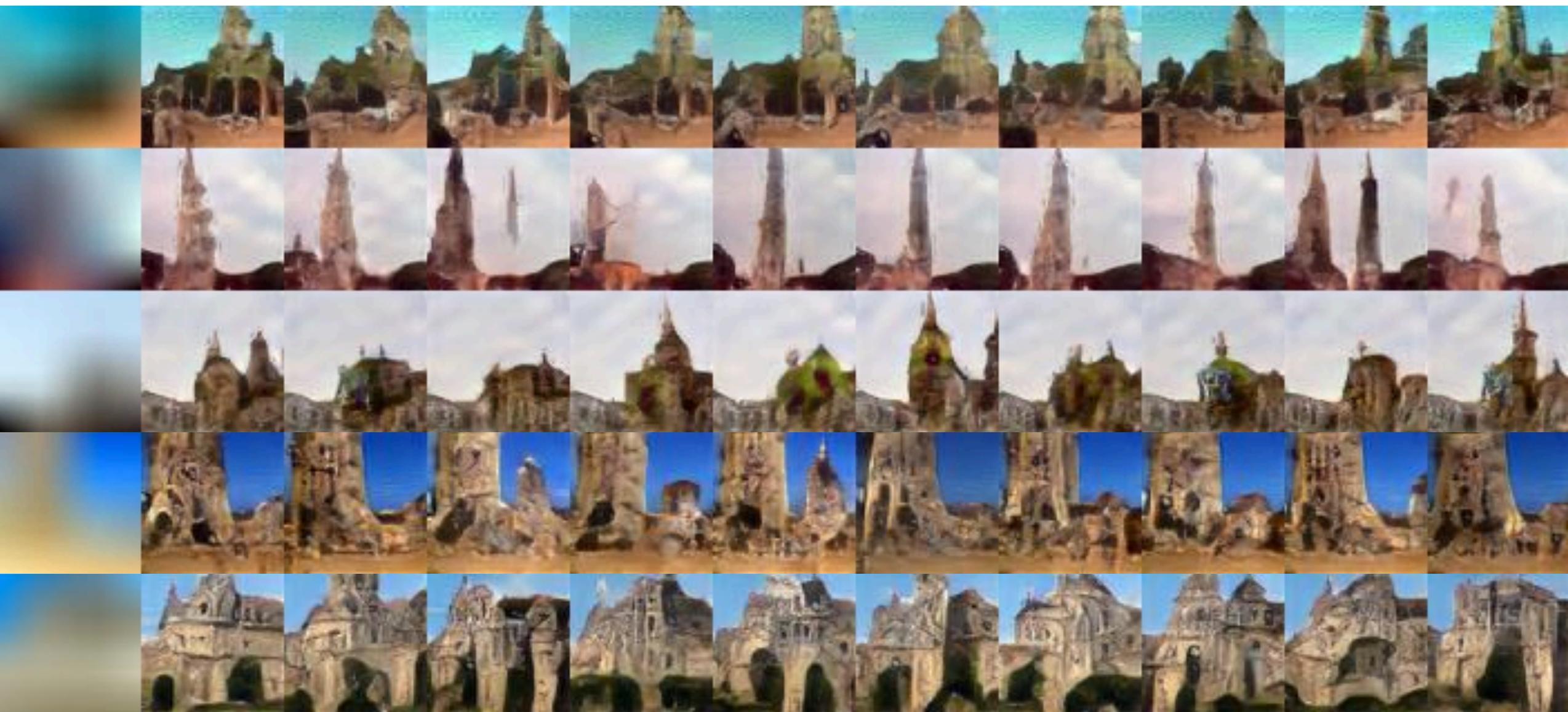


LAPGAN for LSUN Towers

64×64 pixels

~700K images

(Denton et al., 2015)



LAPGAN for LSUN Bedrooms 64×64 pixels ~3M images (Denton et al., 2015)

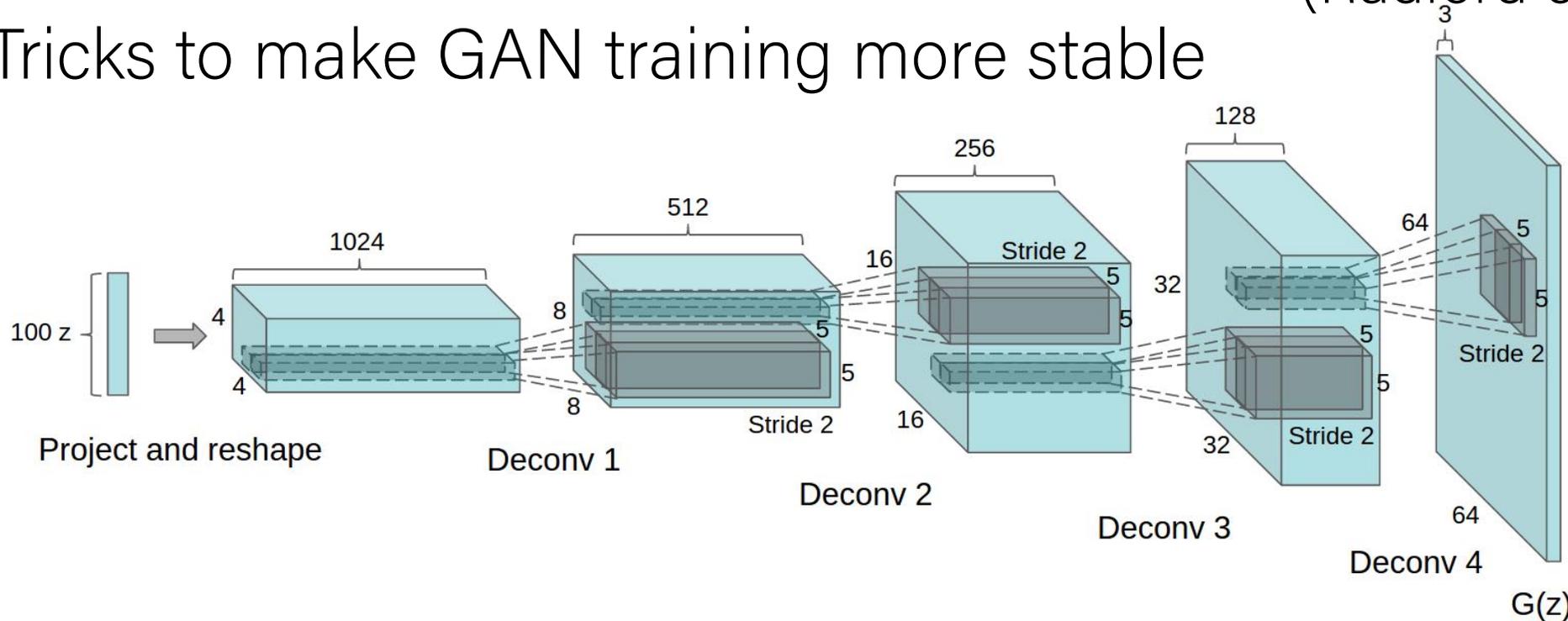


Deep Convolutional GANs (DCGAN)



(Radford et al., 2015)

- **Idea:** Tricks to make GAN training more stable



- No fully connected layers
- Batch Normalization (Ioffe and Szegedy, 2015)
- Leaky Rectifier in D
- Use Adam (Kingma and Ba, 2015)
- Tweak Adam hyperparameters a bit ($\text{lr}=0.0002$, $\text{b1}=0.5$)

DCGAN for LSUN Bedrooms

64×64 pixels

~3M images

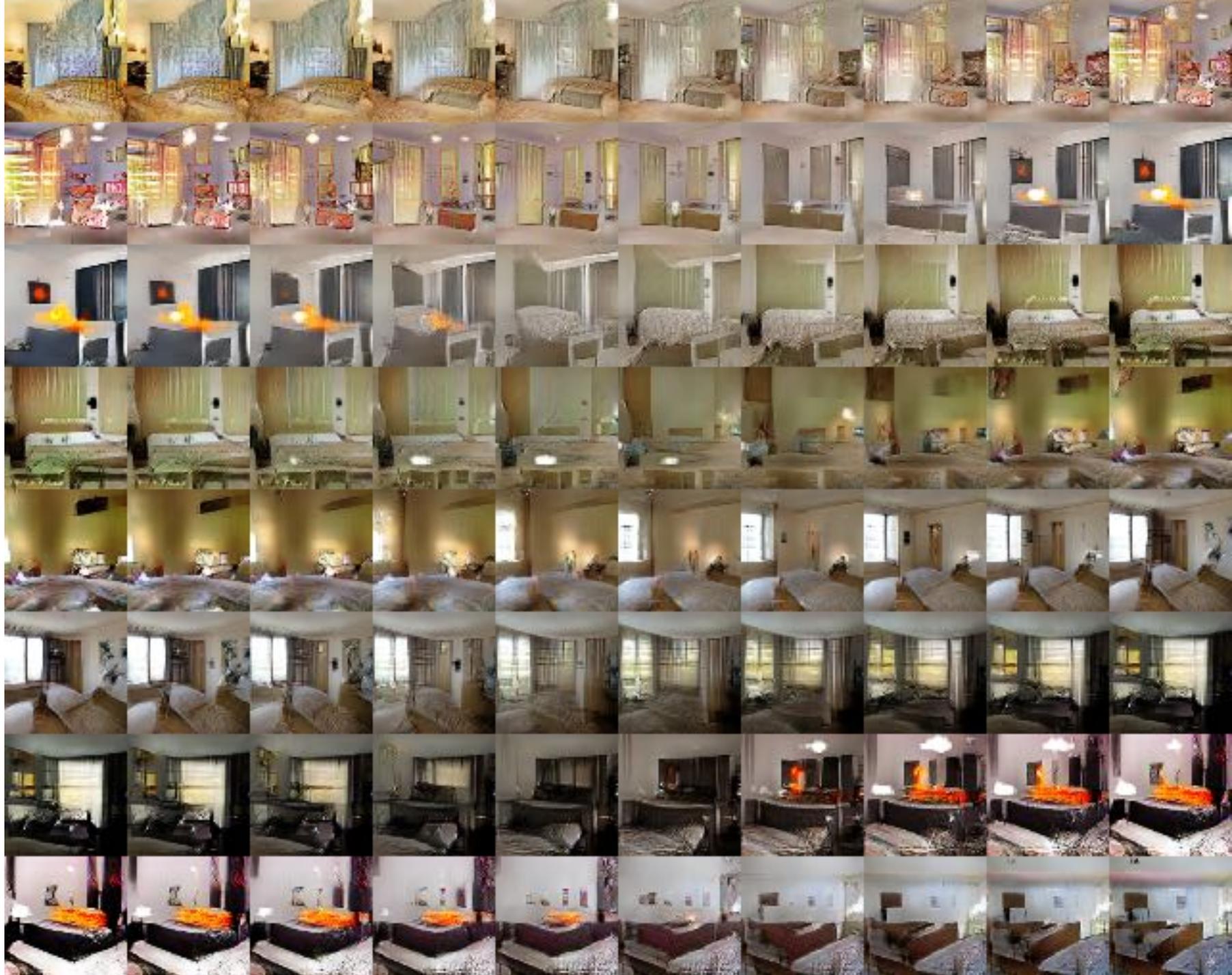
(Radford et al., 2015)



Walking over the latent space

(Radford et al., 2015)

- Interpolation suggests non-overfitting behavior



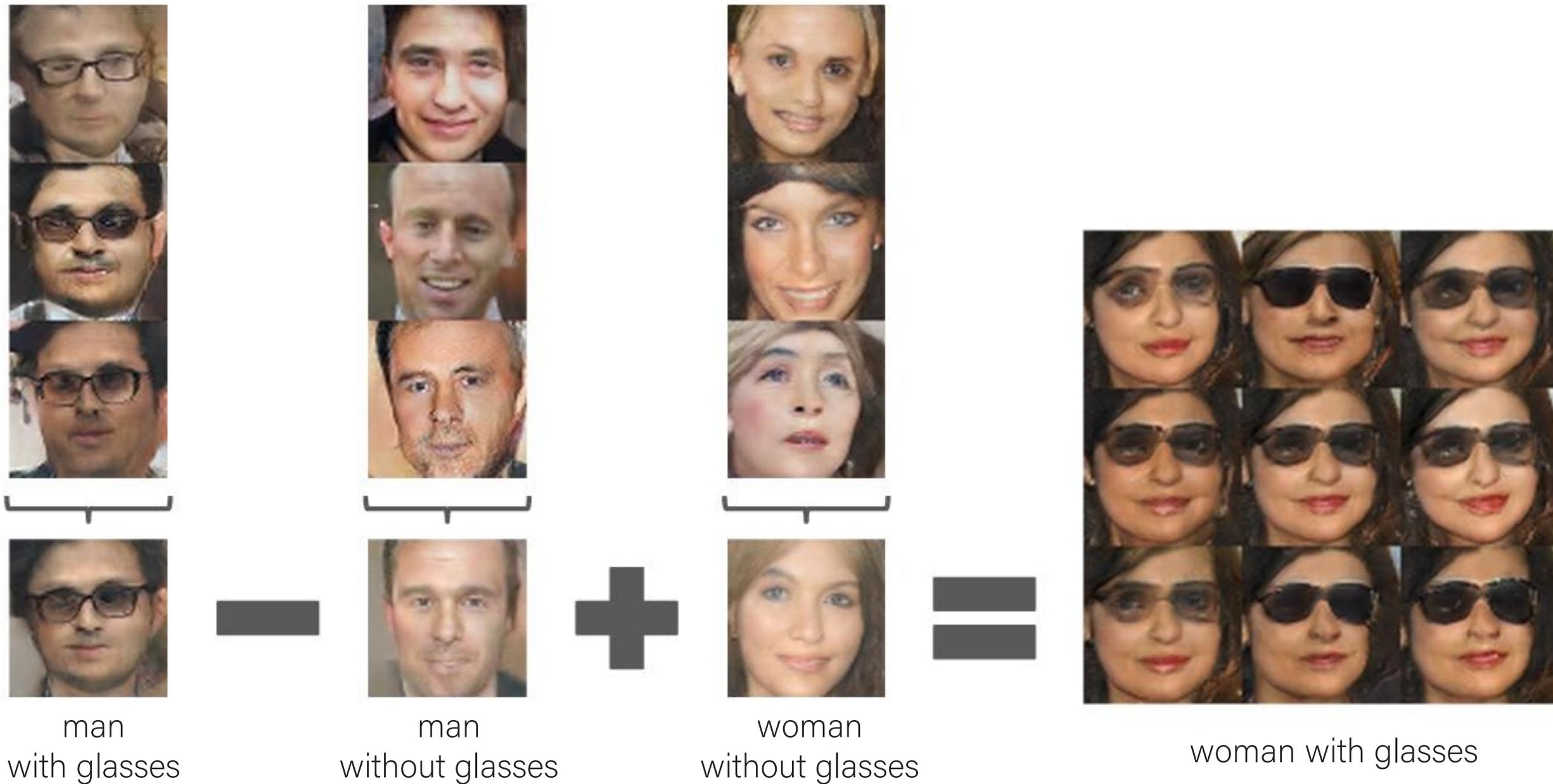
Walking over the latent space

(Radford et al., 2015)



Vector Space Arithmetic

(Radford et al., 2015)



Vector Space Arithmetic

(Radford et al., 2015)



smiling woman



neutral woman



neutral man

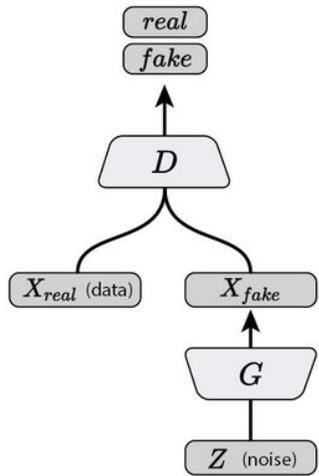


smiling man

Subclasses of GANs

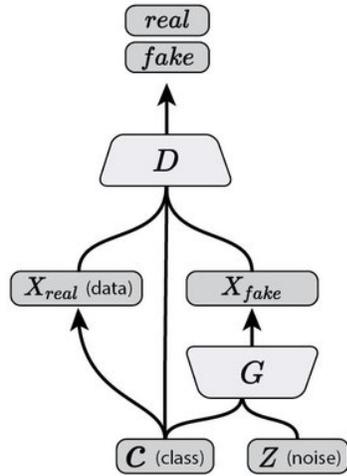
Vanilla GAN

Vanilla GAN
(Goodfellow, et al., 2014)

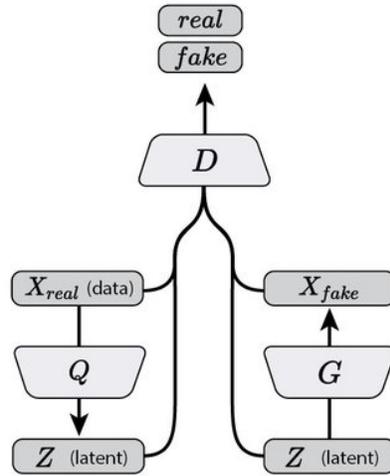


Discriminator Looks at Latent Variables

Conditional GAN
(Mirza & Osindero, 2014)

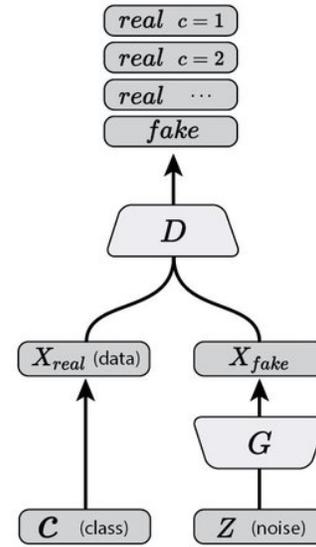


Bidirectional GAN
(Donahue, et al., 2016; Dumoulin, et al., 2016)

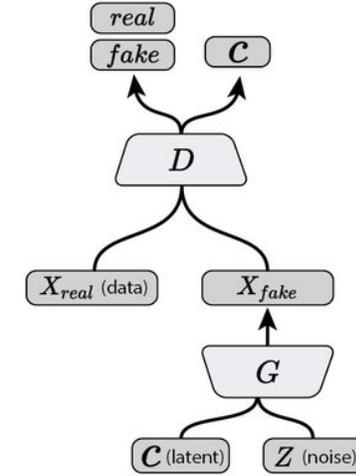


Discriminator Predicts Latent Variables

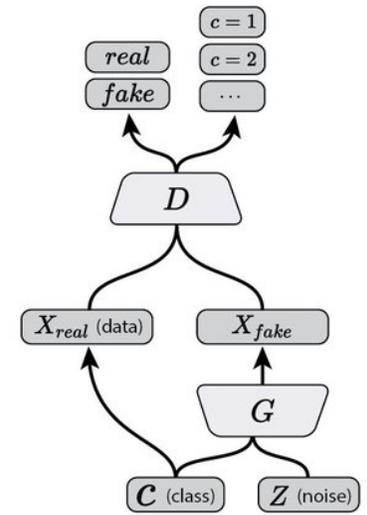
Semi-Supervised GAN
(Odena, 2016; Salimans, et al., 2016)



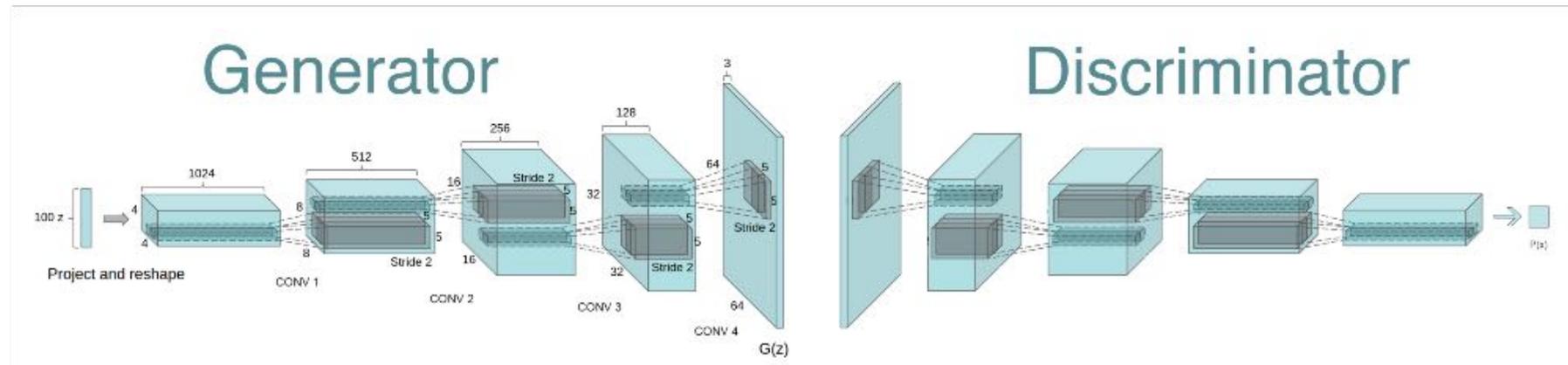
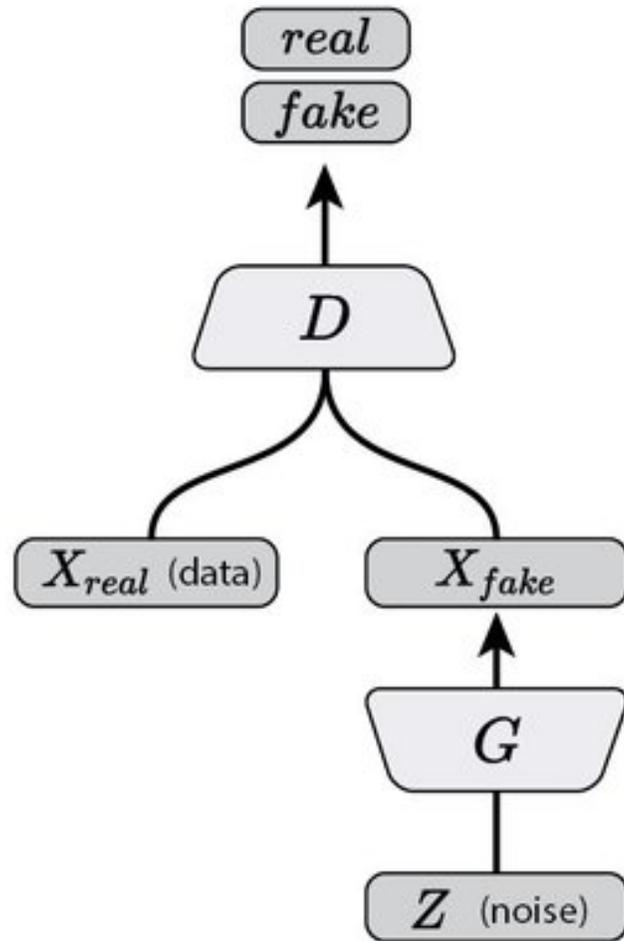
InfoGAN
(Chen, et al., 2016)



Auxiliary Classifier GAN
(Odena, et al., 2016)



Vanilla GAN (Goodfellow et al., 2014)

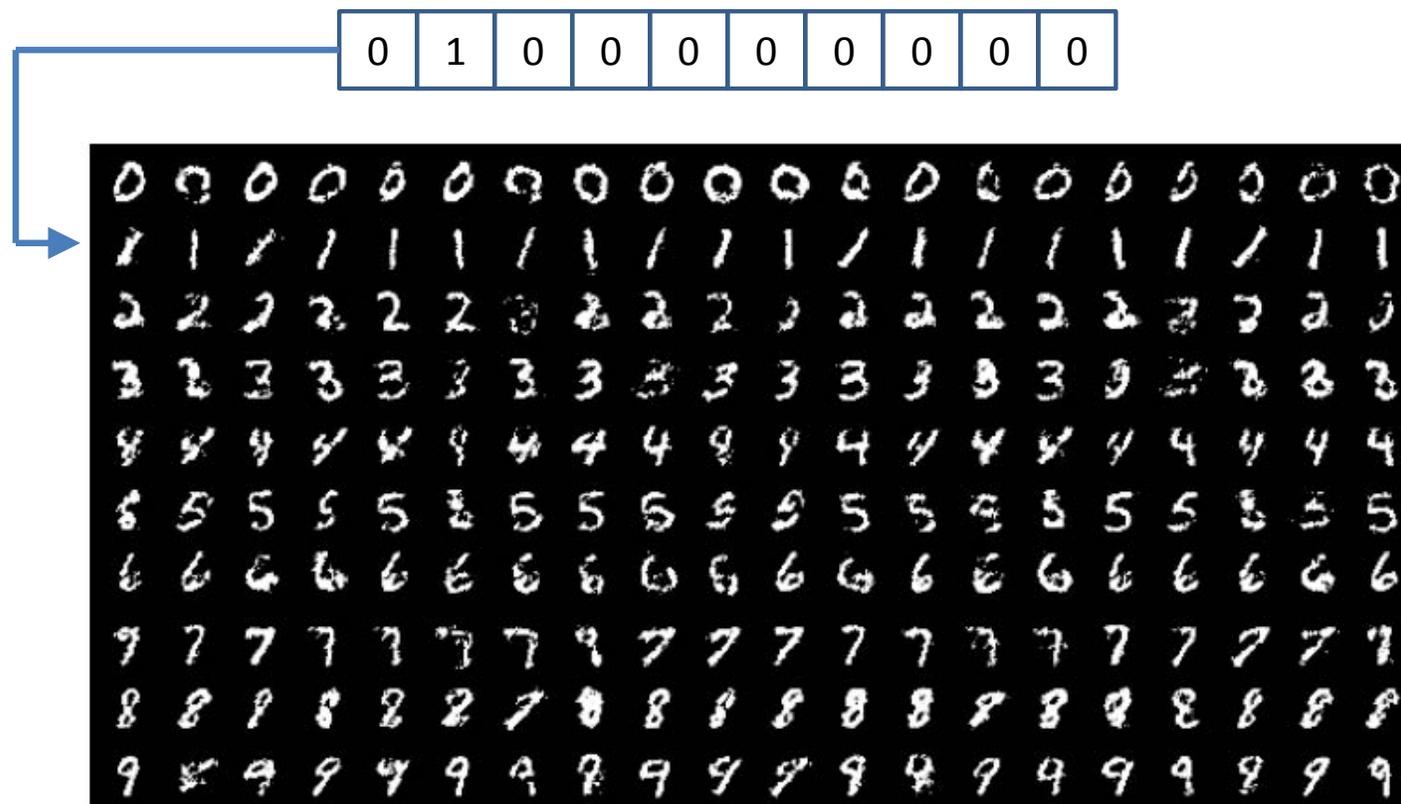
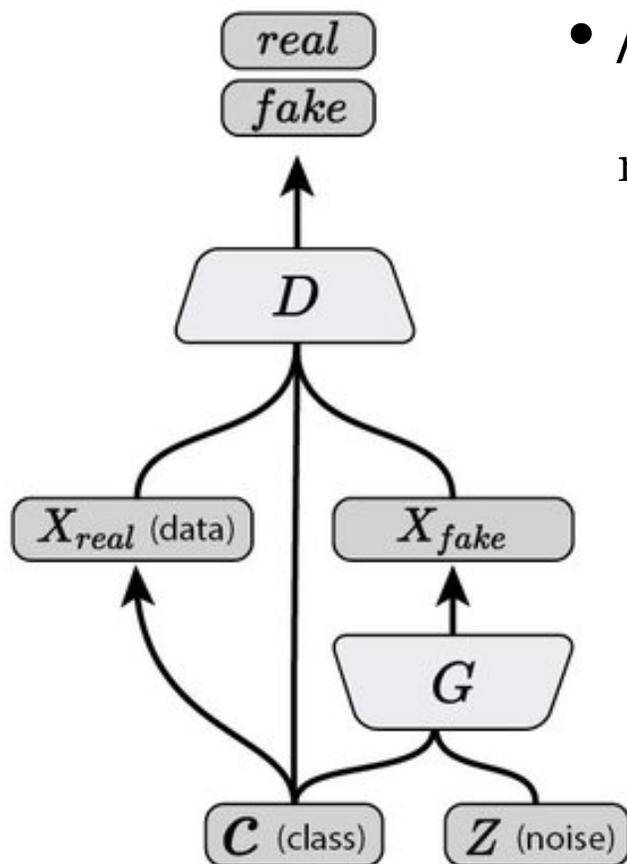


DCGAN (Radford et al., 2015)

Conditional GAN (Mirza and Osindero, 2014)

- Add conditional variables \mathbf{y} into G and D

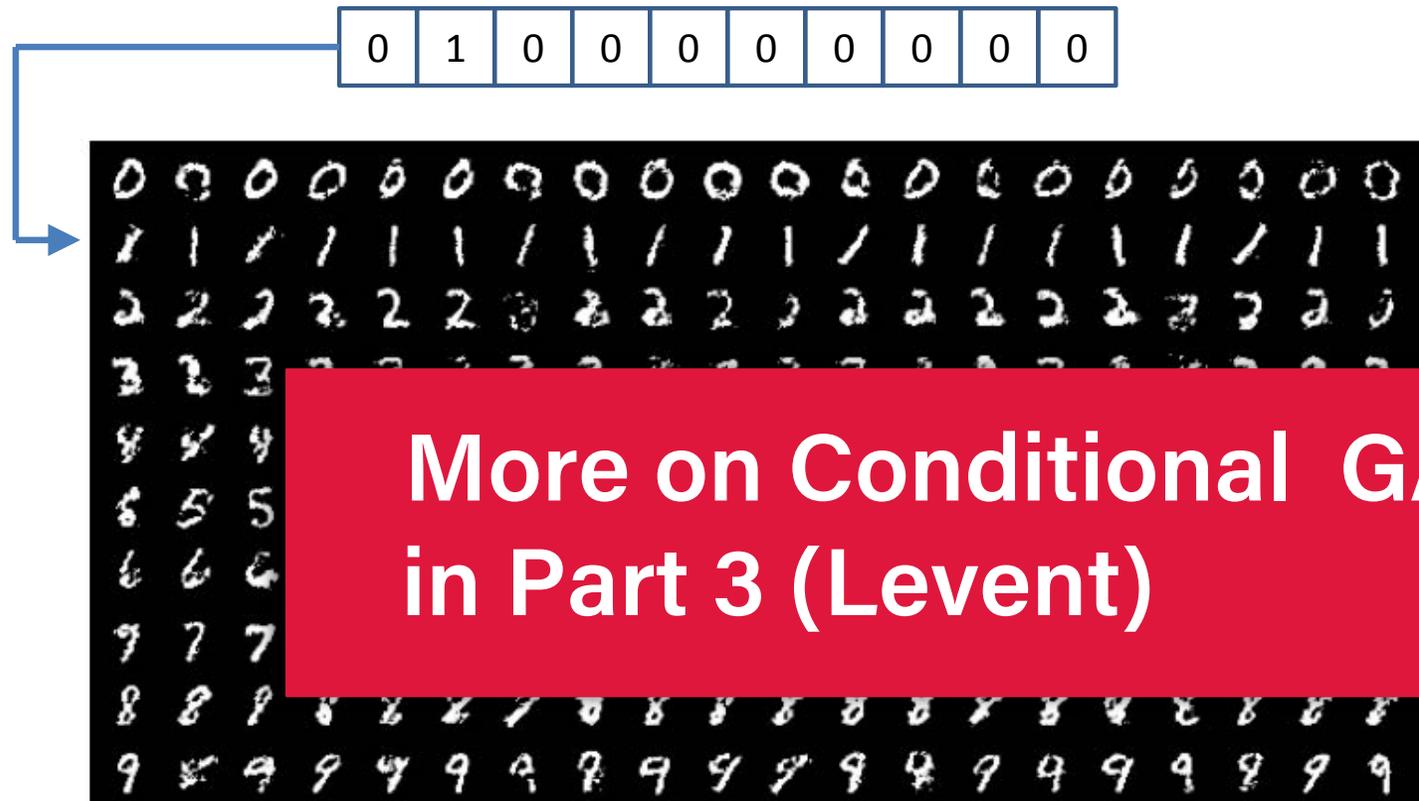
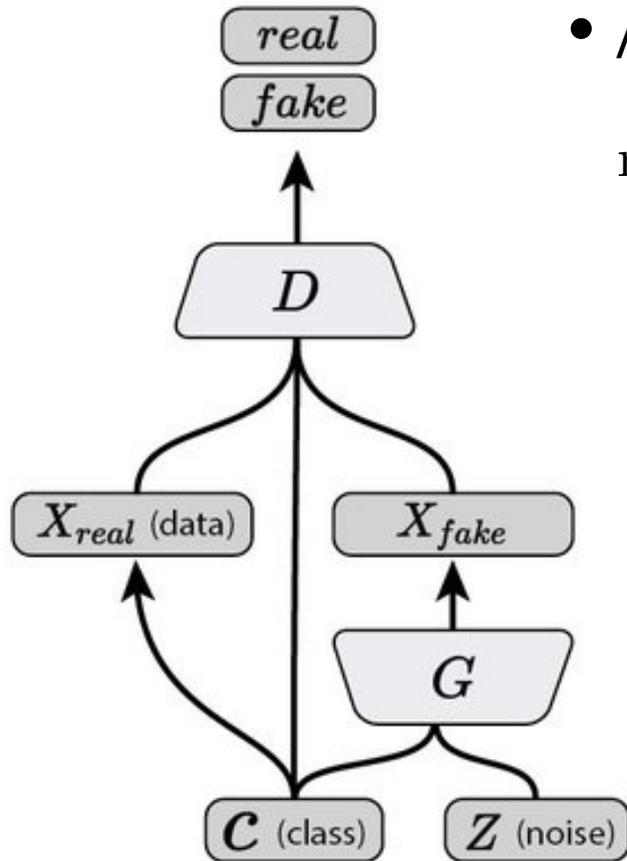
$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))]$$



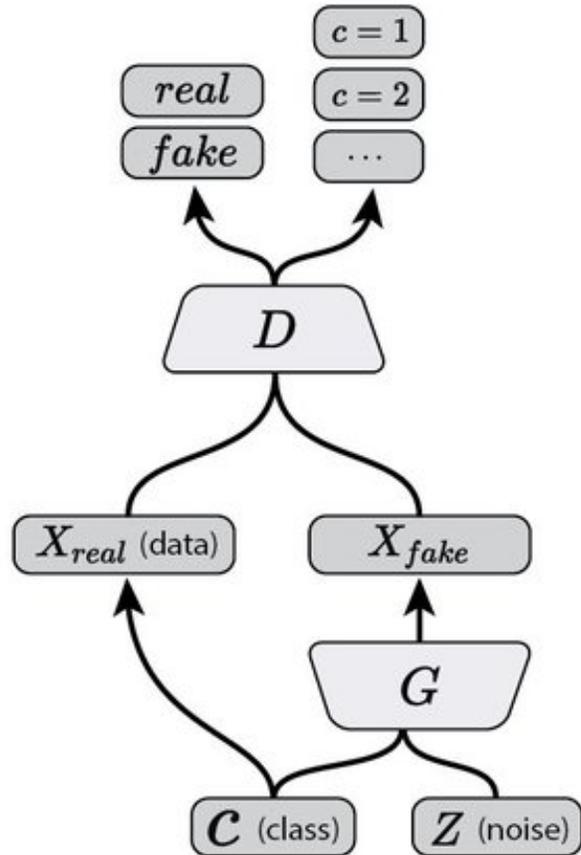
Conditional GAN (Mirza and Osindero, 2014)

- Add conditional variables \mathbf{y} into G and D

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x} | \mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z} | \mathbf{y})))]$$



Auxiliary Classifier GAN (Odena et al., 2016)



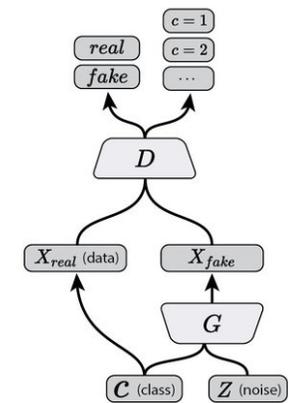
- Every generated sample has a corresponding class label

$$L_S = E[\log P(S = real \mid X_{real})] + E[\log P(S = fake \mid X_{fake})]$$

$$L_C = E[\log P(C = c \mid X_{real})] + E[\log P(C = c \mid X_{fake})]$$

- D is trained to maximize $L_S + L_C$
- G is trained to maximize $L_C - L_S$
- Learns a representation for z that is independent of class label

Auxiliary Classifier GAN (Odena et al., 2016)



128×128 resolution samples from 5 classes taken from an AC-GAN trained on the ImageNet



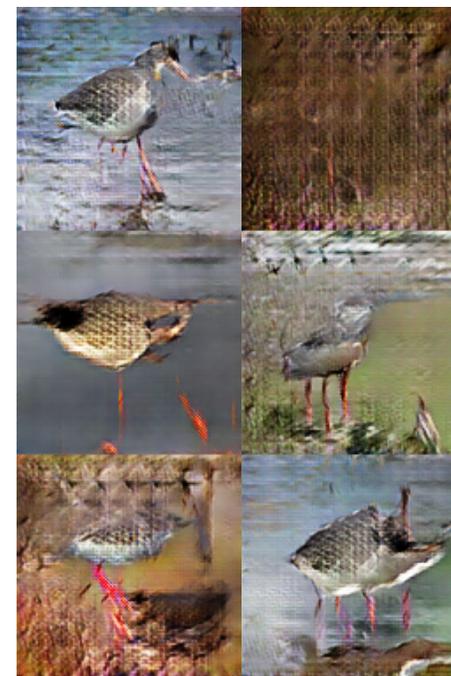
monarch butterfly



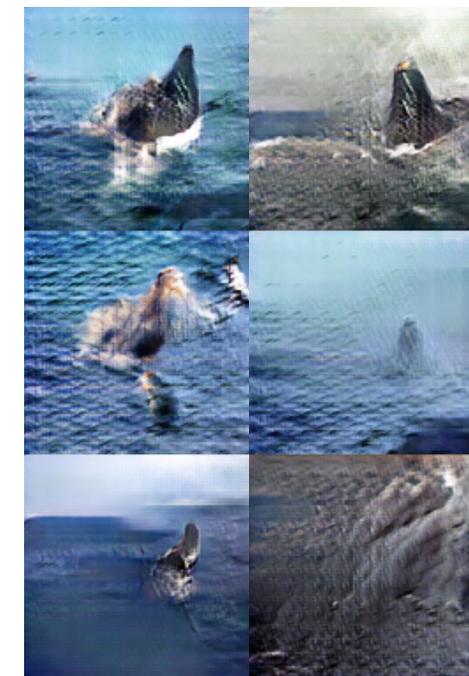
goldfinch



daisy



redshank

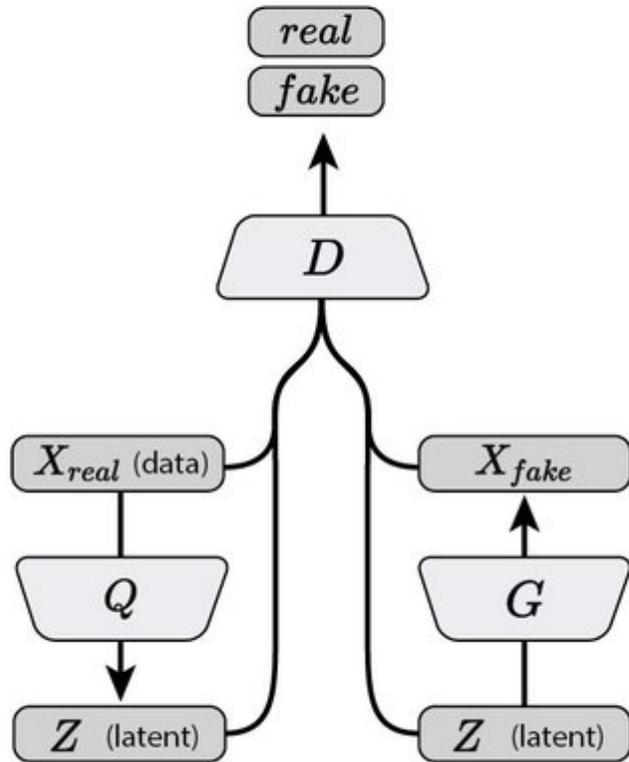


grey whale

Bidirectional GAN (Donahue et al., 2016; Dumoulin et al., 2016)

- Jointly learns a generator network and an inference network using an adversarial process.

$$\begin{aligned} \min_G \max_D V(D, G) &= \mathbb{E}_{q(\mathbf{x})} [\log(D(\mathbf{x}, G_z(\mathbf{x})))] + \mathbb{E}_{p(\mathbf{z})} [\log(1 - D(G_x(\mathbf{z}), \mathbf{z}))] \\ &= \iint q(\mathbf{x})q(\mathbf{z} | \mathbf{x}) \log(D(\mathbf{x}, \mathbf{z})) d\mathbf{x}d\mathbf{z} \\ &+ \iint p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) \log(1 - D(\mathbf{x}, \mathbf{z})) d\mathbf{x}d\mathbf{z}. \end{aligned}$$

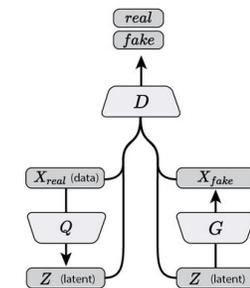


CelebA reconstructions

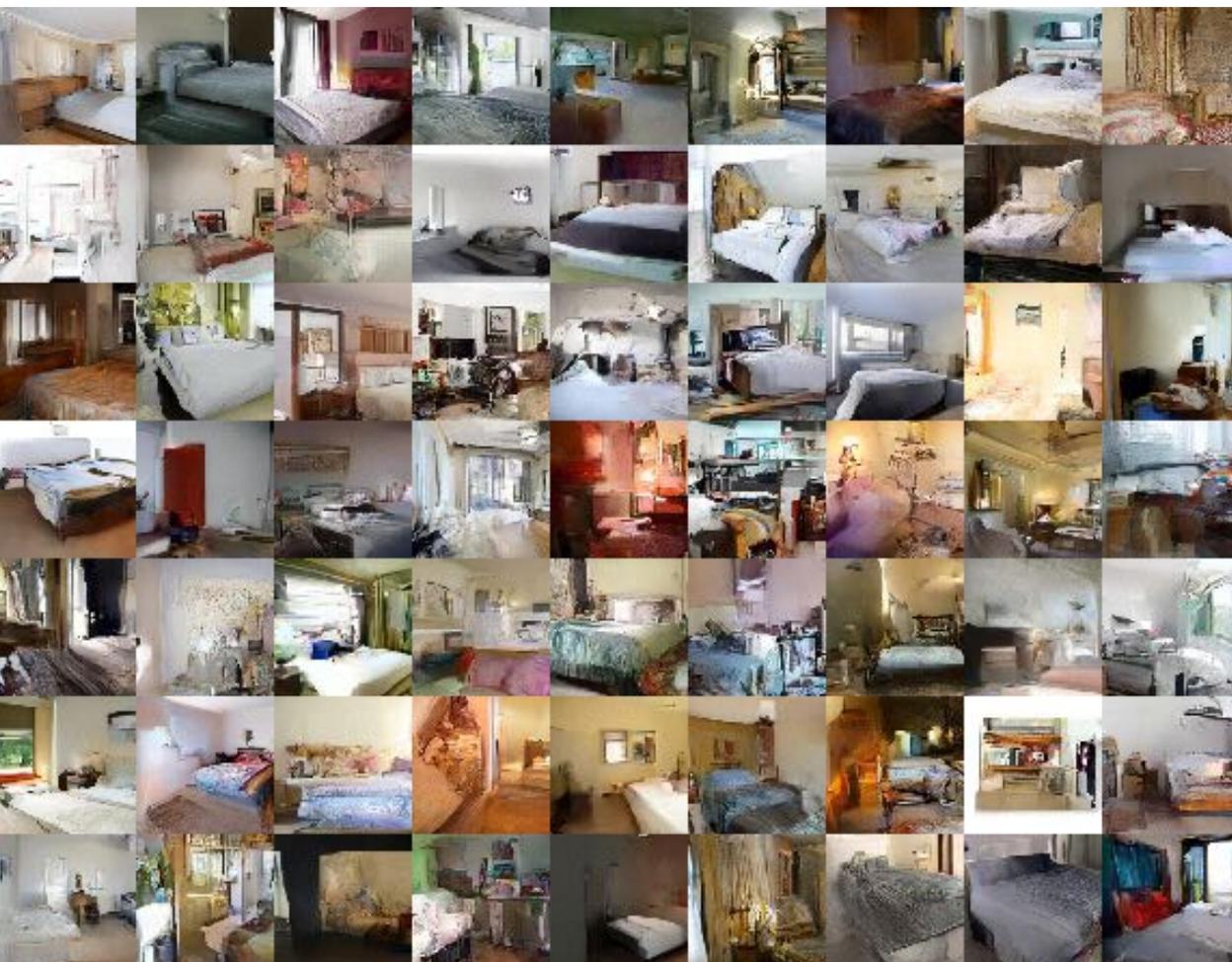


SVNH reconstructions

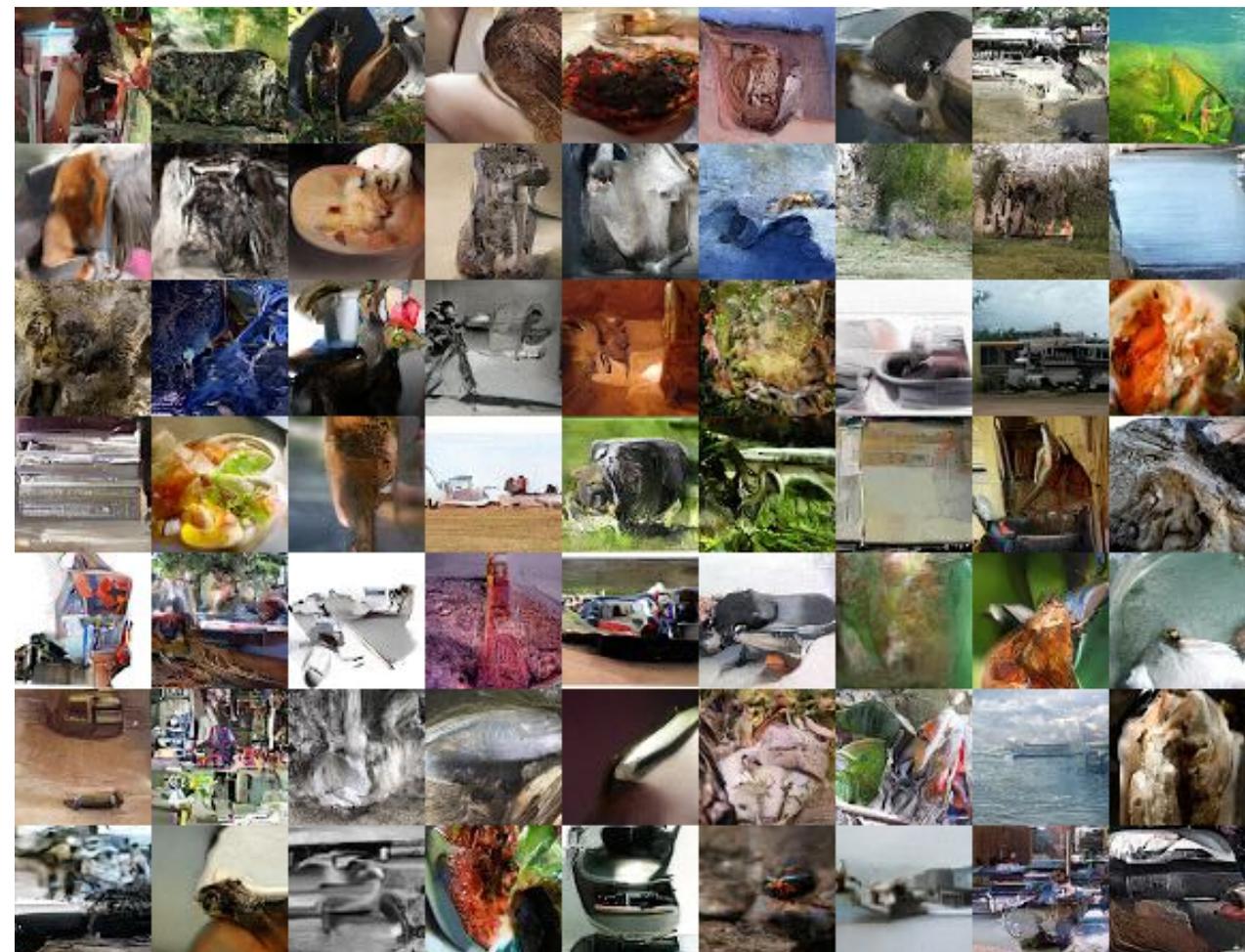
Bidirectional GAN (Donahue et al., 2016; Dumoulin et al., 2016)



LSUN bedrooms



Tiny ImageNet



Applications of GANs

Semi-supervised Classification

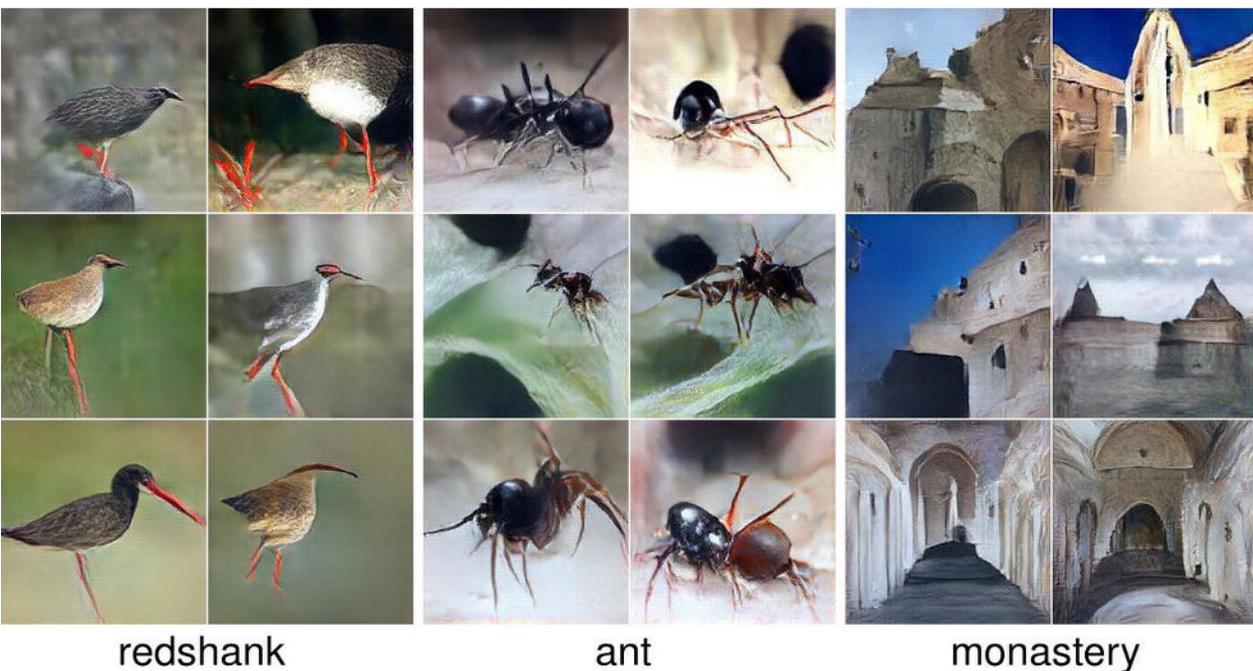
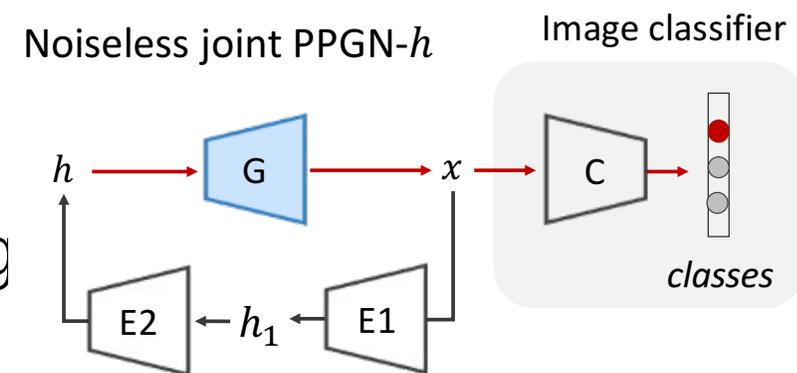
(Salimans et al., 2016;
Dumoulin et al., 2016)

SVNH

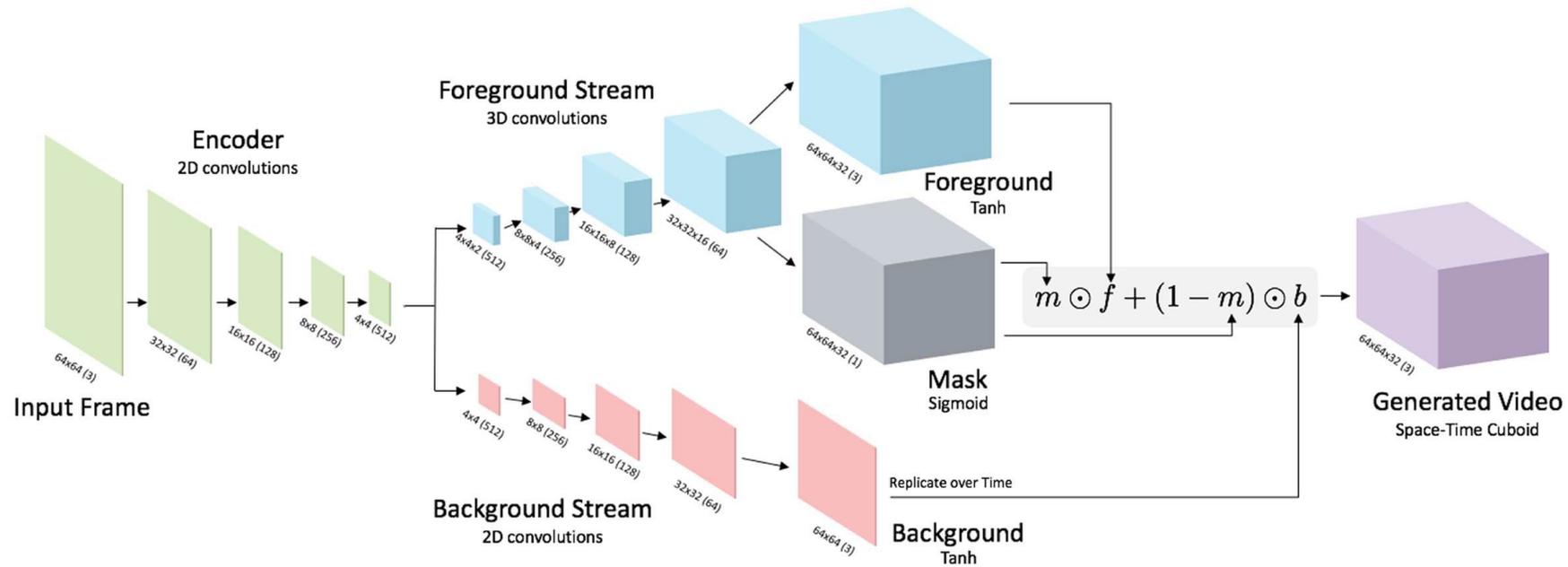
Model	Misclassification rate
VAE (M1 + M2) (Kingma et al., 2014)	36.02
SWWAE with dropout (Zhao et al., 2015)	23.56
DCGAN + L2-SVM (Radford et al., 2015)	22.18
SDGM (Maaløe et al., 2016)	16.61
GAN (feature matching) (Salimans et al., 2016)	8.11 ± 1.3
ALI (ours, L2-SVM)	19.14 ± 0.50
ALI (ours, no feature matching)	7.42 ± 0.65

Class-specific Image Generation (Nguyen et al., 2016)

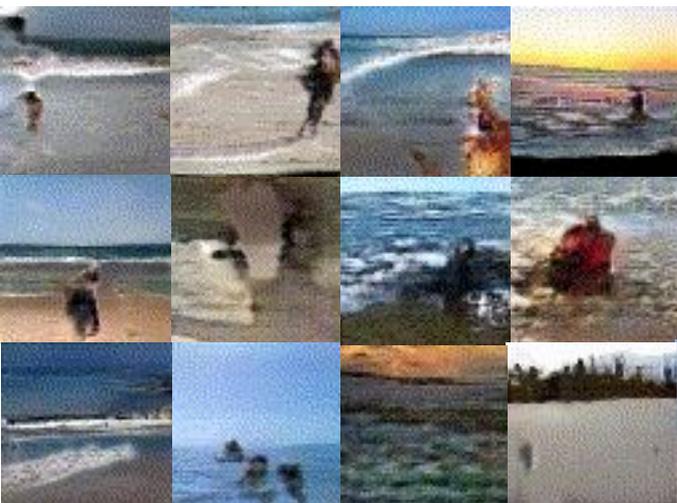
- Generates 227x227 realistic images from all ImageNet classes
- Combines adversarial training, moment matching denoising autoencoders, and Langevin sampling



Video Generation (Vondrick et al., 2016)



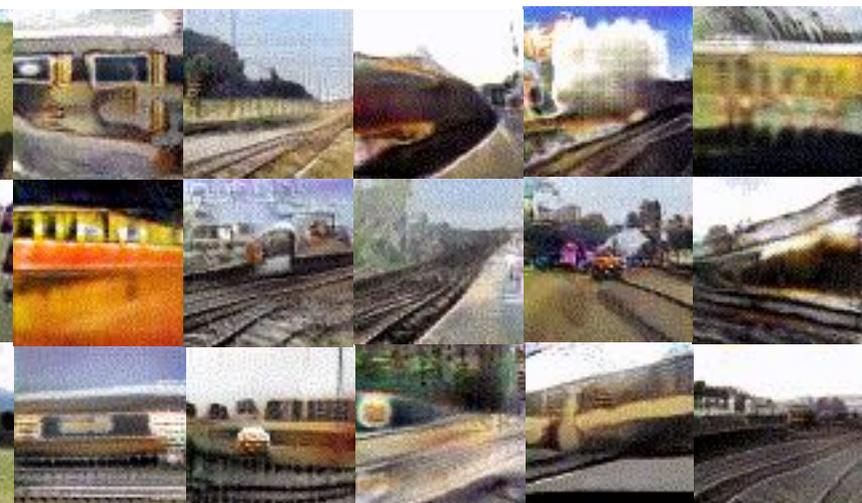
Beach



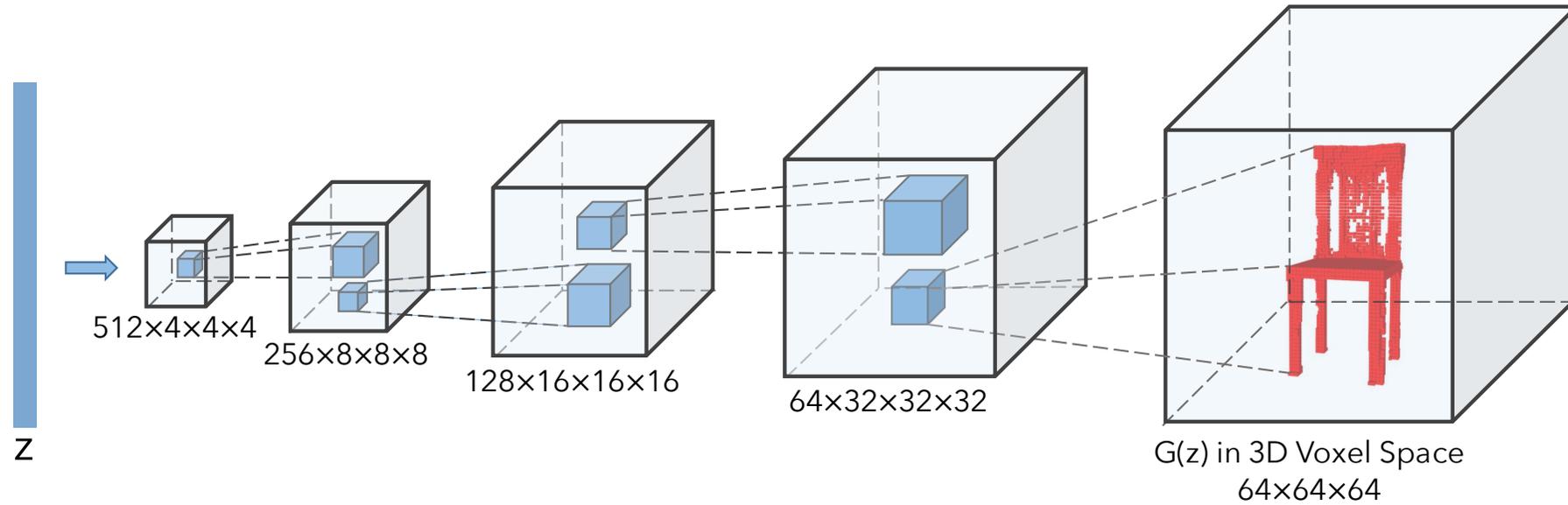
Golf



Train Station



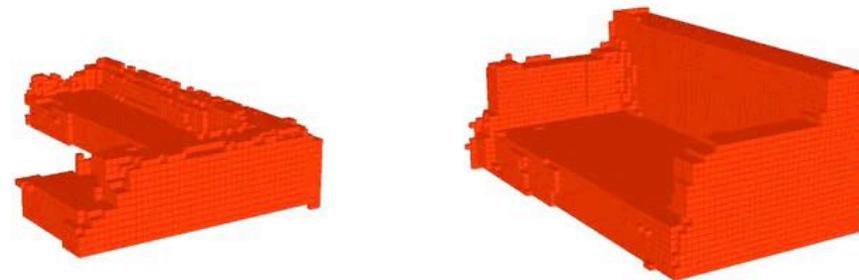
Generative Shape Modeling (Wu et al., 2016)



Chairs



Sofas



Text-to-Image Synthesis (Zhang et al., 2016)

The small bird has a red head with feathers that fade from red to gray from head to tail



The petals of this flower are white with a large stigma



A unique yellow flower with no visible pistils protruding from the center



This flower is pink and yellow in color, with petals that are oddly shaped



This is a light colored flower with many different petals on a green stem



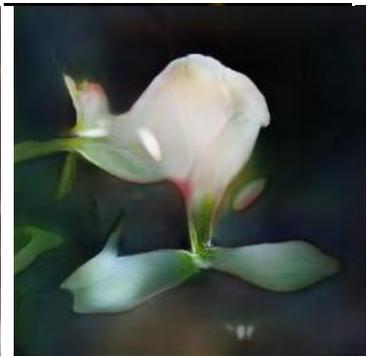
This flower is yellow and green in color, with petals that are ruffled



The flower have large petals that are pink with yellow on some of the petals



A flower that has white petals with some tones of yellow and green filaments



Single Image Super-Resolution (Ledig et al., 2016)

- Combine content loss with adversarial loss

bicubic



SRResNet



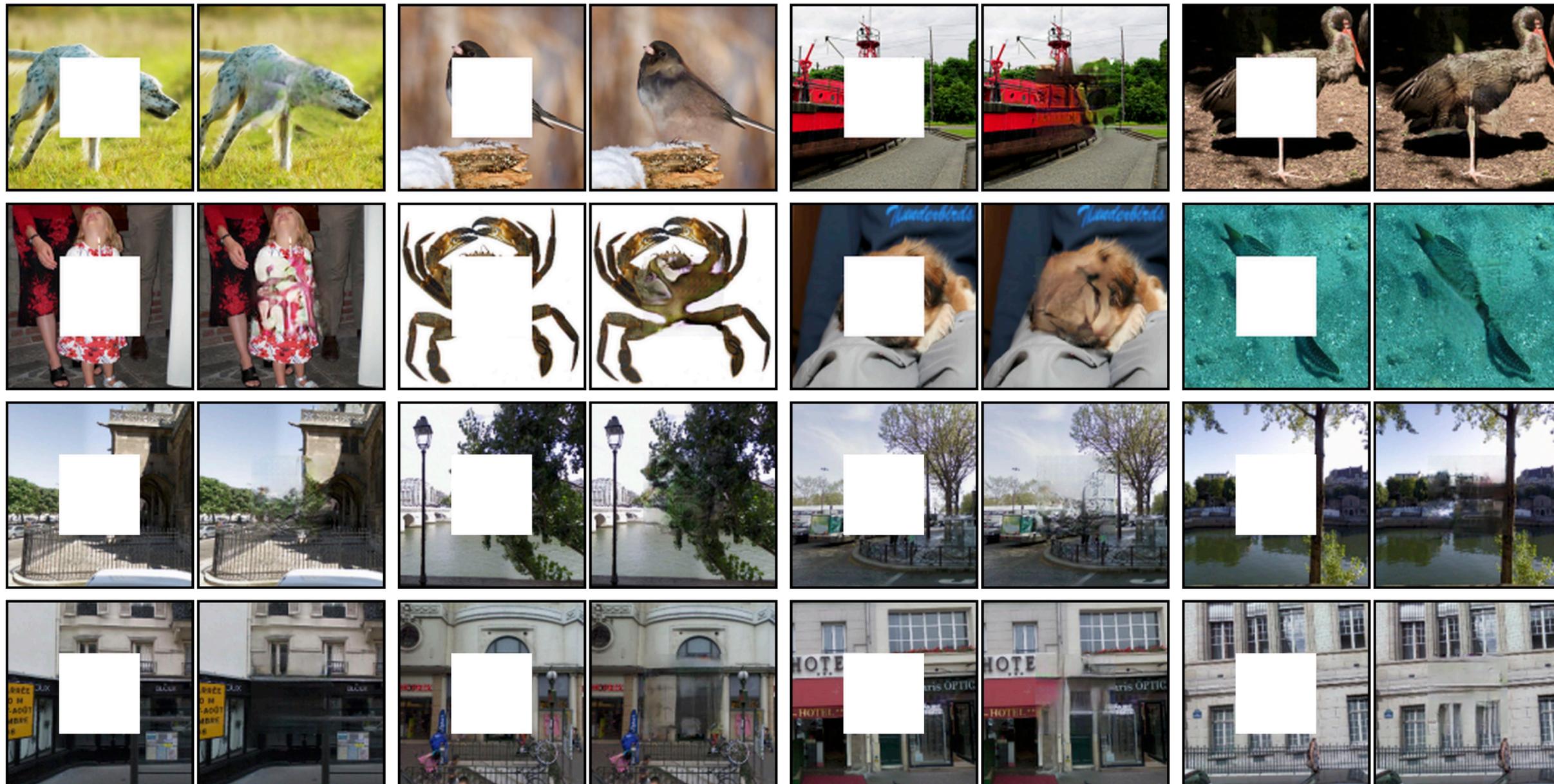
SRGAN



original



Image Inpainting (Pathak et al., 2016)



Unsupervised Domain Adaptation (Bousmalis et al., 2016)

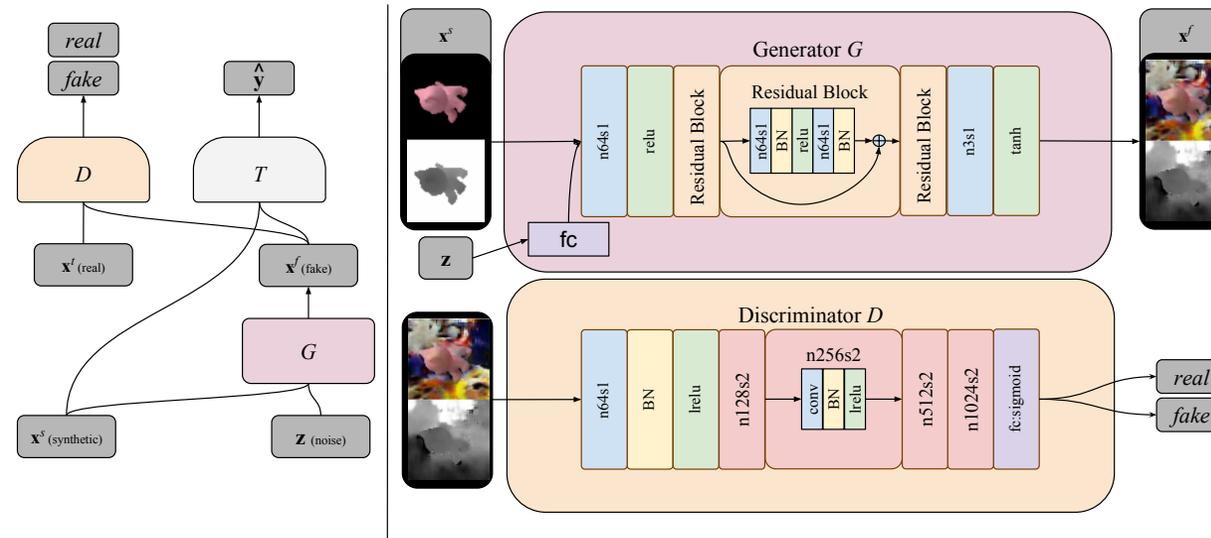


Image examples from the Linemod dataset



RGDB image samples
(conditioned on a synthetic image)

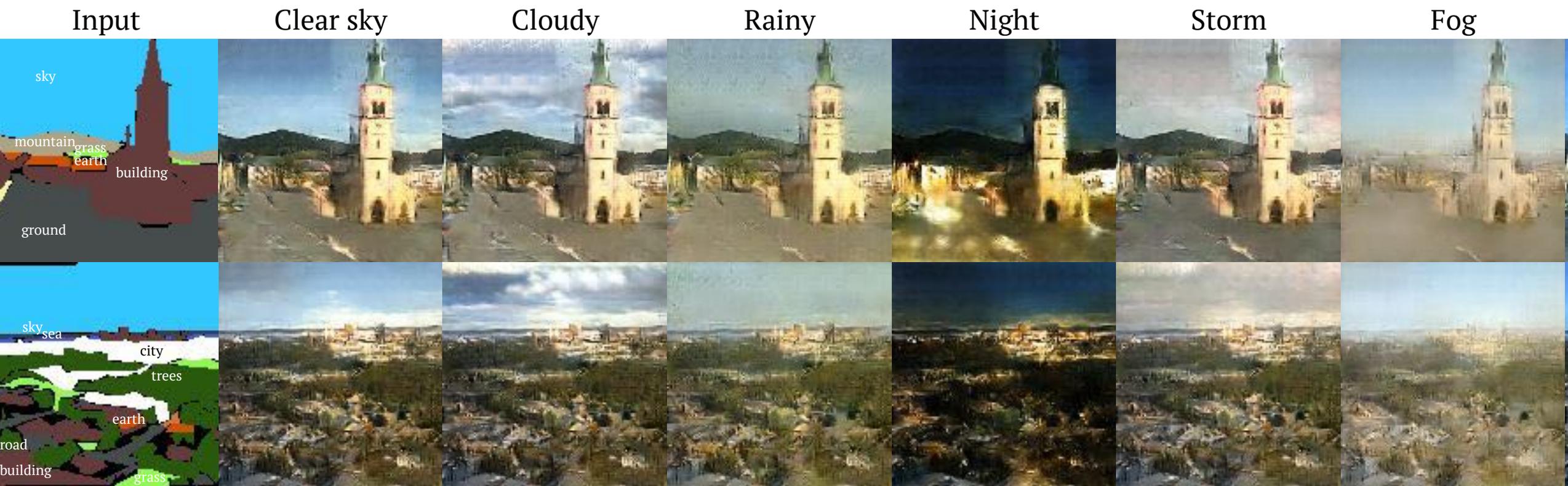


Image Editing

(Karacan et al., 2016)



“Maybe in our world lives a happy little tree
over there.”
— *Bob Ross*



How to Evaluate GANs?

Human Study

In this task, we present you computer generated pictures of outdoor scenes generated by different computer programs. Your task is to compare them and determine which is more realistic and natural looking. See the below table for some examples.



Steps

1. Analyze both images and consider their features carefully
2. Determine which computer generated image (image A or image B) is more realistic than the other.



A



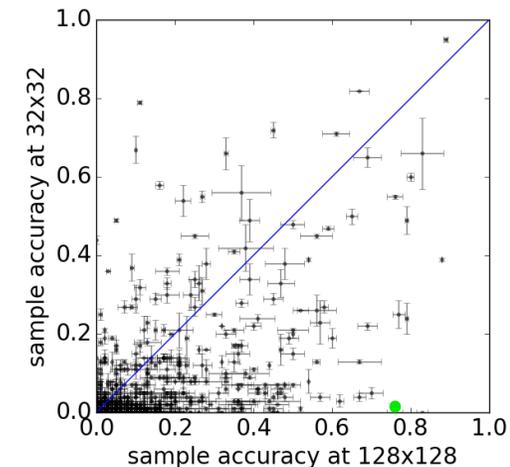
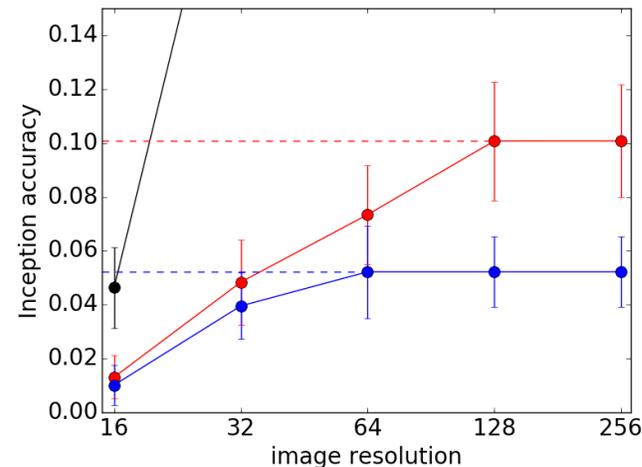
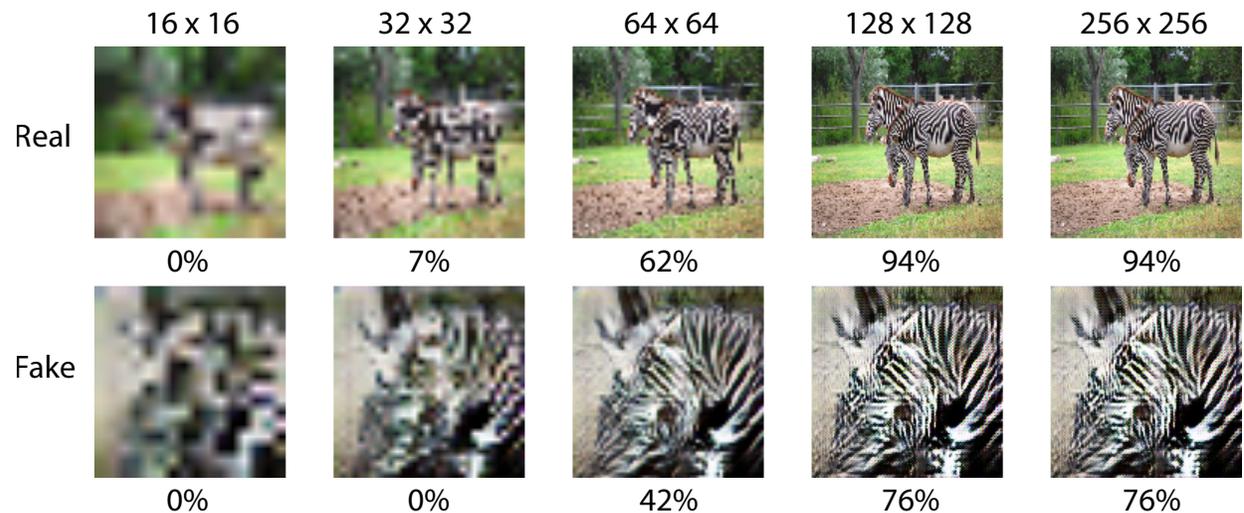
B

Which image seems more real?

- A
 B

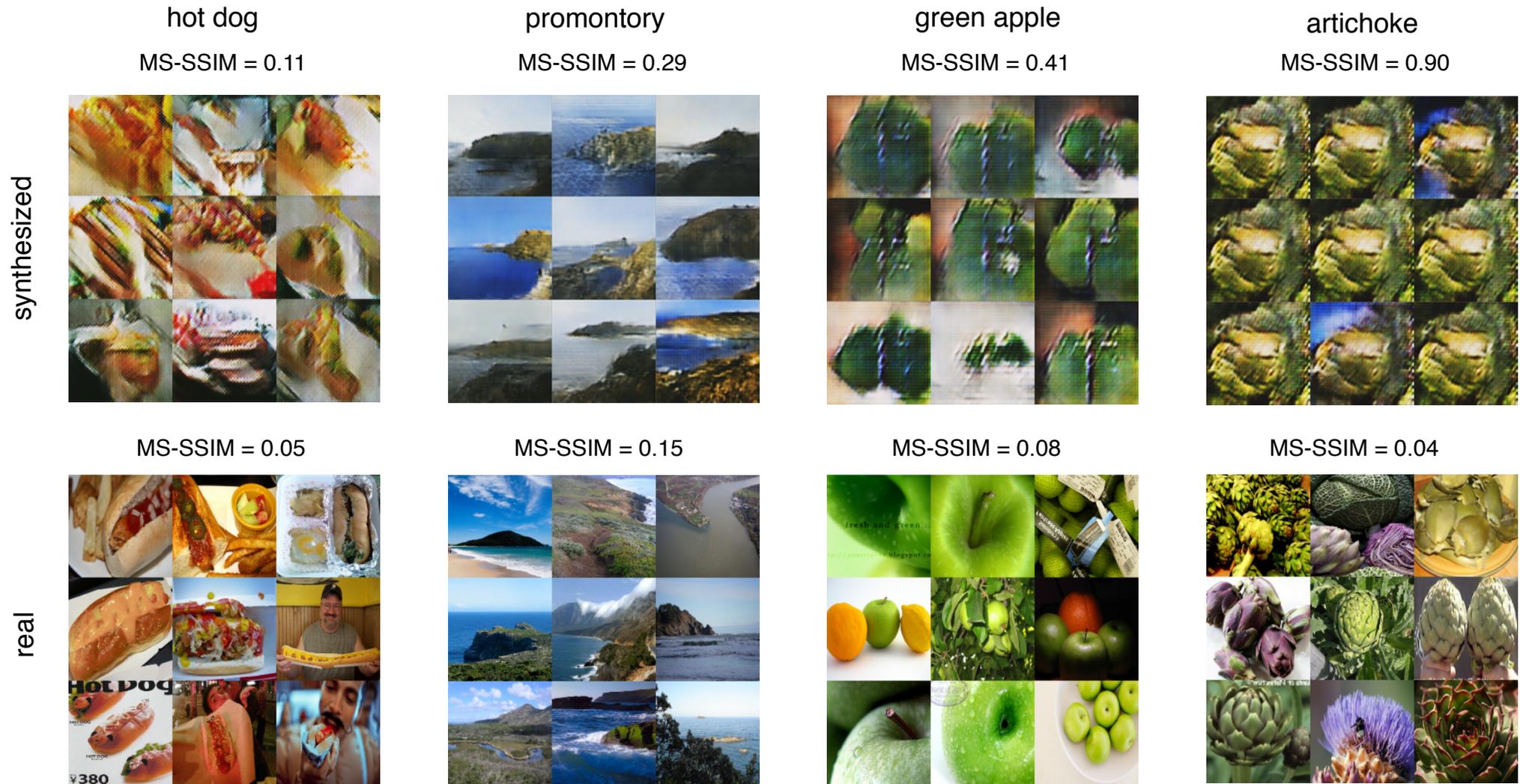
Evaluating Quality

- Hard to tell if progress is being made by looking at losses
- Inception Score (Salimans et al., 2016)
- Inception Accuracy (Odena et al., 2016)
 - Report the fraction of the samples for which the Inception network assigned the correct label



Measuring Diversity (Odena et al., 2016)

- MS-SSIM scores [between randomly chosen pairs of images within a given class]



Searching for Overfitting

- Nearest Neighbor Analysis (Odena et al., 2016)



Synthesized Samples

Corresponding Nearest Neighbors
In The Training Set

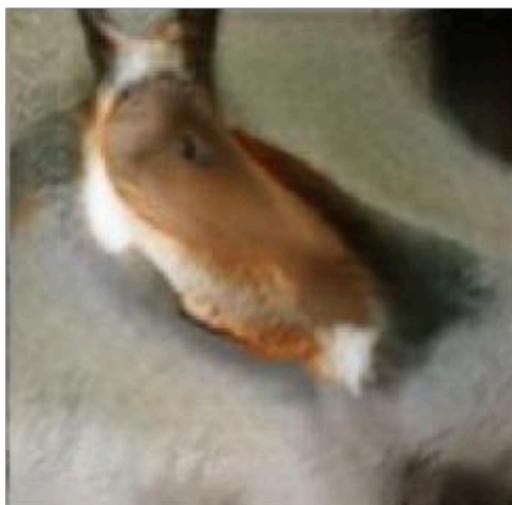
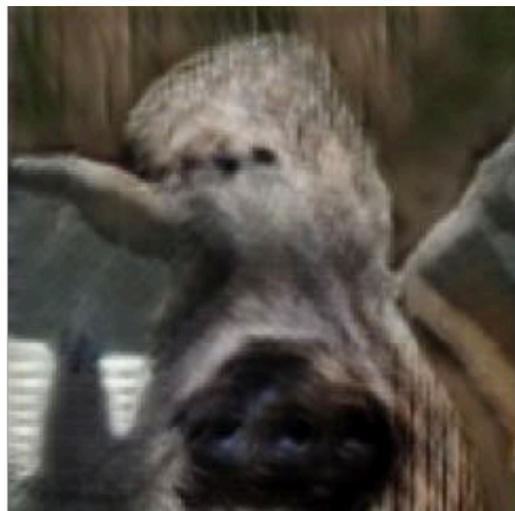
- Latent Space Interpolations

Image: (Dumoulin et al., 2016)

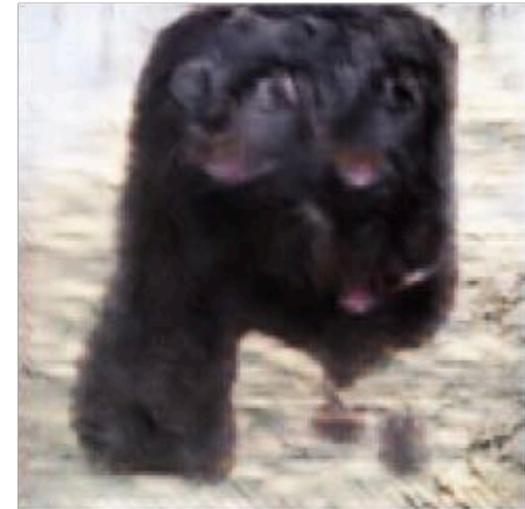


Limitations

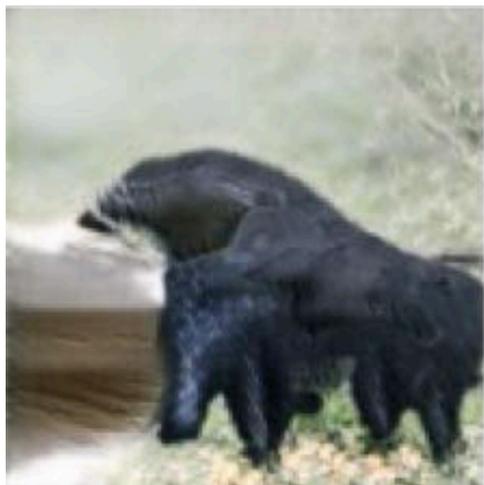
Cherry-Picked Results



Problems with Counting



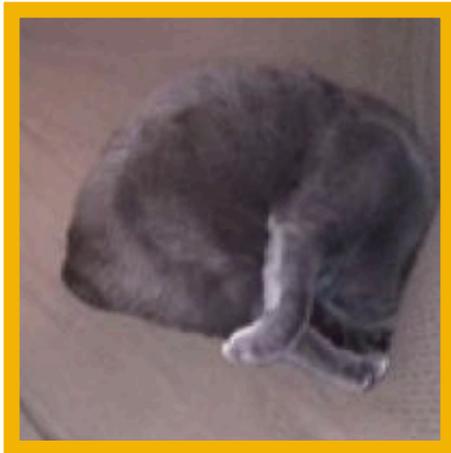
Problems with Perspective



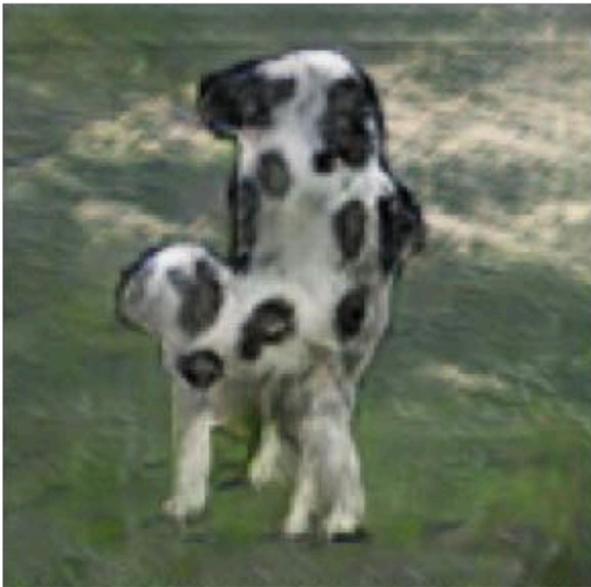
Problems with Perspective



This one
was real



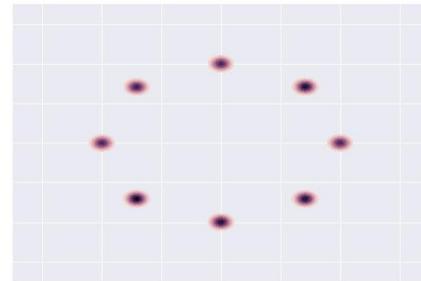
Problems with Global Structure



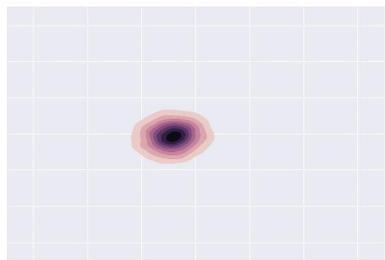
Mode Collapse (Metz et al., 2016)

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

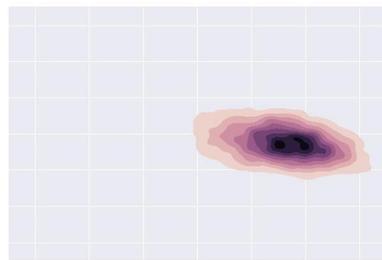
- D in inner loop: convergence to correct distribution
- G in inner loop: place all mass on most likely point



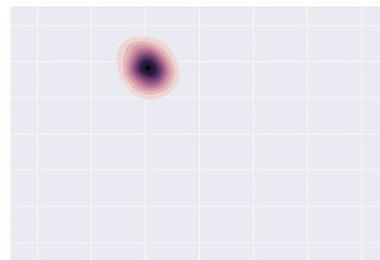
Target



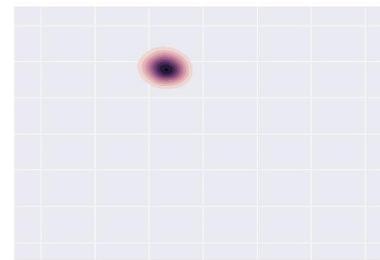
Step 0



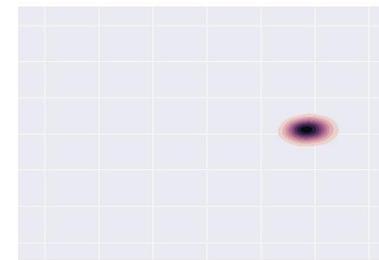
Step 5k



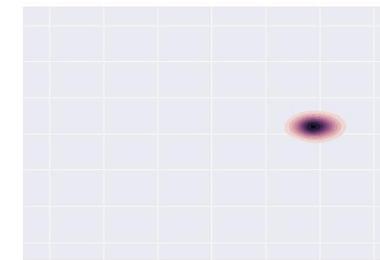
Step 10k



Step 15k



Step 20k



Step 25k

Mode collapse causes low output diversity

this small bird has a pink breast and crown, and black primaries and secondaries.



the flower has petals that are bright pinkish purple with white stigma



(Reed et al. 2016)

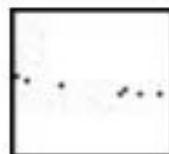
this magnificent fellow is almost all black with a red crest, and white cheek patch.



this white and yellow flower have thin white petals and a round yellow stamen



Key-points



GAN (Reed 2016b)

A man in a orange jacket with sunglasses and a hat ski down a hill.



This guy is in black trunks and swimming underwater.



A tennis player in a blue polo shirt is looking down at the green court.



This work



(Reed et al., 2017)

Non-convergence

- Optimization algorithms often approach a saddle point or local minimum rather than a global minimum
- Game solving algorithms may not approach an equilibrium at all

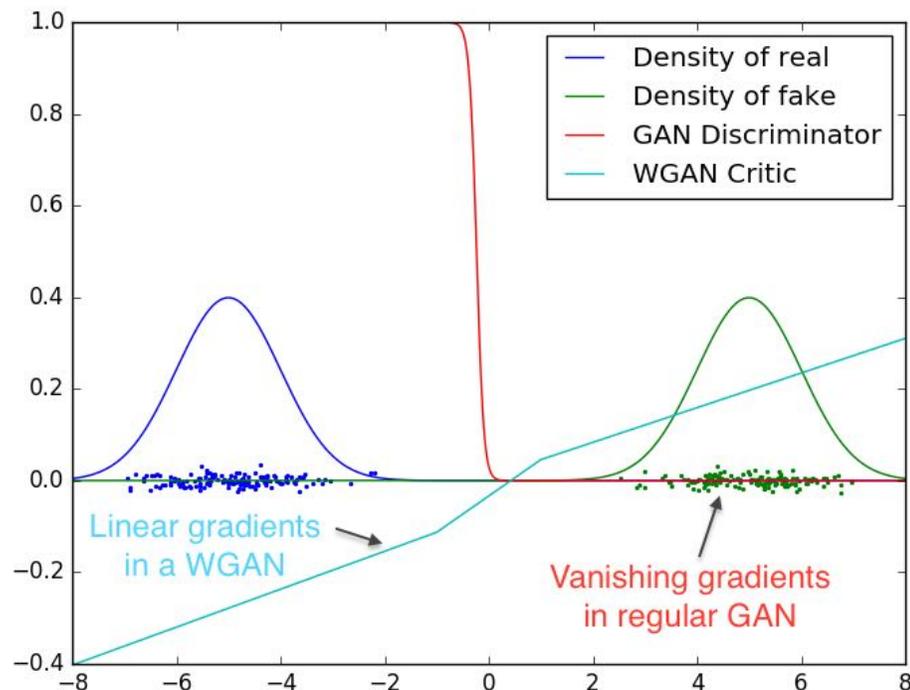
Frontiers

Wasserstein GAN (Arjovsky et al., 2016)

- Objective based on Earth-Mover or Wasserstein distance:

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [D_{\omega}(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [D_{\omega}(G_{\theta}(\mathbf{z}))]$$

- Provides nice gradients over real and fake samples



WGAN

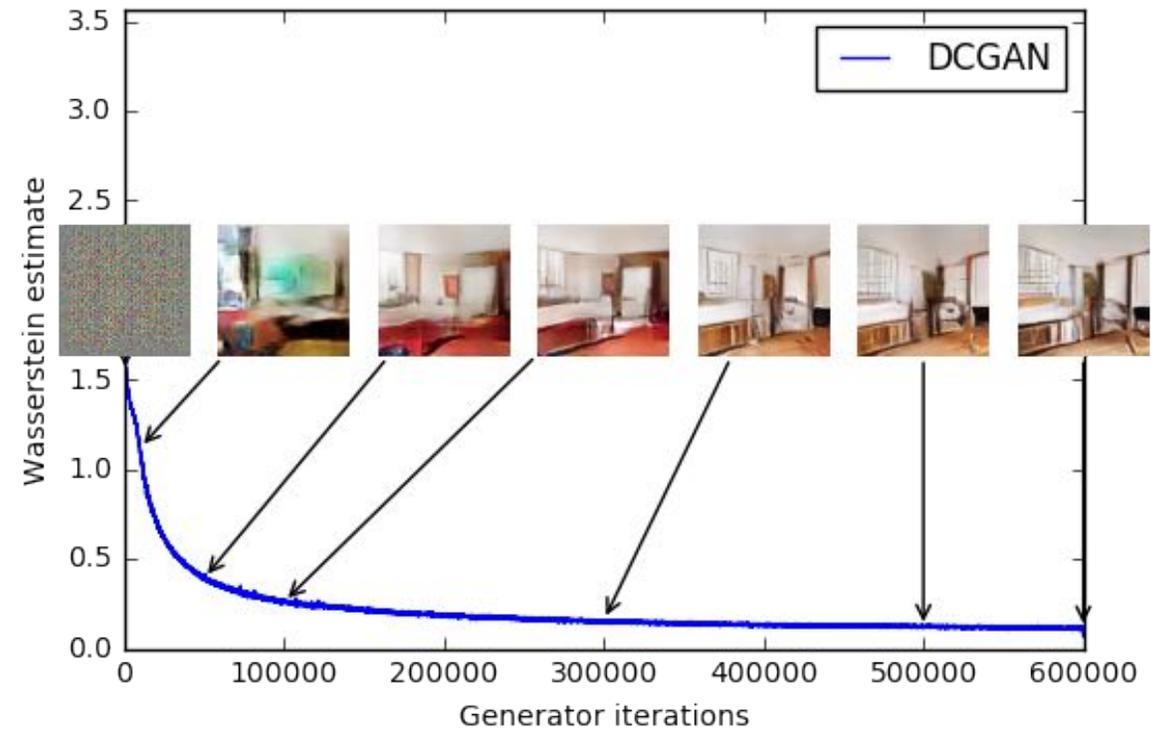
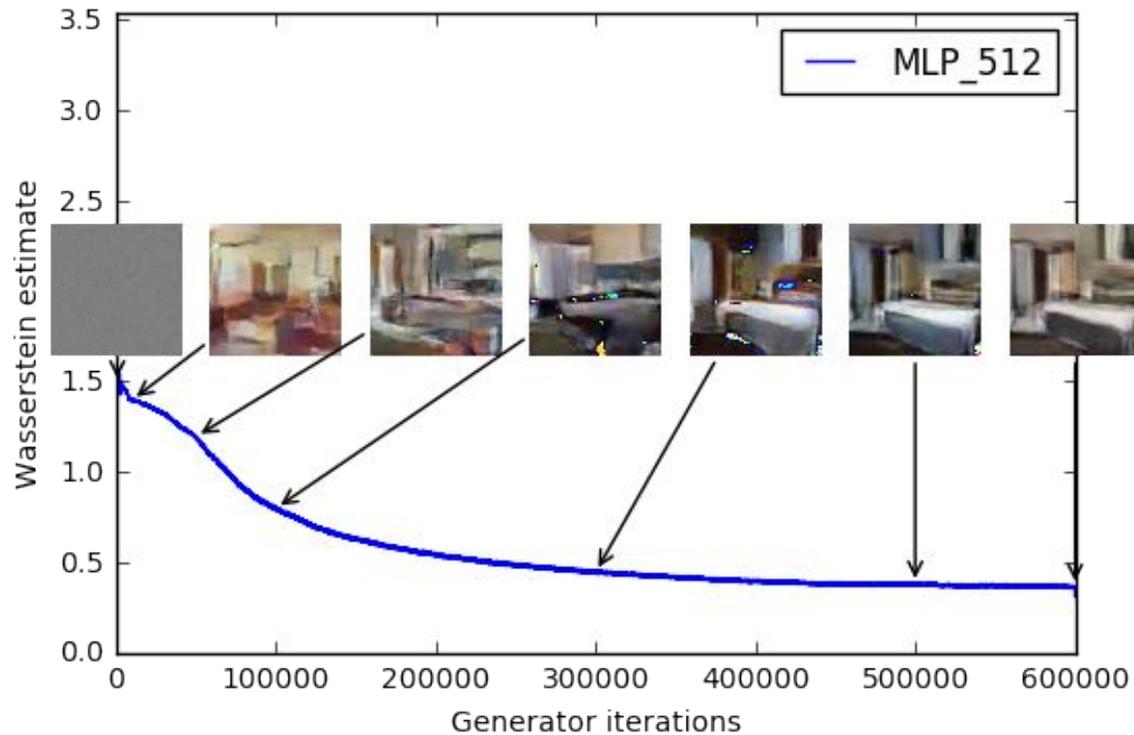


DCGAN



Wasserstein GAN (Arjovsky et al., 2016)

- Wasserstein loss seems to correlate well with image quality.



WGAN with gradient penalty (Gulraani et al., 2017)

$$L = \underbrace{\mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}} [(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2]}_{\text{Our gradient penalty}}$$

- Faster convergence and higher-quality samples than WGAN with weight clipping
- Train a wide variety of GAN architectures with almost no hyperparameter tuning, including discrete models

Samples from a character-level GAN language model on Google Billion Word

WGAN with gradient penalty

Busino game camperate spent odea
 In the bankaway of smarling the
 SingersMay , who kill that invic
 Keray Pents of the same Reagun D
 Manging include a tudancs shat "
 His Zuith Dudget , the Denmbern
 In during the Uitational questio
 Divos from The ' noth ronkies of
 She like Monday , of macunsuer S
 The investor used ty the present
 A papees are country congress oo
 A few year inom the group that s
 He said this syenn said they wan
 As a world 1 88 ,for Autouries
 Foand , th Word people car , Il
 High of the upseader homing pull
 The guipe is worly move dogsfor
 The 1874 incidested he could be
 The allo tooks to security and c

Solice Norkedin pring in since
 ThiS record (31.) UBS) and Ch
 It was not the annuas were plogr
 This will be us , the ect of DAN
 These leaded as most-worsd p2 a0
 The time I paid0a South Cubry i
 Dour Fraps higs it was these del
 This year out howneed allowed lo
 Kaulna Seto consficutes to repor
 A can teal , he was schoon news
 In th 200. Pesish picriers rega
 Konney Panice rimimber the teami
 The new centuct cut Denester of
 The near , had been one injustie
 The incestion to week to shorted
 The company the high product of
 20 - The time of accomplete , wh
 John WVuderenson seqiivic spends
 A ceetens in indestedly the Wat

Standard GAN objective

dddddddddddddddddddddddddddddd
 ddddddddddddddddddddddddddd

dddddddddddddddddddddddddd
 ddddddddddddddddddddddd

Boundary Equilibrium GAN (BEGAN)

(Berthelot et al., 2017)

- A loss derived from the Wasserstein distance for training auto-encoder based GANs

$$\mathcal{L}(v) = |v - D(v)|^\eta \text{ where } \begin{cases} D : \mathbb{R}^{N_x} \mapsto \mathbb{R}^{N_x} & \text{is the autoencoder function.} \\ \eta \in \{1, 2\} & \text{is the target norm.} \\ v \in \mathbb{R}^{N_x} & \text{is a sample of dimension } N_x. \end{cases}$$

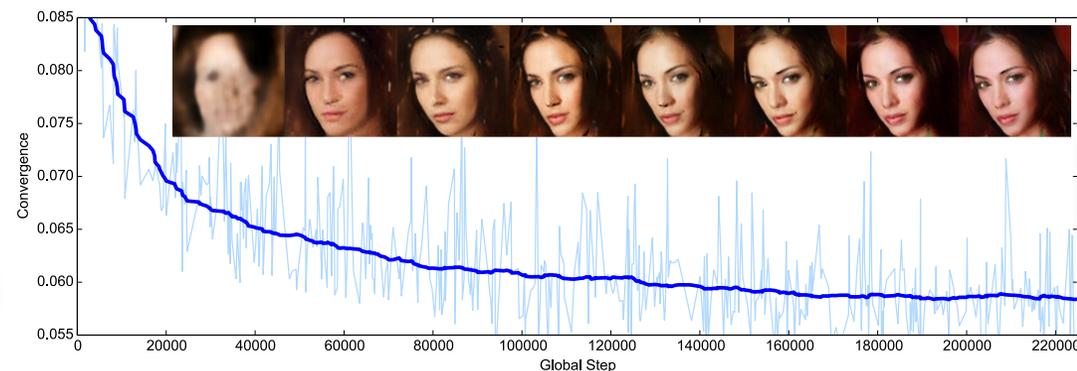
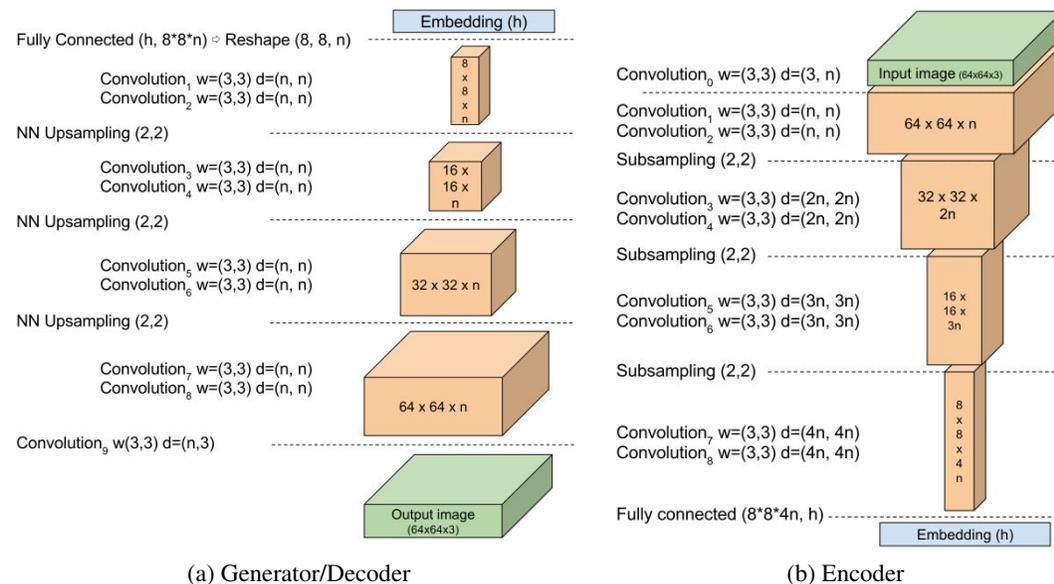
- Wasserstein distance btw. the reconstruction losses of real and generated data

- Convergence measure:

$$\mathcal{M}_{global} = \mathcal{L}(x) + |\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))|$$

- Objective:

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x) - k_t \cdot \mathcal{L}(G(z_D)) & \text{for } \theta_D \\ \mathcal{L}_G = \mathcal{L}(G(z_G)) & \text{for } \theta_G \\ k_{t+1} = k_t + \lambda_k (\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))) & \text{for each training step } t \end{cases}$$



BEGANs for CelebA

360K celebrity face images
128x128 with 128 filters

(Berthelot et al., 2017)



Interpolations in the latent space



Mirror interpolation example