

Introduction to Generative Adversarial Network (GAN)

Hongsheng Li

Department of Electronic Engineering

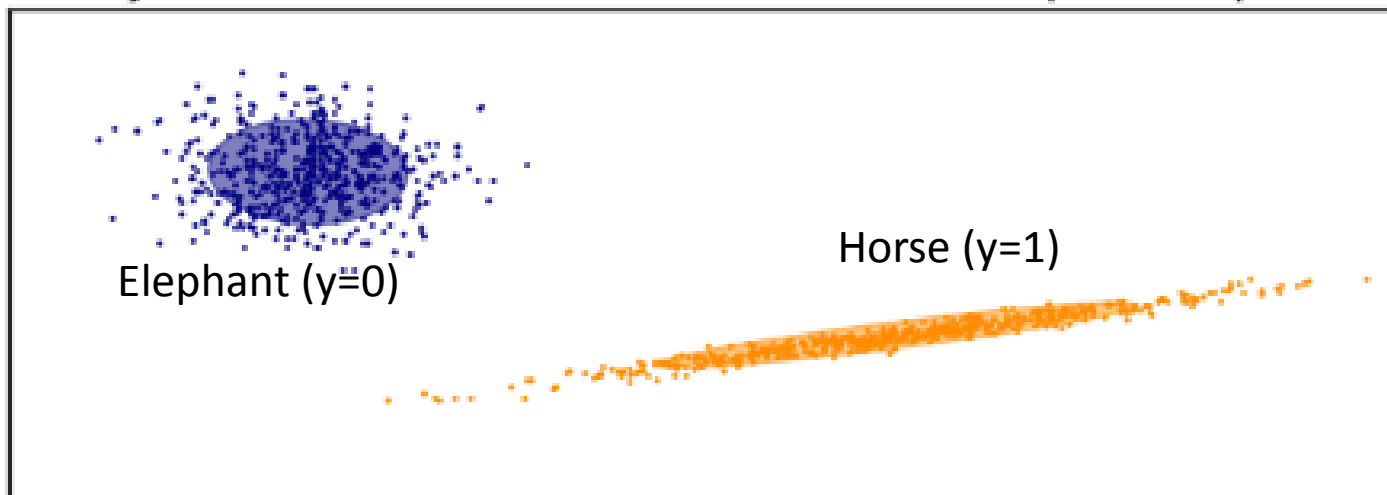
Chinese University of Hong Kong

Generative Models

- Density Estimation
 - Discriminative model: $p(y | x)$
 - $y=0$ for elephant, $y=1$ for horse
 - Generative model: $p(x | y)$

$$p(x | y = 0)$$

$$p(x | y = 1)$$



Generative Models

- $$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$
$$\arg \max_y p(y|x) = \arg \max_y \frac{p(x|y)p(y)}{p(x)}$$
$$= \arg \max_y p(x|y)p(y).$$

- Sample Generation



Generative Models

- Sample Generation

Training samples



Model samples

Training samples



Generative Models

- Generative model



- GAN is a generative model
 - Mainly focuses on sample generation
 - Possible to do both

Why Worth Studying?

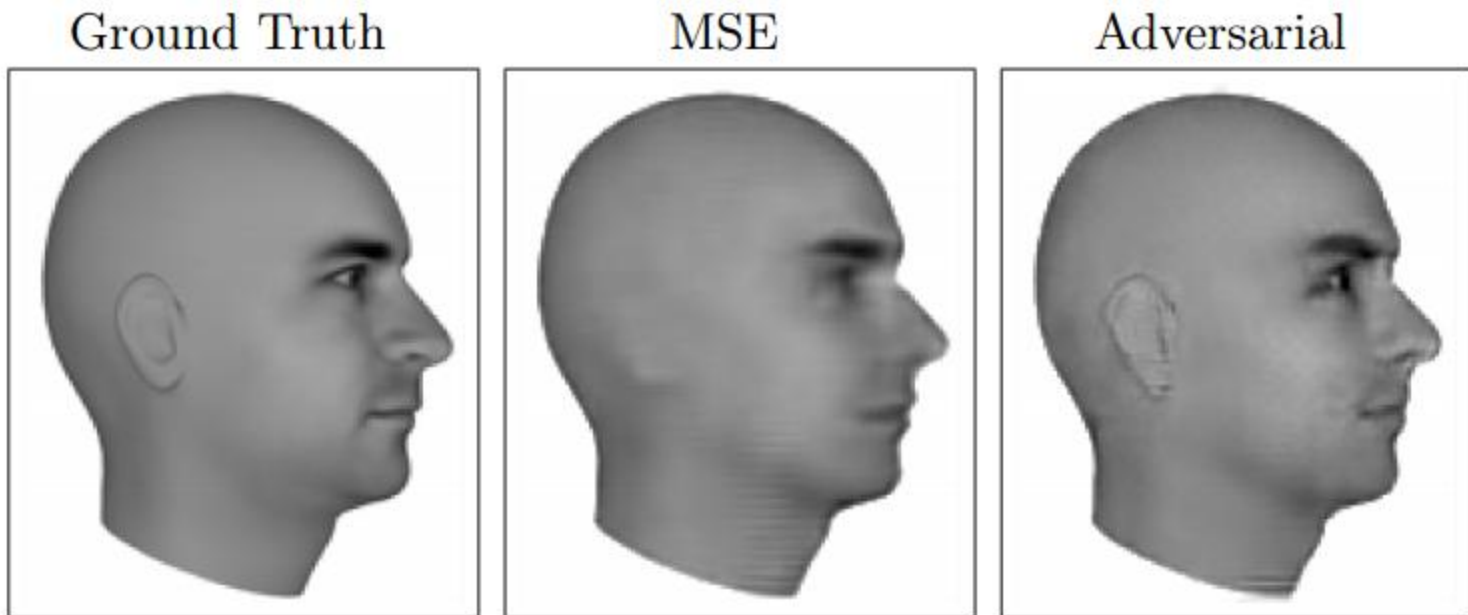
- Excellent test of our ability to use high-dimensional, complicated probability distributions

P_{model}  Sample generation

- Missing data
 - Semi-supervised learning

Why Worth Studying?

- Multi-modal outputs
 - Example: next frame prediction



Lotter et al. 2015

Why Worth Studying?

- Image generation tasks
 - Example: single-image super-resolution

original



bicubic
(21.59dB/0.6423)



SRResNet
(23.44dB/0.7777)



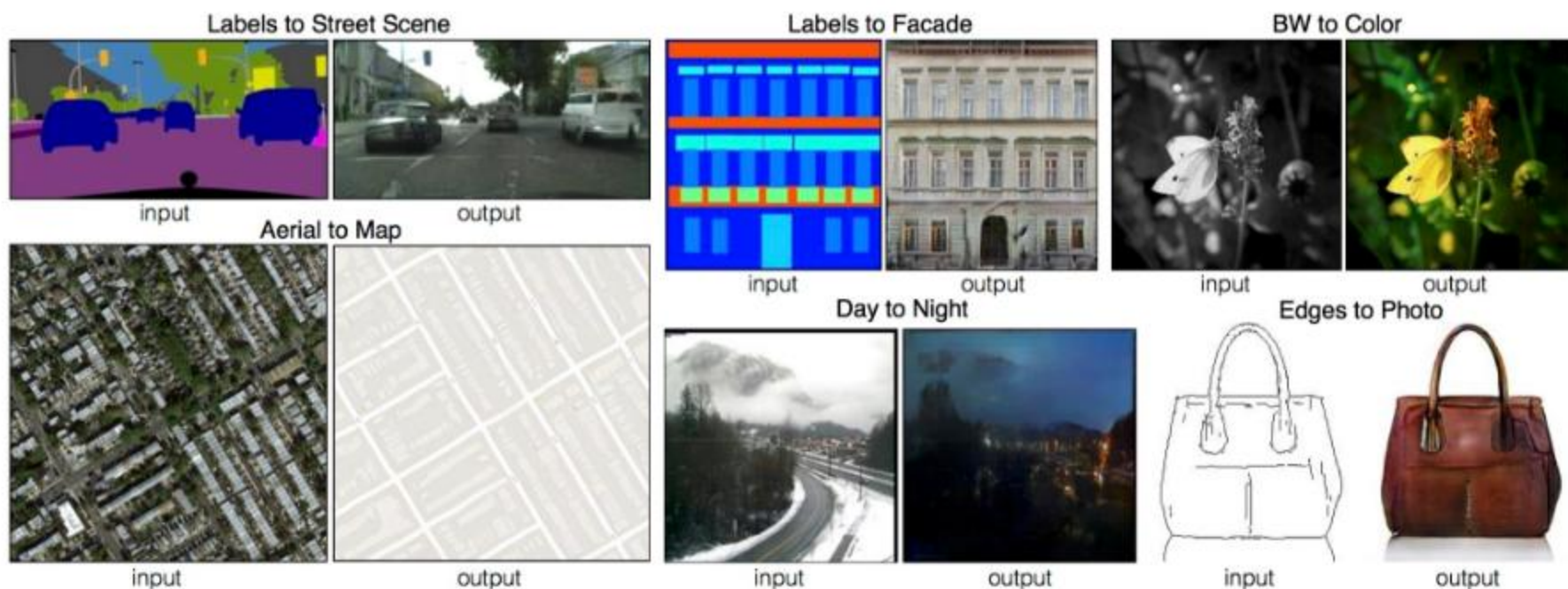
SRGAN
(20.34dB/0.6562)



Ledig et al 2015

Why Worth Studying?

- Image generation tasks
 - Example: Image-to-Image Translation
 - <https://affinelayer.com/pixsrv/>



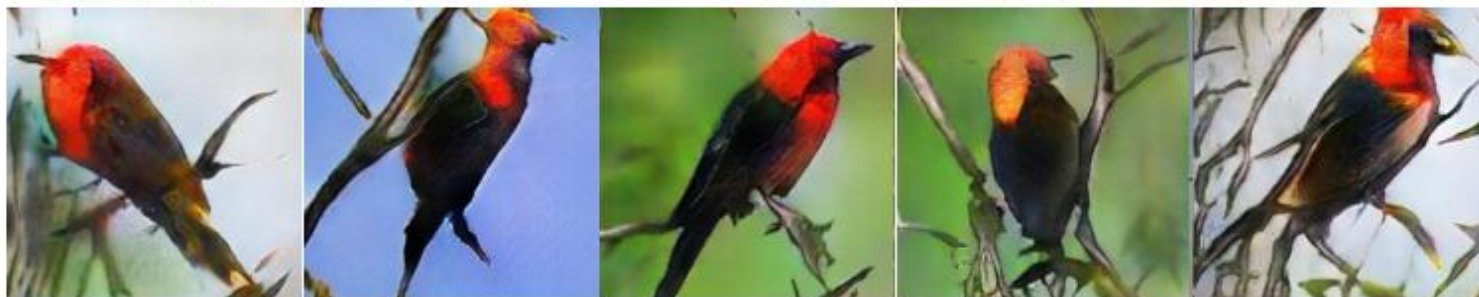
Why Worth Studying?

- Image generation tasks
 - Example: Text-to-Image Generation

This small blue bird has a short pointy beak and brown on its wings



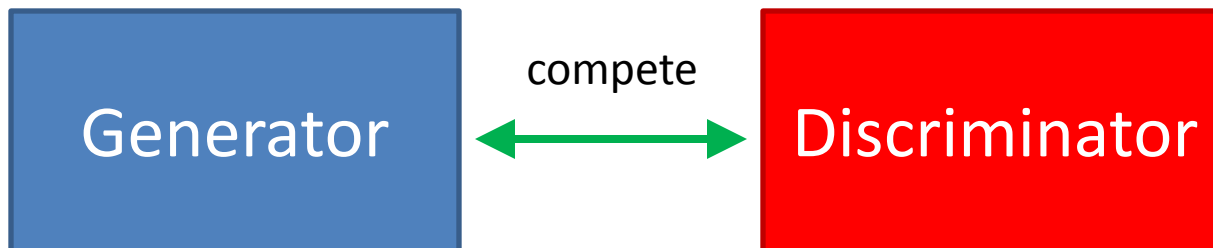
This bird is completely red with black wings and pointy beak



Zhang et al 2016

How does GAN Work?

- Adversarial – adj. 對抗的
- Two networks:
 - Generator G : creates (fake) samples that the discriminator cannot distinguish
 - Discriminator D : determine whether samples are fake or real

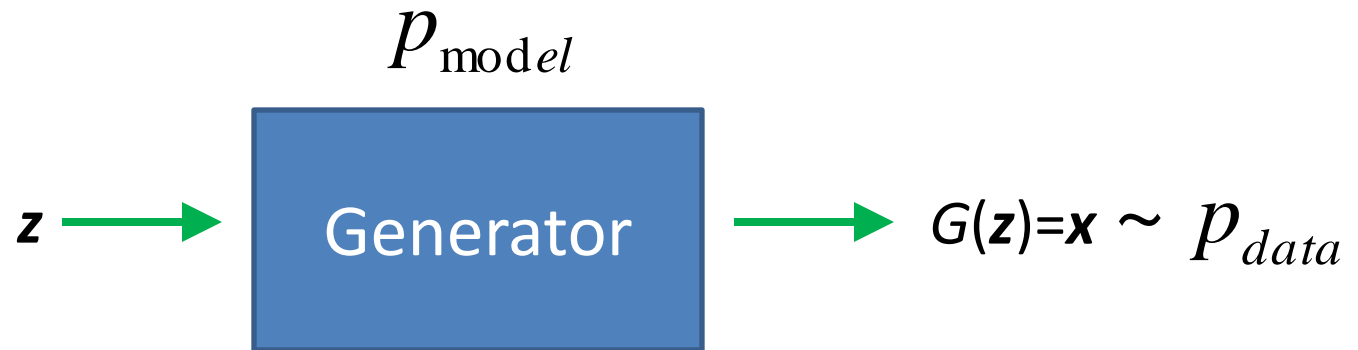


The Generator

- G : a differentiable function
 - modeled as a neural network
- Input:
 - \mathbf{z} : random noise vector from some simple prior distribution
- Output:
 - $\mathbf{x} = G(\mathbf{z})$: generated samples



The Generator



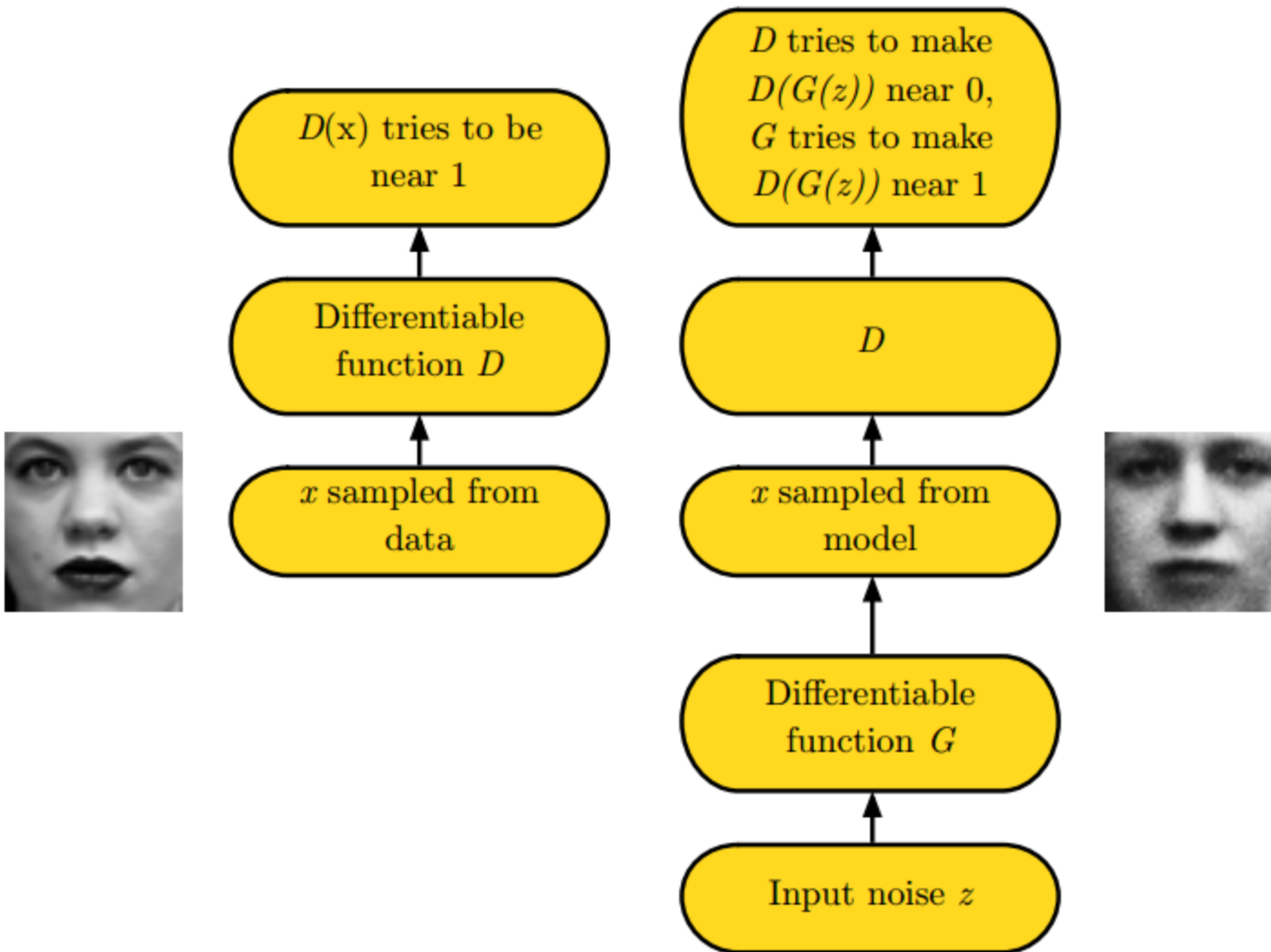
- The dimension of \mathbf{z} should be at least as large as that of \mathbf{x}

The Discriminator

- D : modeled as a neural network
- Input:
 - Real sample
 - Generated sample x
- Output:
 - 1 for real samples
 - 0 for fake samples



Generative Adversarial Networks



Cost Functions

- The discriminator outputs a value $D(x)$ indicating the chance that x is a real image
- For real images, their ground-truth labels are 1. For generated images their labels are 0.
- Our objective is to maximize the chance to recognize real images as real and generated images as fake
- The objective for generator can be defined as

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

recognize real images better

recognize generated images better

Cost Functions

- For the generator G , its objective function wants the model to generate images with the highest possible value of $D(x)$ to fool the discriminator
- The cost function is

$$\min_G V(G) = \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

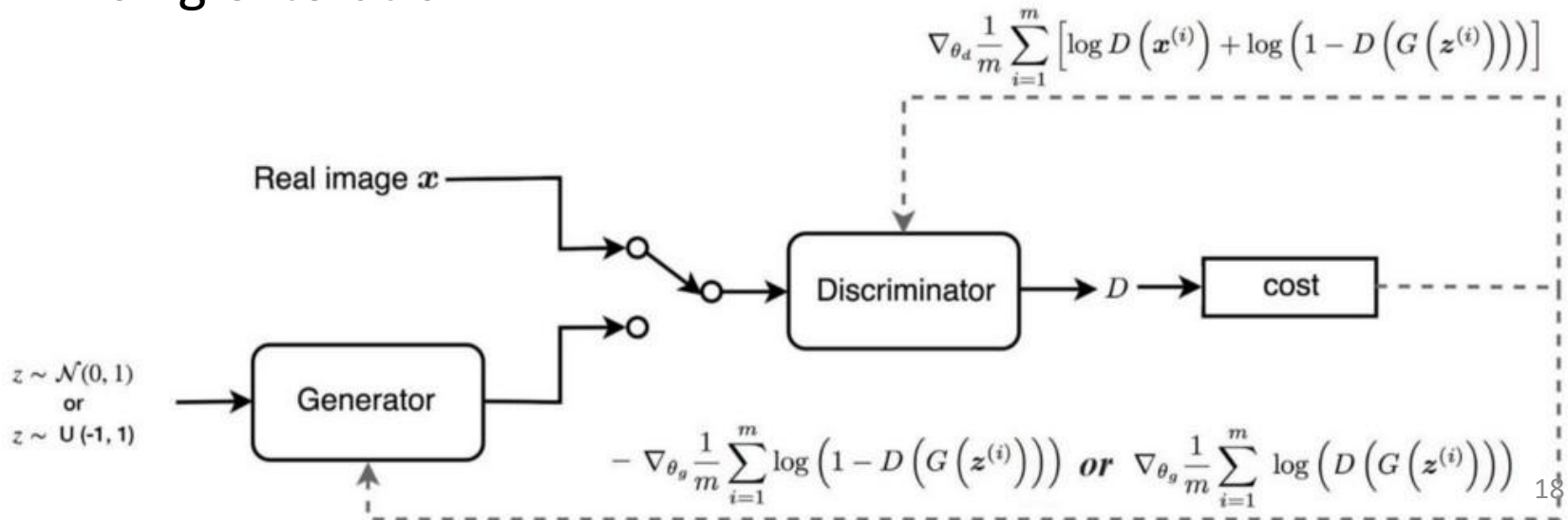
Optimize G that can fool the discriminator the most.

- The overall GAN training is therefore a min-max game

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

Training Procedure

- The generator and the discriminator are learned jointly by the alternating gradient descent
 - Fix the generator's parameters and perform a single iteration of gradient descent on the discriminator using the real and the generated images
 - Fix the discriminator and train the generator for another single iteration



The Algorithm

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log \left(1 - D(G(\mathbf{z}^{(i)})) \right) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

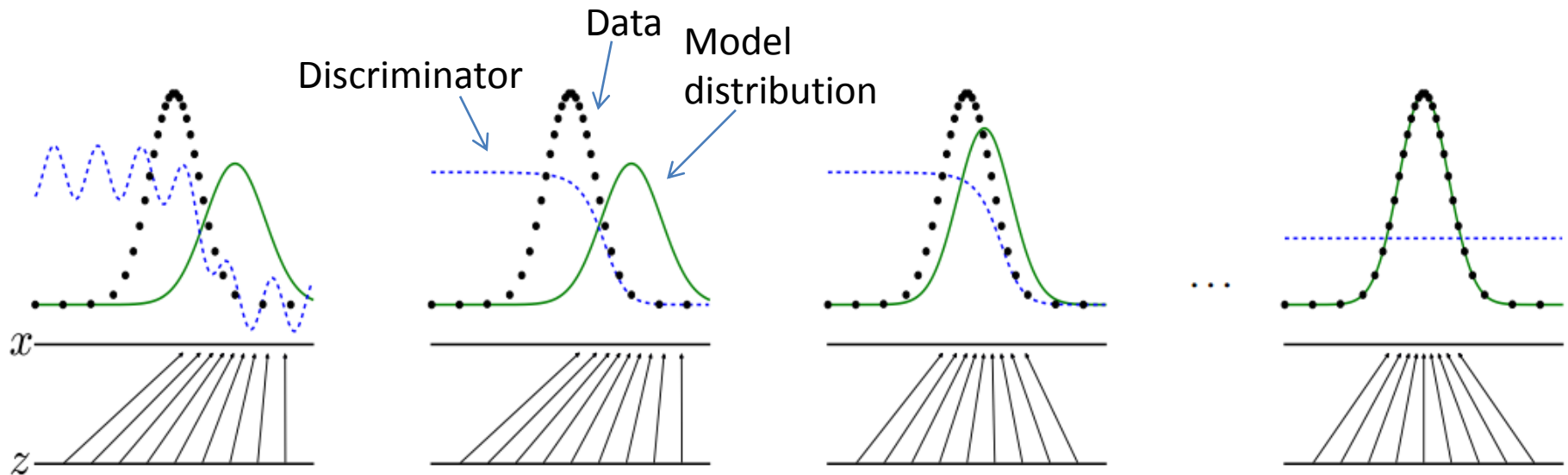
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D(G(\mathbf{z}^{(i)})) \right).$$

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Illustration of the Learning

- Generative adversarial learning aims to learn a model distribution that matches the actual data distribution

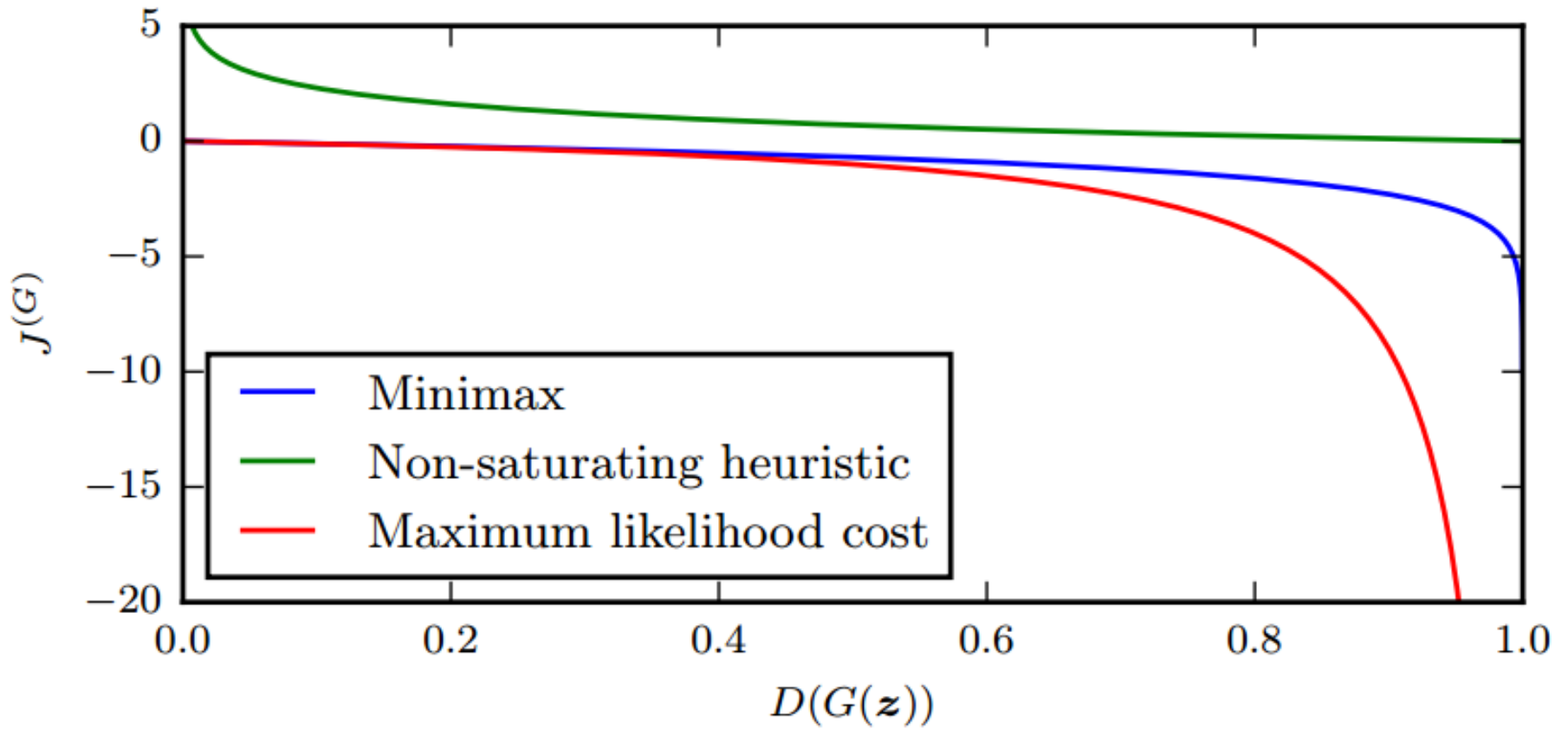


Generator diminished gradient

- However, we encounter a gradient diminishing problem for the generator. The discriminator usually wins early against the generator
- It is always easier to distinguish the generated images from real images in early training. That makes *cost function* approaches 0. i.e. $-\log(1 - D(G(z))) \rightarrow 0$
- The gradient for the generator will also vanish which makes the gradient descent optimization very slow
- To improve that, the GAN provides an alternative function to backpropagate the gradient to the generator

$$\begin{array}{ccc} \textit{minimize} & & \textit{maximize} \\ -\nabla_{\theta_g} \log \left(1 - D \left(G \left(z^{(i)} \right) \right) \right) \rightarrow 0 & \textit{change to} & \nabla_{\theta_g} \log \left(D \left(G \left(z^{(i)} \right) \right) \right) \end{array}$$

Comparison between Two Losses

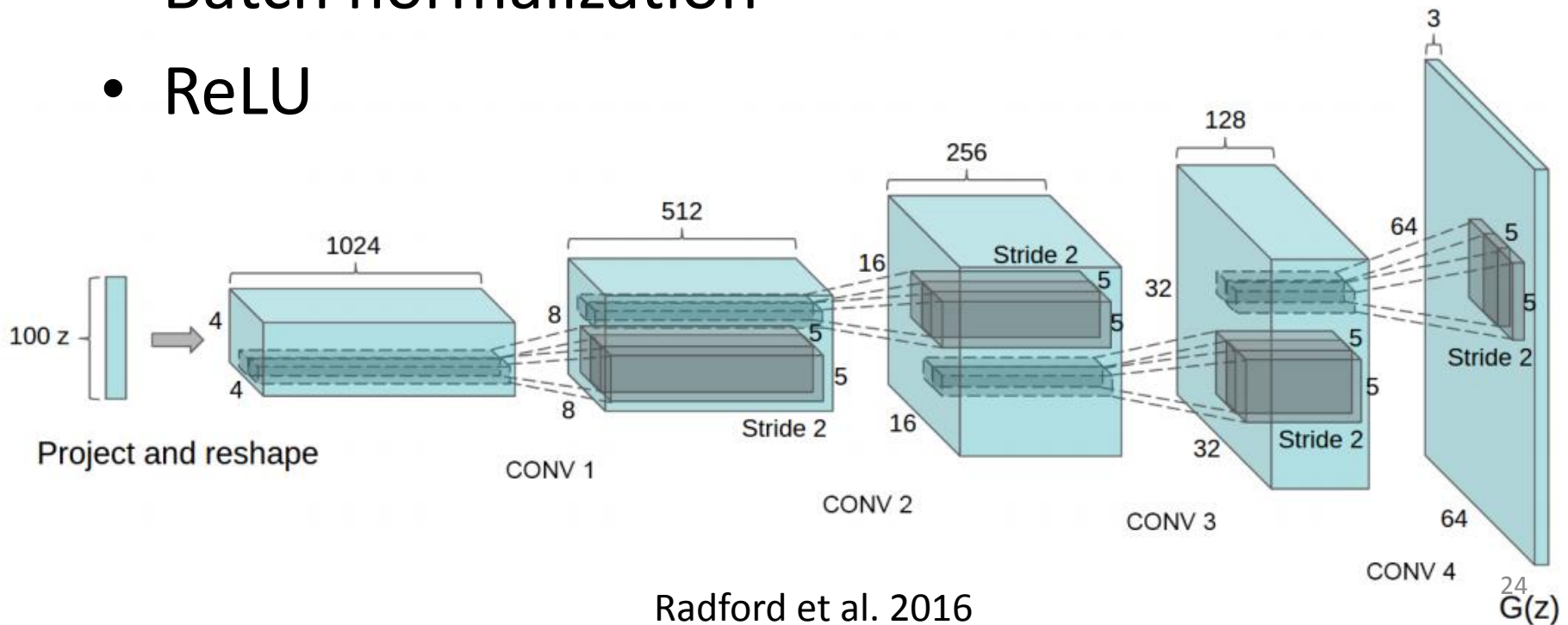


Non-Saturating Game

- $J^{(D)} = -\frac{1}{2}\mathbb{E}_{\mathbf{x}\sim p_{\text{data}}}\log D(\mathbf{x}) - \frac{1}{2}\mathbb{E}_{\mathbf{z}}\log(1 - D(G(\mathbf{z})))$
- $J^{(G)} = -\frac{1}{2}\mathbb{E}_{\mathbf{z}}\log D(G(\mathbf{z}))$
- In the min-max game, the generator maximizes the same cross-entropy
- Now, generator maximizes the log-probability of the discriminator being mistaken
- Heuristically motivated; generator can still learn even when discriminator successfully rejects all generator samples

Deep Convolutional Generative Adversarial Networks (DCGAN)

- All convolutional nets
- No global average pooling
- Batch normalization
- ReLU

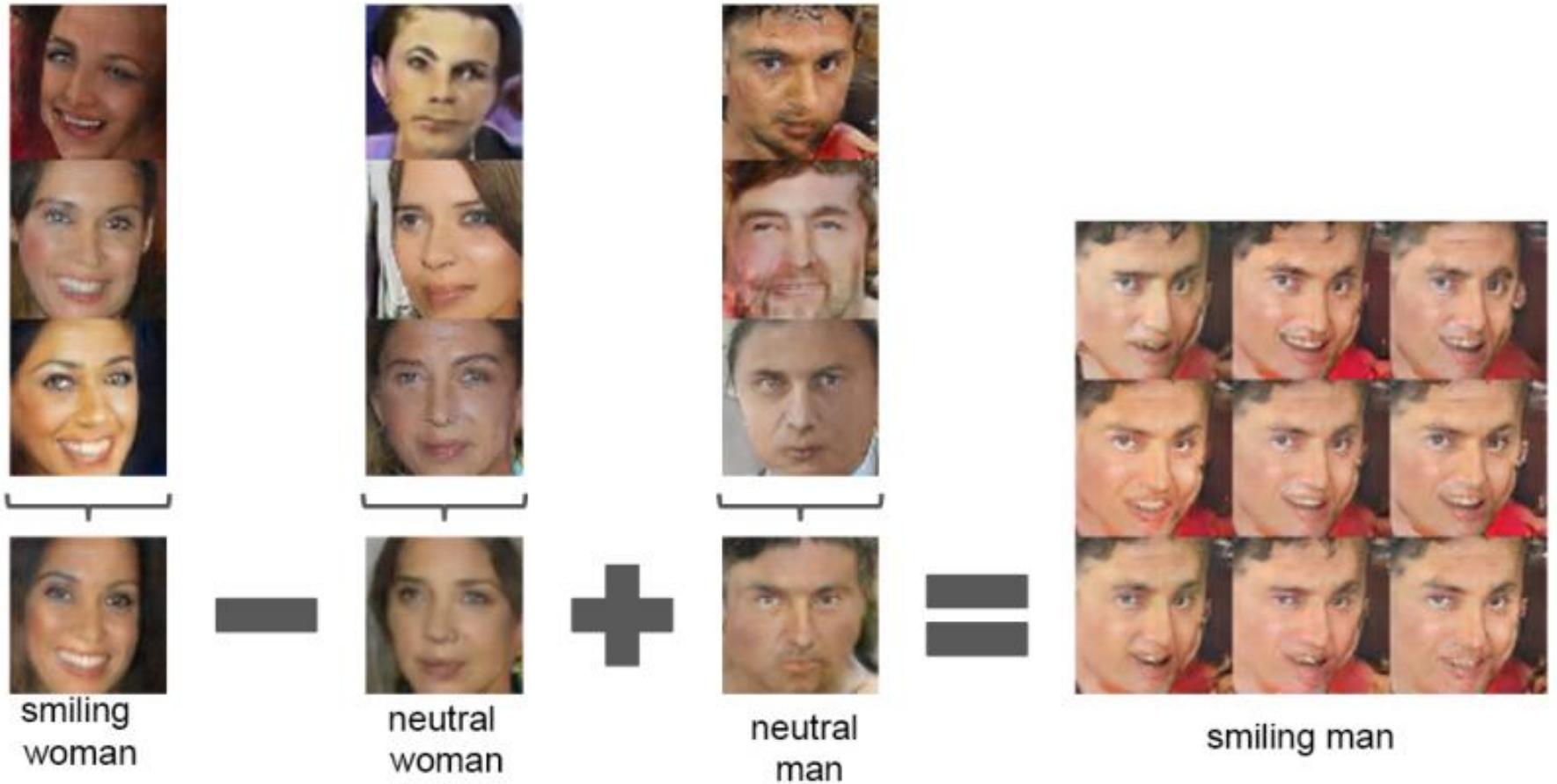


Deep Convolutional Generative Adversarial Networks (DCGAN)

- LSUN bedroom (about 3m training images)



Manipulating Learned z



Manipulating Learned z

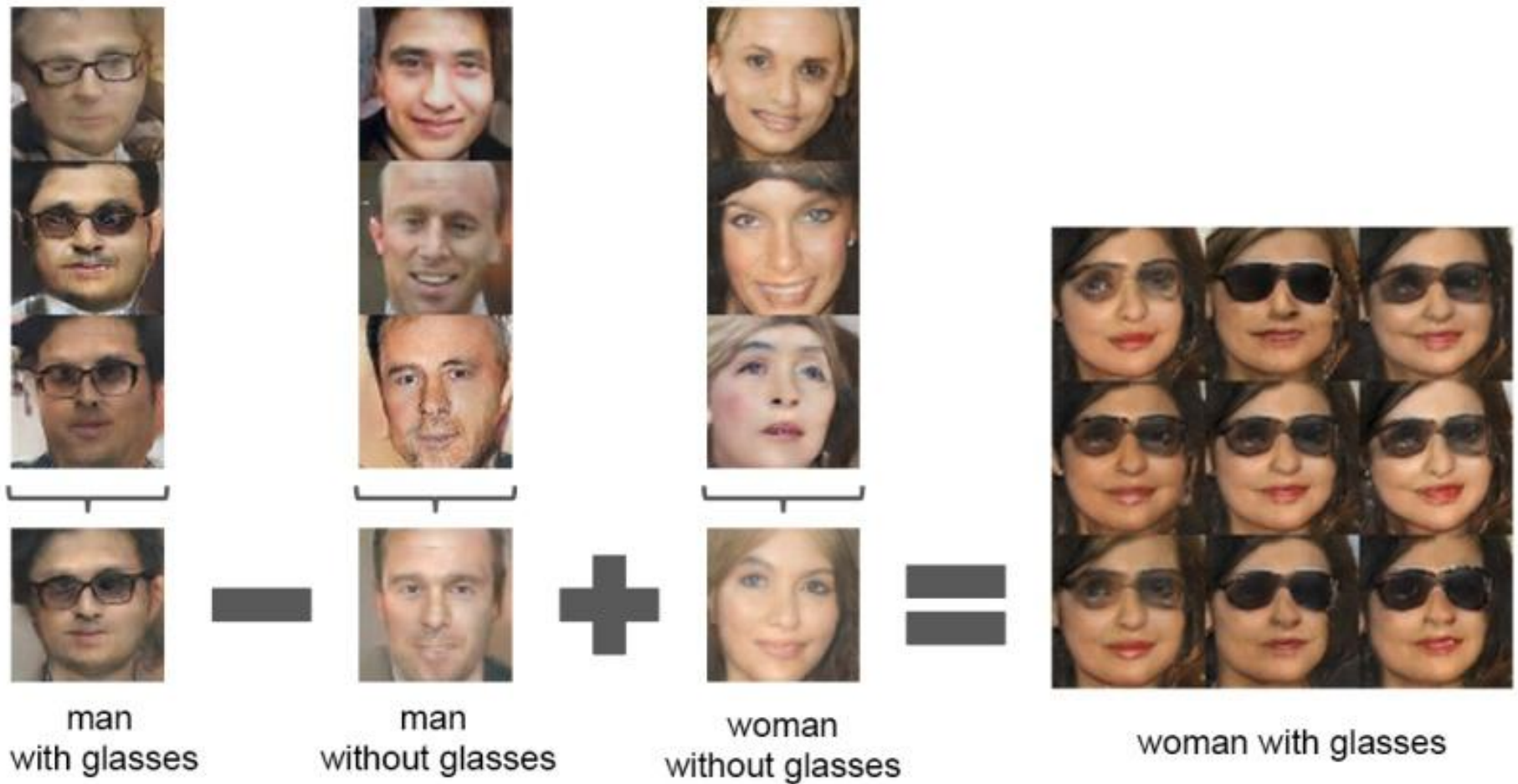
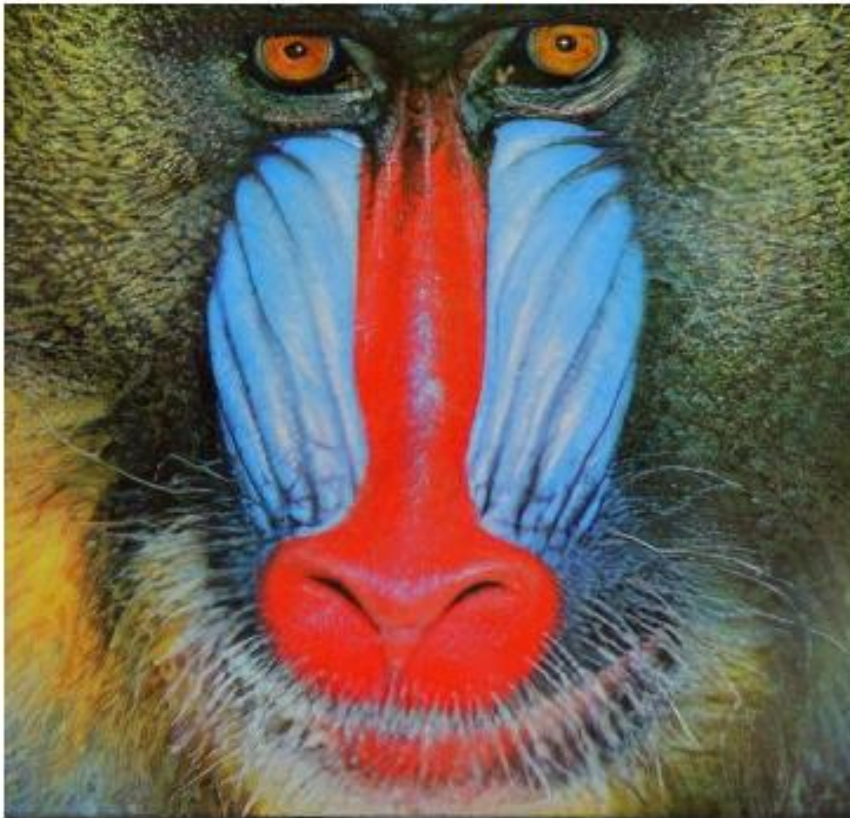


Image Super-resolution with GAN

4× SRGAN (proposed)



original



Image Super-resolution with GAN

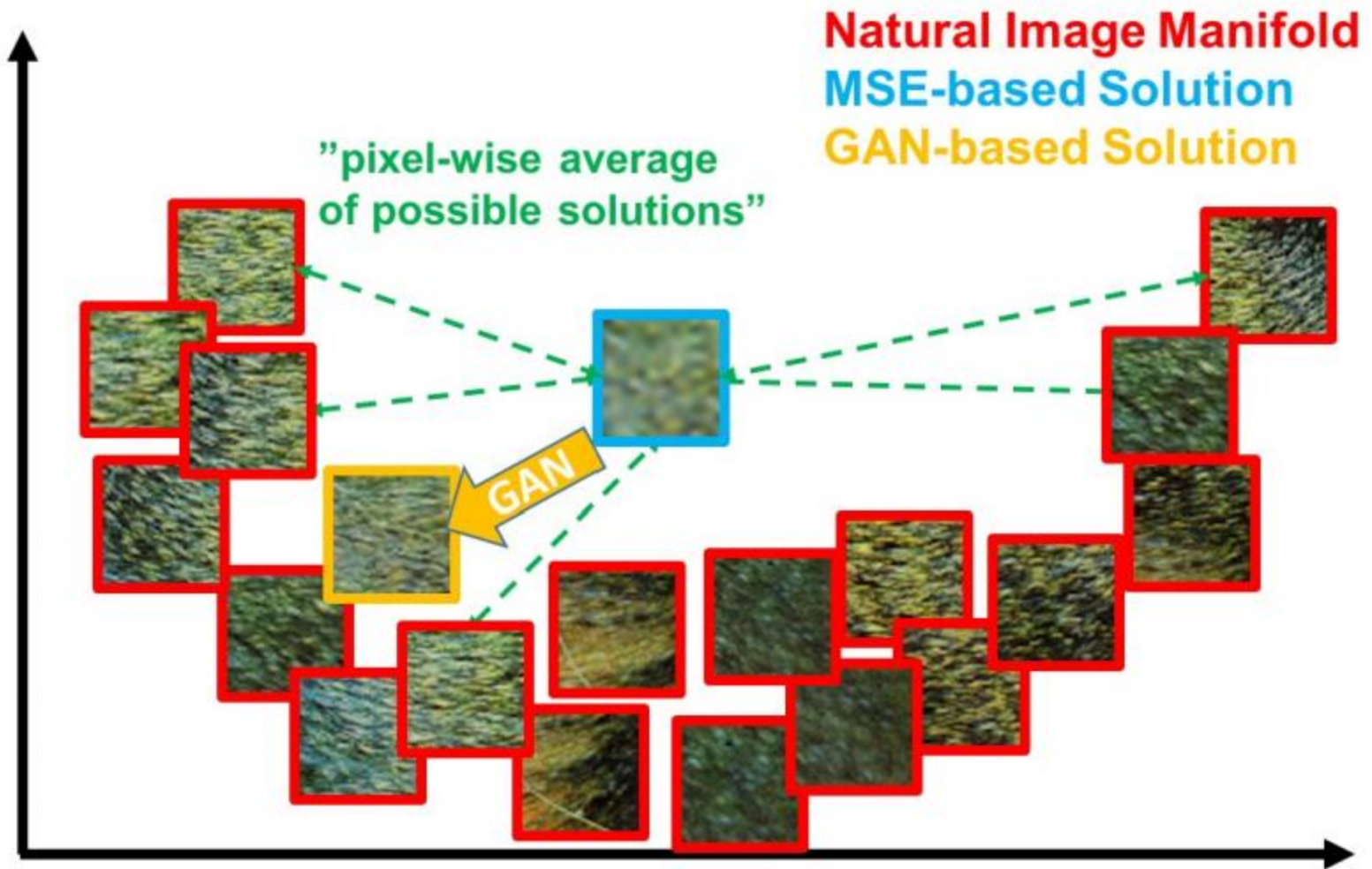
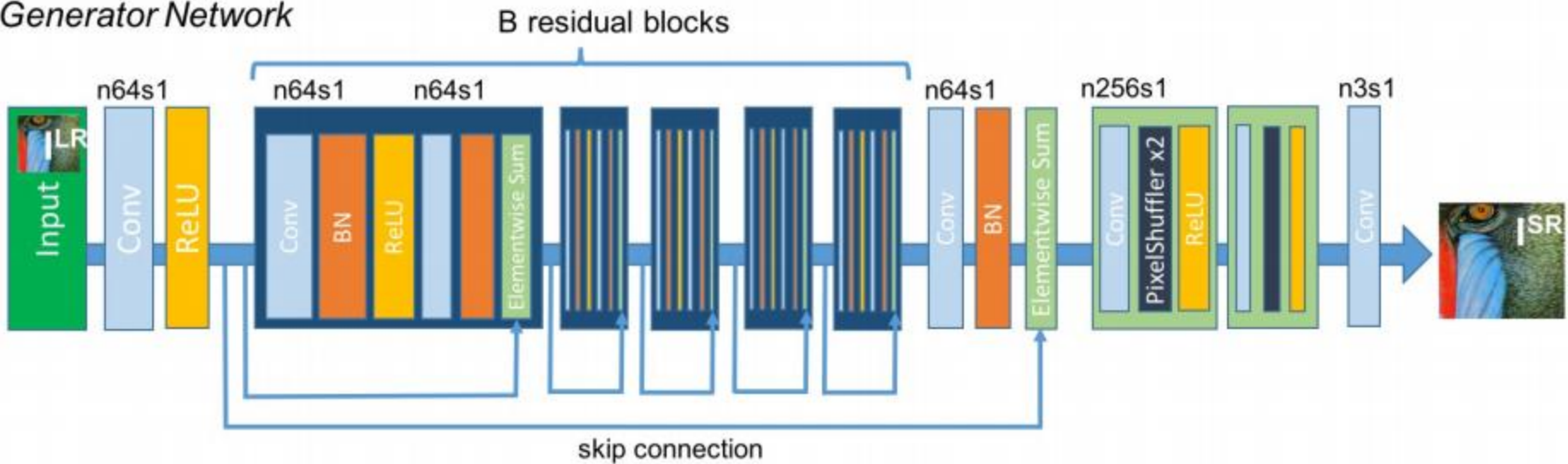


Image Super-resolution with GAN

Generator Network



Discriminator Network

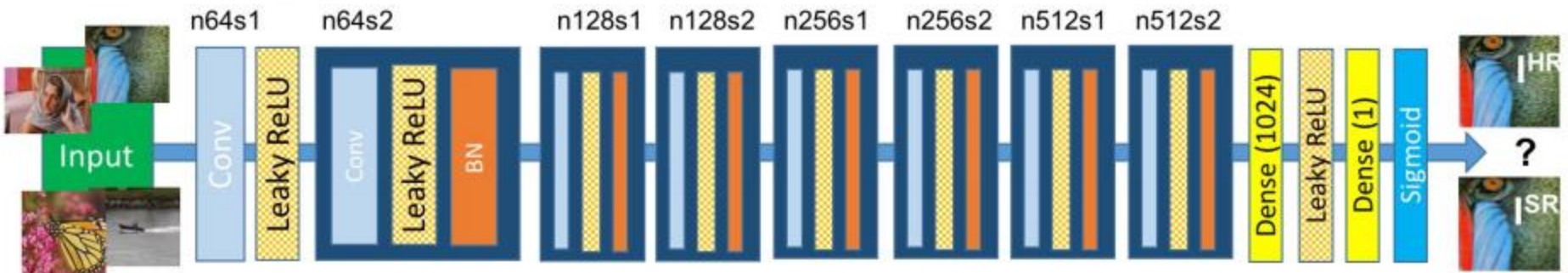


Image Super-resolution with GAN

$$l^{SR} = \underbrace{l_X^{SR}}_{\text{content loss}} + \underbrace{10^{-3} l_{Gen}^{SR}}_{\text{adversarial loss}}$$

bicubic



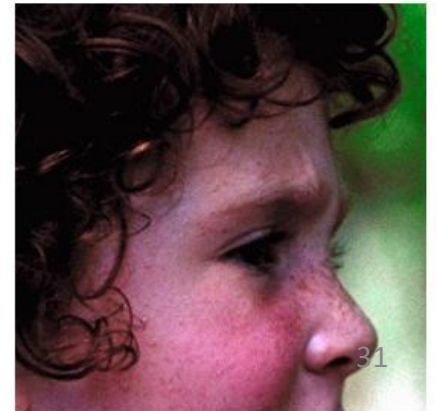
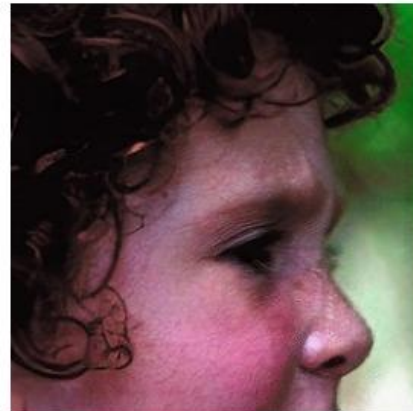
SRResNet



SRGAN

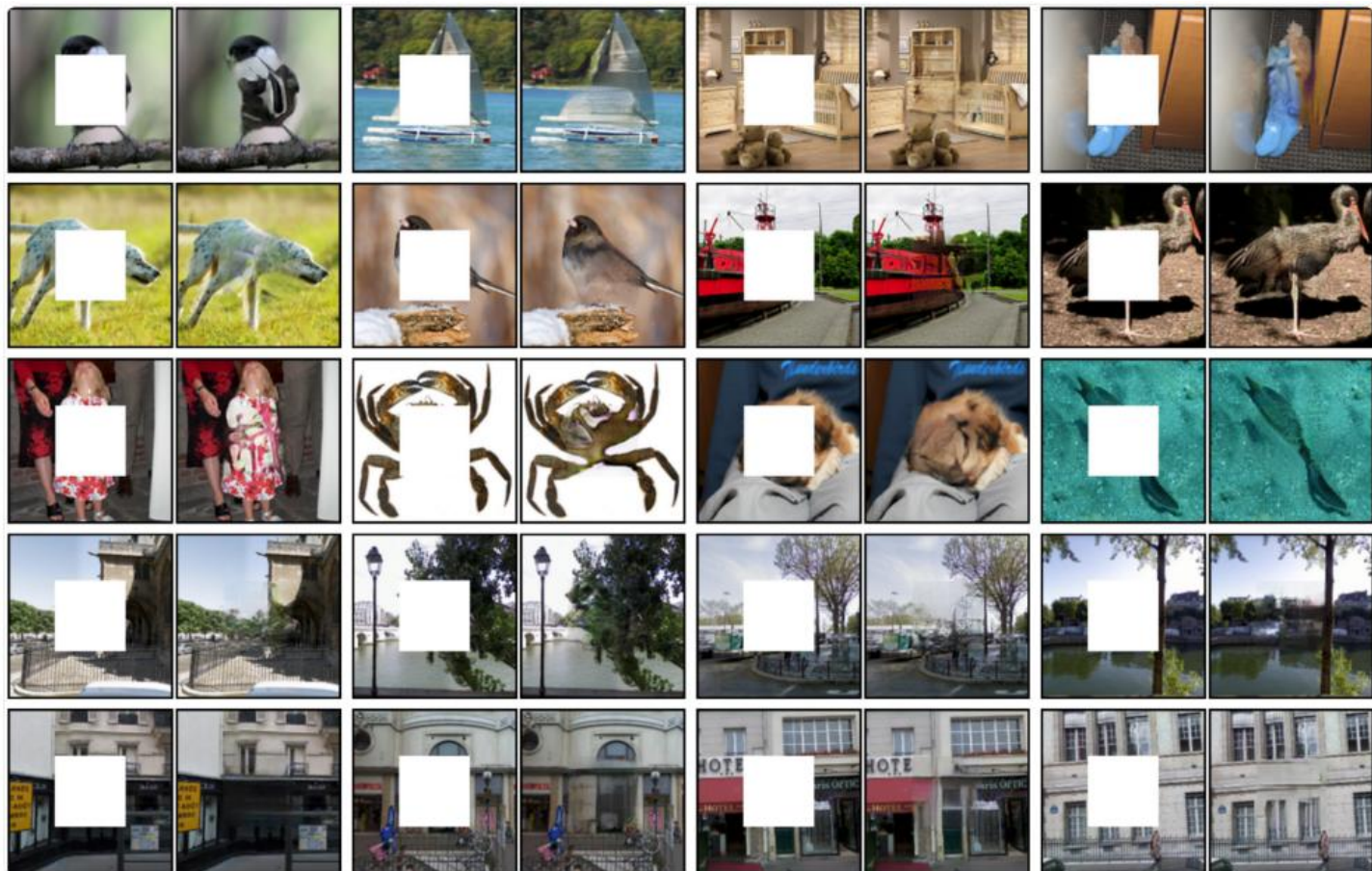


original



Context-Encoder for Image Inpainting

- For a pre-defined region, synthesize the image contents



Context-Encoder for Image Inpainting

- For a pre-defined region, synthesize the image contents



(a) Input context

(b) Human artist

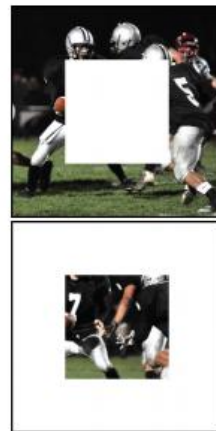
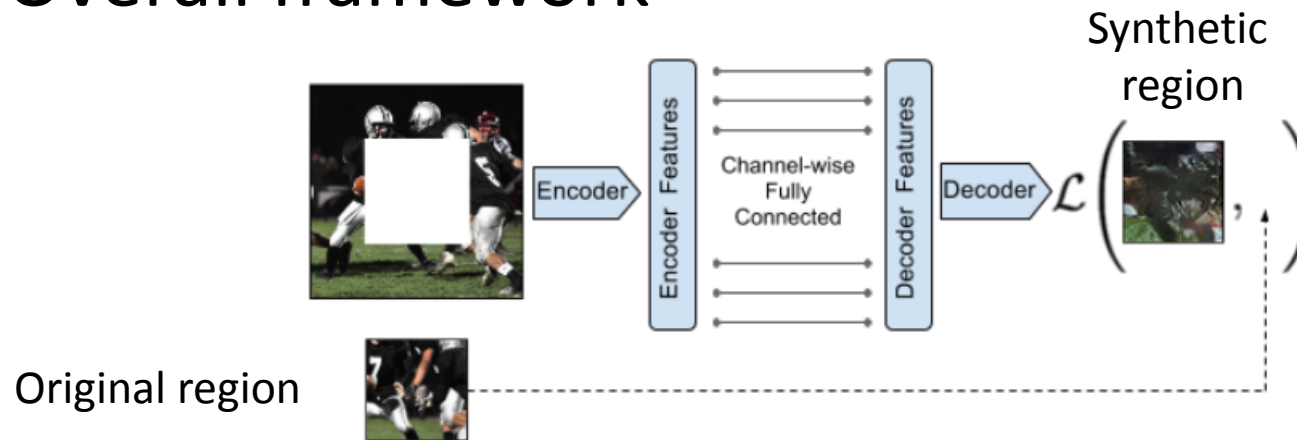


(c) Context Encoder
(L_2 loss)

(d) Context Encoder
($L_2 + \text{Adversarial loss}$)

Context-Encoder for Image Inpainting

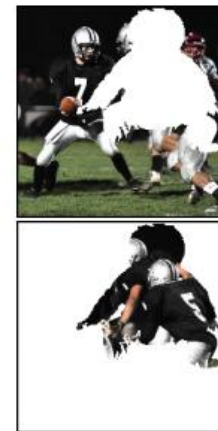
- Overall framework



(a) Central region



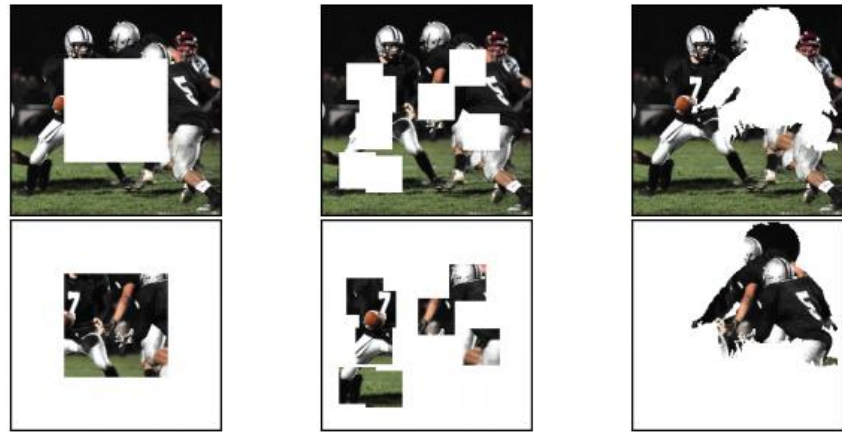
(b) Random block



(c) Random region

Context-Encoder for Image Inpainting

- The objective



(a) Central region

(b) Random block

(c) Random region

$$\mathcal{L}_{rec}(x) = \|\hat{M} \odot (x - F((1 - \hat{M}) \odot x))\|_2^2,$$

$$\mathcal{L}_{adv} = \max_D \mathbb{E}_{x \in \mathcal{X}} [\log(D(x)) + \log(1 - D(F((1 - \hat{M}) \odot x)))],$$

$$\mathcal{L} = \lambda_{rec} \mathcal{L}_{rec} + \lambda_{adv} \mathcal{L}_{adv}.$$

Context-Encoder for Image Inpainting

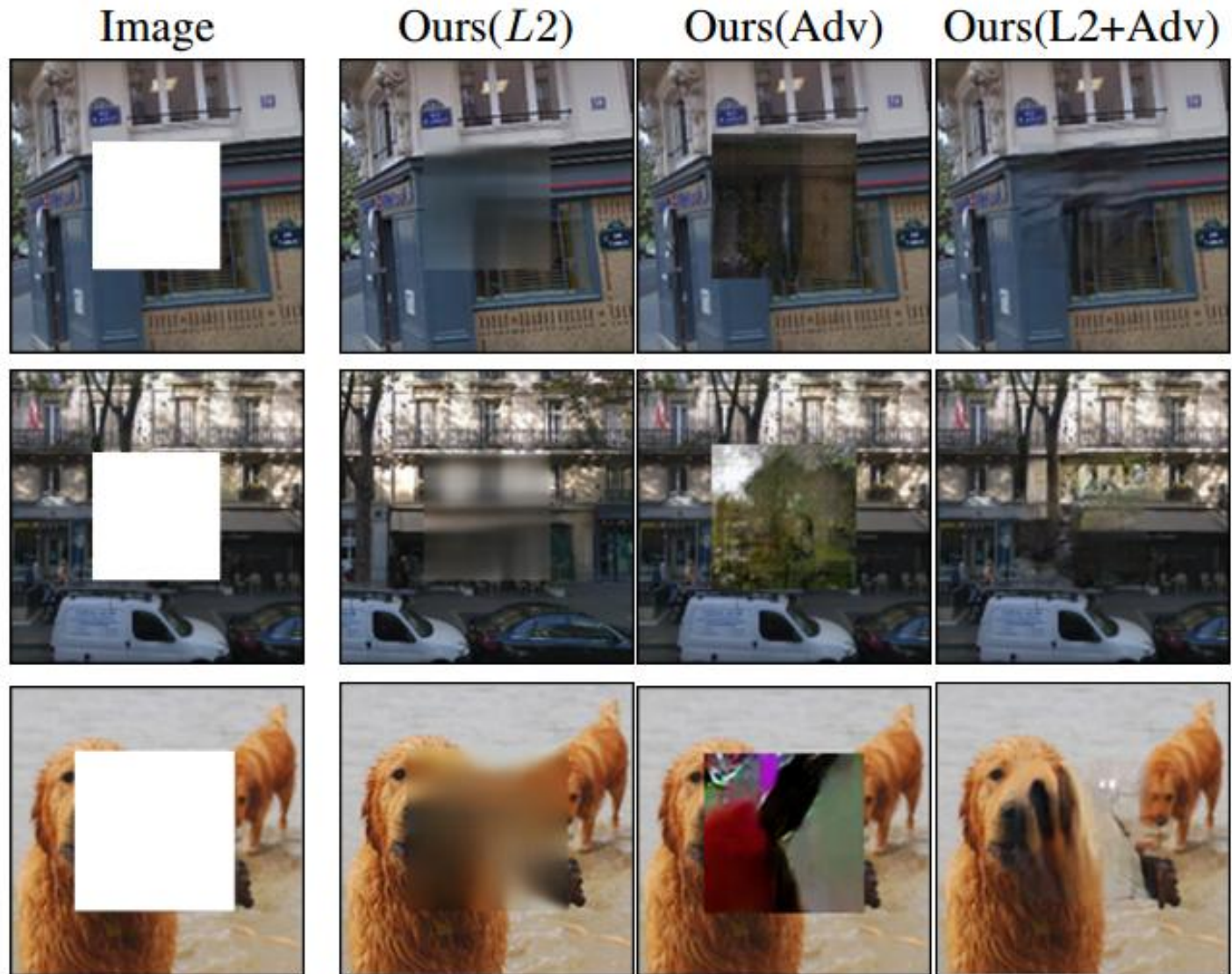
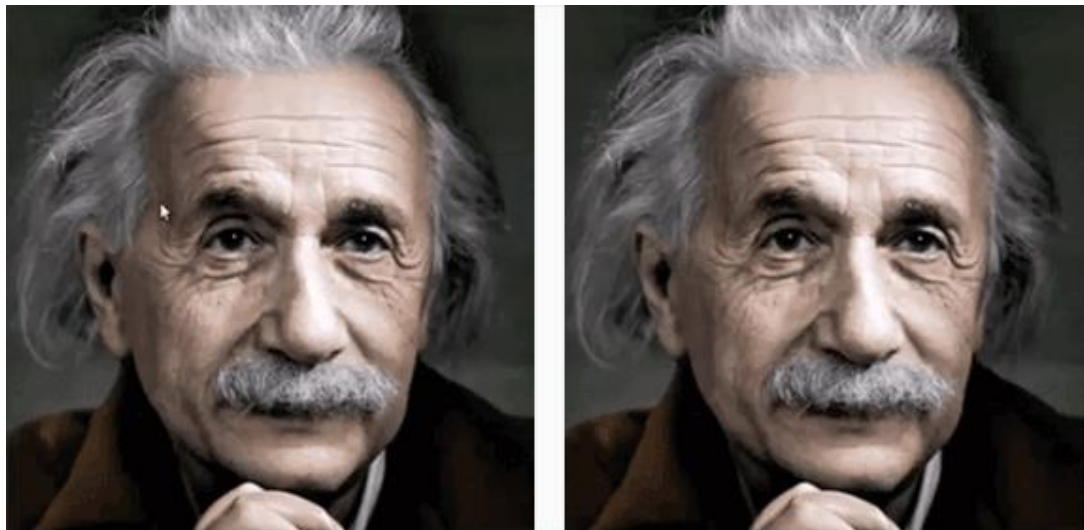


Image Inpainting with Partial Convolution

- Partial convolution for handling missing data
- L1 loss: minimizing the pixel differences between the generated image and their ground-truth images
- Perceptual loss: minimizing the VGG features of the generated images and their ground-truth images
- Style loss (Gram matrix): minimizing the gram matrices of the generated images and their ground-truth images



Image Inpainting with Partial Convolution: Results

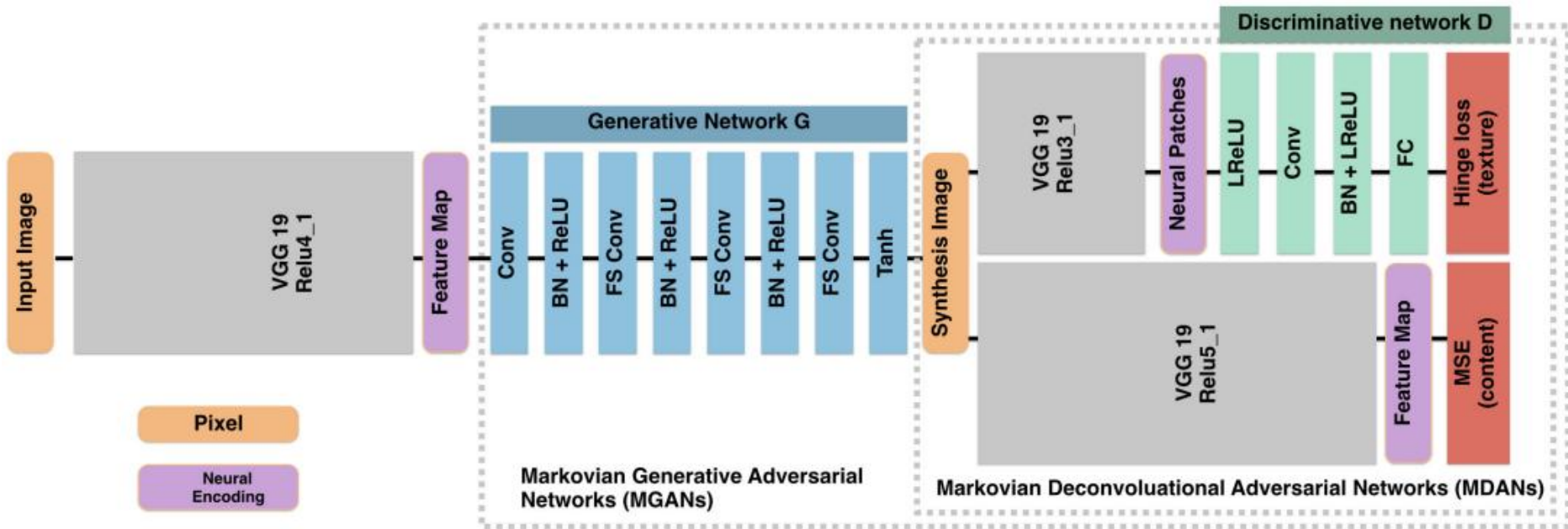


Texture Synthesis with Patch-based GAN

- Synthesize textures for input images



Texture Synthesis with Patch-based GAN



$$\mathbf{x} = \arg \min E_t(\Phi(\mathbf{x}), \Phi(\mathbf{x}_t)) + \alpha_1 E_c(\Phi(\mathbf{x}), \Phi(\mathbf{x}_c)) + \alpha_2 \mathcal{Y}(\mathbf{x})$$

MSE Loss

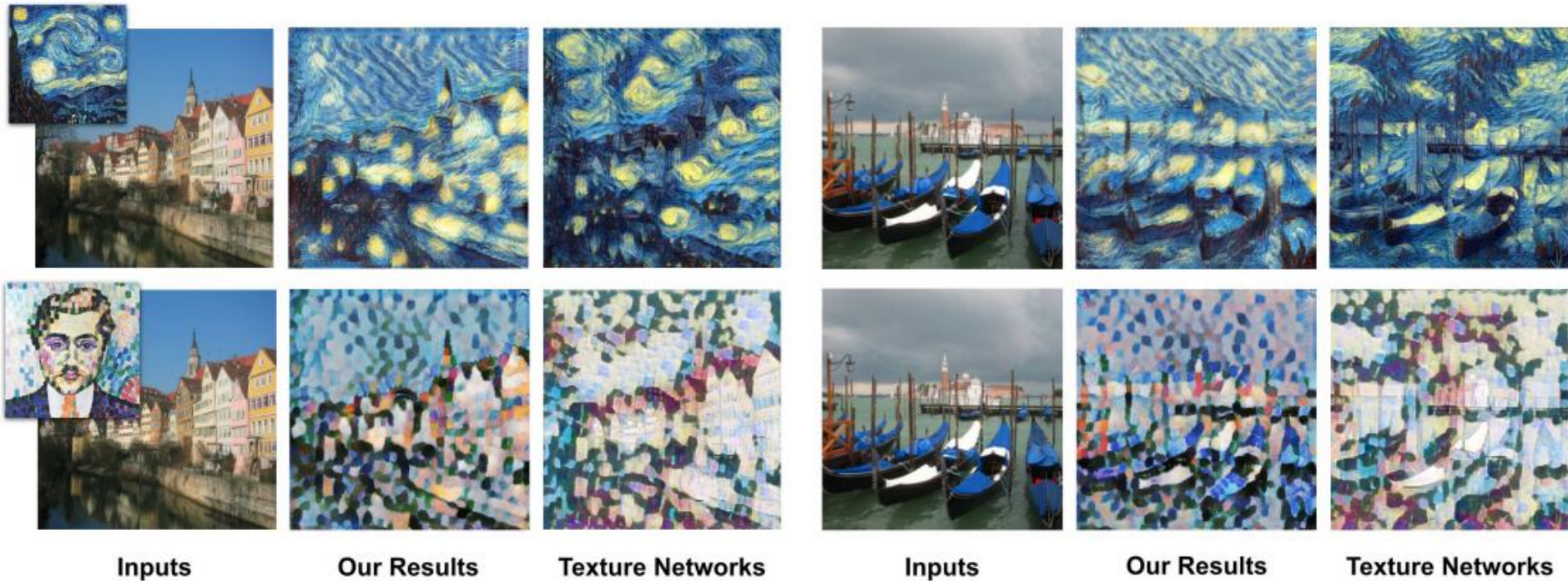
Adv loss

Texture Synthesis with Patch-based GAN

$$\mathbf{x} = \arg \min E_t(\Phi(\mathbf{x}), \Phi(\mathbf{x}_t)) + \alpha_1 E_c(\Phi(\mathbf{x}), \Phi(\mathbf{x}_c)) + \alpha_2 \Upsilon(\mathbf{x})$$

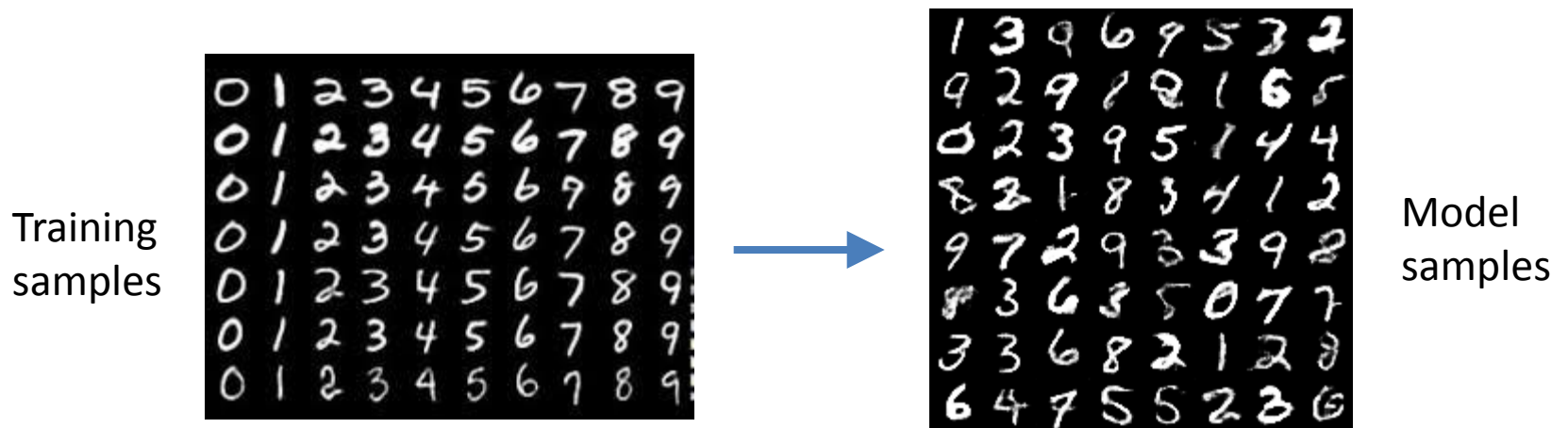


Texture Synthesis with Patch-based GAN



Conditional GAN

- GAN is too free. How to add some constraints?
- Add conditional variables \mathbf{y} into the generator



Conditional GAN

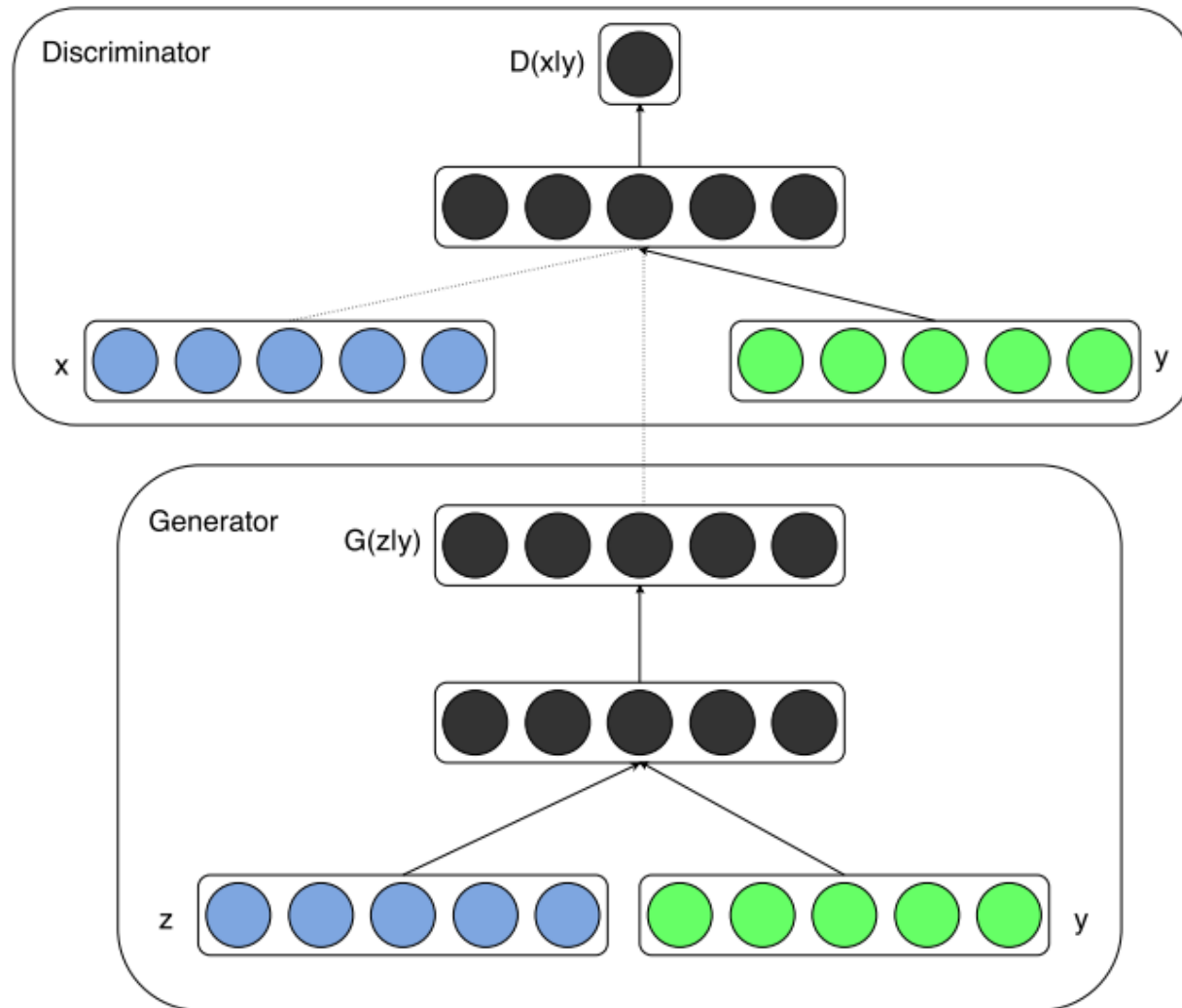
- GAN is too free. How to add some constraints?
- Add conditional variables \mathbf{y} into G and D

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

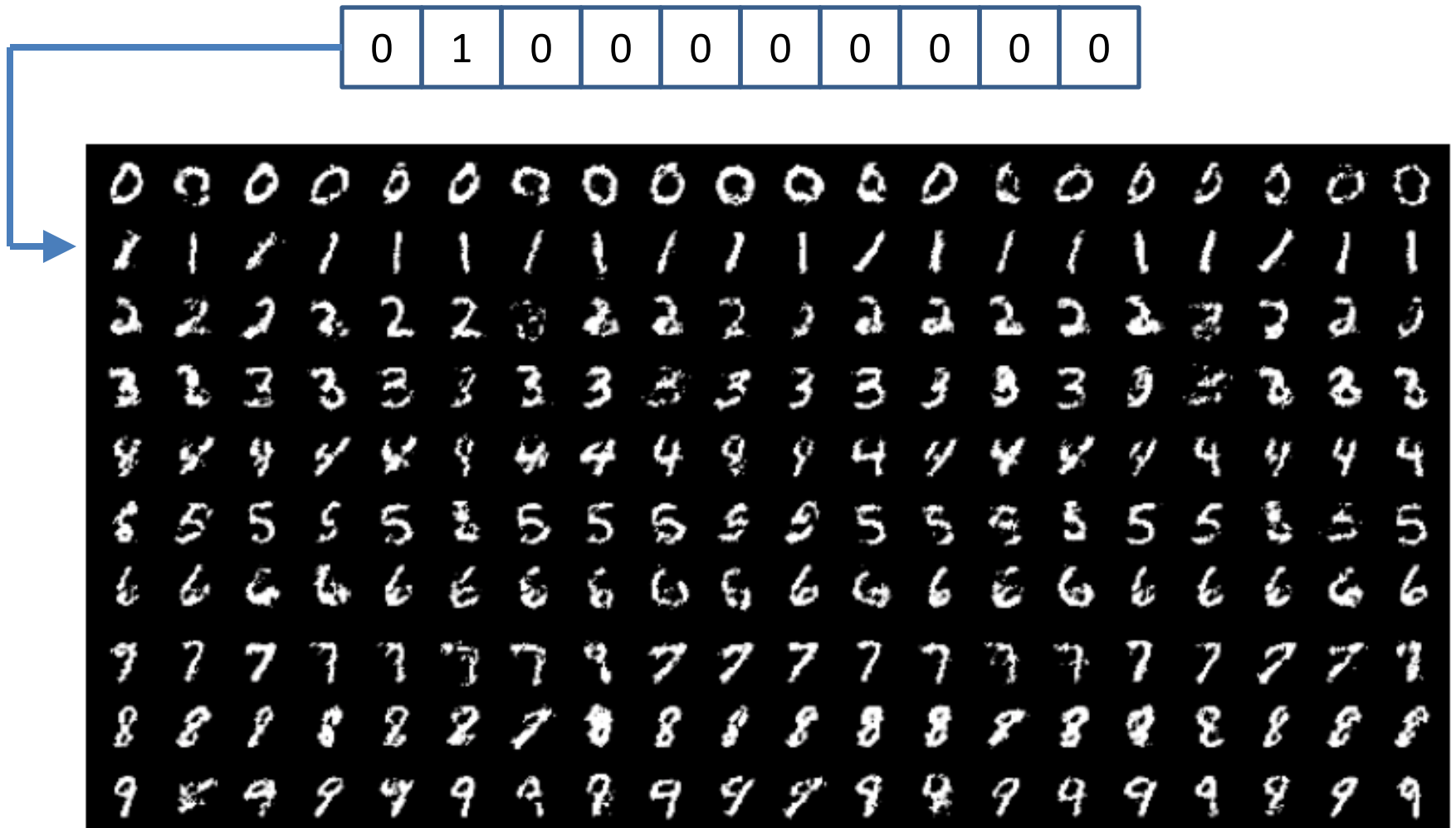


$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))].$$

Conditional GAN



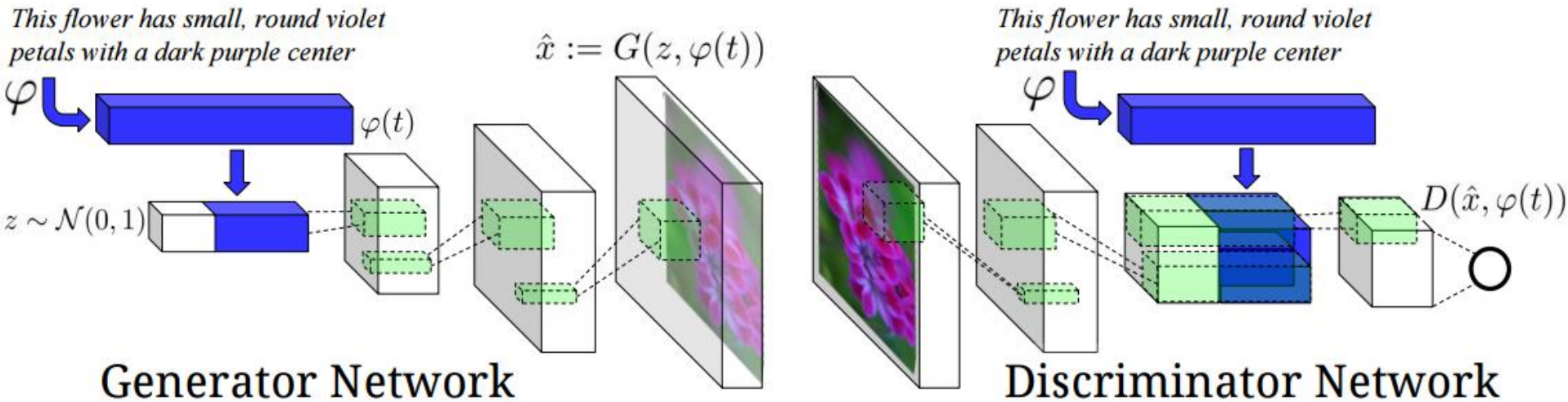
Conditional GAN



Conditional GAN

- Positive samples for D
 - True data + corresponding conditioning variable
- Negative samples for D
 - Synthetic data + corresponding conditioning variable
 - True data + non-corresponding conditioning variable

Text-to-Image Synthesis



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

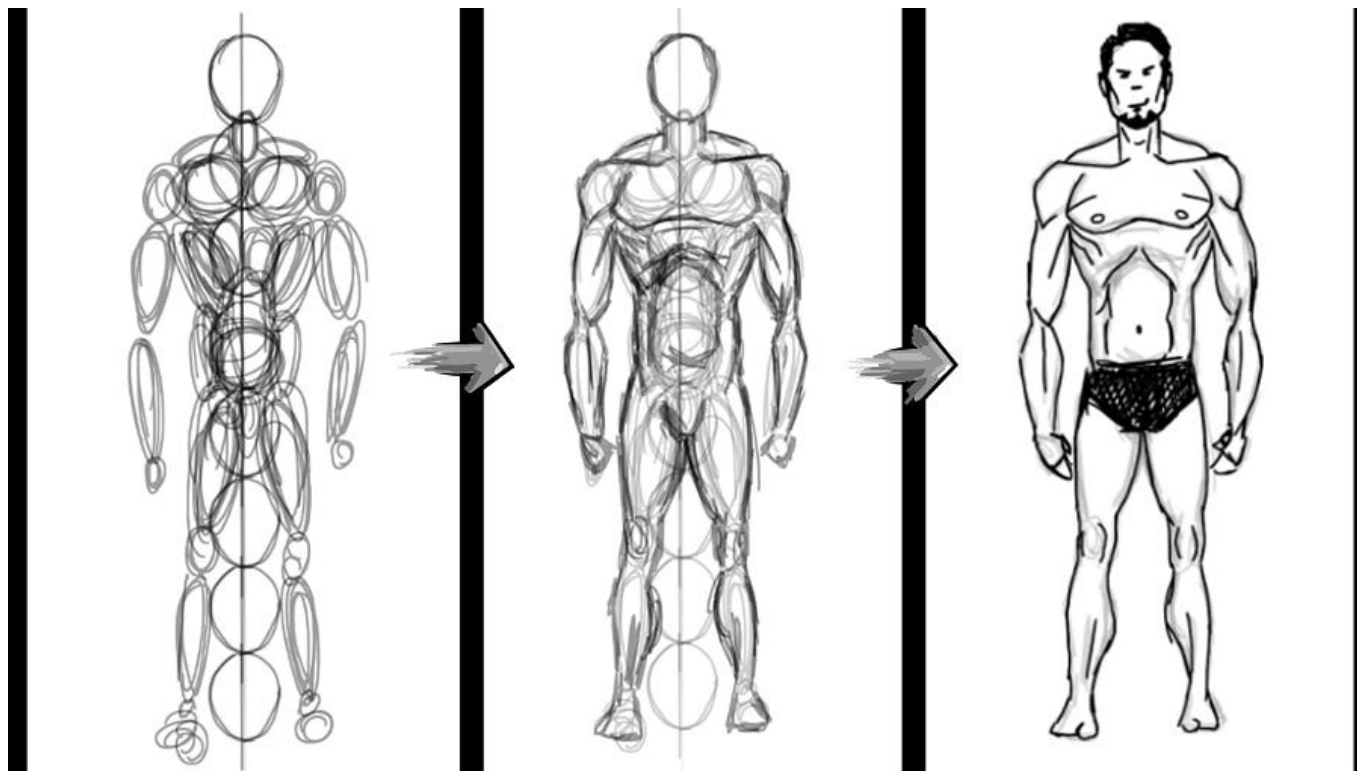


this magnificent fellow is almost all black with a red crest, and white cheek patch.



StackGAN: Text to Photo-realistic Images

- How humans draw a figure?
 - A coarse-to-fine manner



Zhang et al. 2016

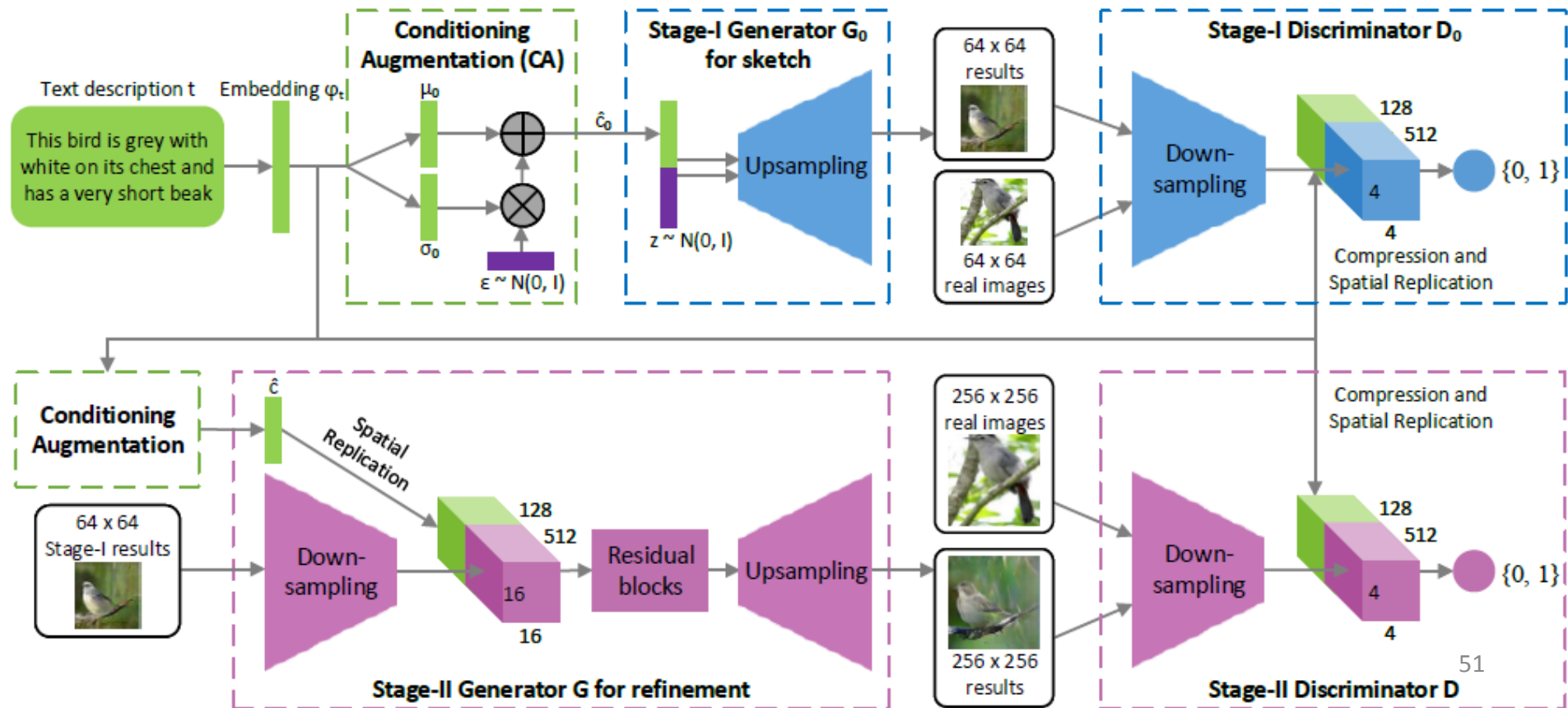
StackGAN: Text to Photo-realistic Images

- Use stacked GAN structure for text-to-image synthesis



StackGAN: Text to Photo-realistic Images

- Use stacked GAN structure for text-to-image synthesis

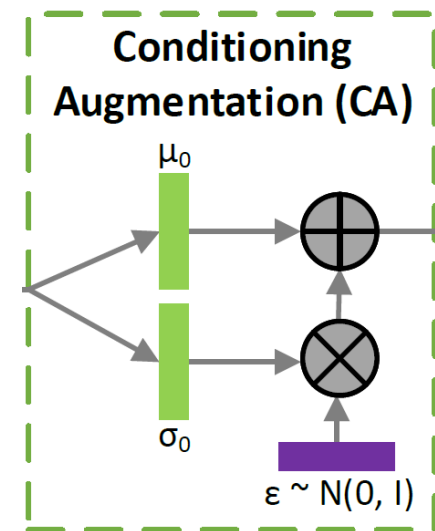


StackGAN: Text to Photo-realistic Images

- Conditioning augmentation
- No random noise vector \mathbf{z} for Stage-2
- Conditioning both stages on text help achieve better results
- Spatial replication for the text conditional variable
- Negative samples for D
 - True images + non-corresponding texts
 - Synthetic images + corresponding texts

Conditioning Augmentation

- How train parameters like the mean and variance of a Gaussian distribution
- $N(\mu_0, \sigma_0)$
- Sample from standard Normal distribution $N(0,1)$
- Multiple with σ_0 and then add with μ_0
- The re-parameterization trick



More StackGAN Results on Flower

Text
description

This flower is pink, white, and yellow in color, and has petals that are striped

This flower has a lot of small purple petals in a dome-like configuration

This flower is white and yellow in color, with petals that are wavy and smooth

This flower has petals that are dark pink with white edges and pink stamen

64x64
GAN-INT-
CLS



256x256
StackGAN



More StackGAN Results on COCO

A picture of a very clean living room



Eggs fruit candy nuts and meat served on white dish



A street sign on a stoplight pole in the middle of a day

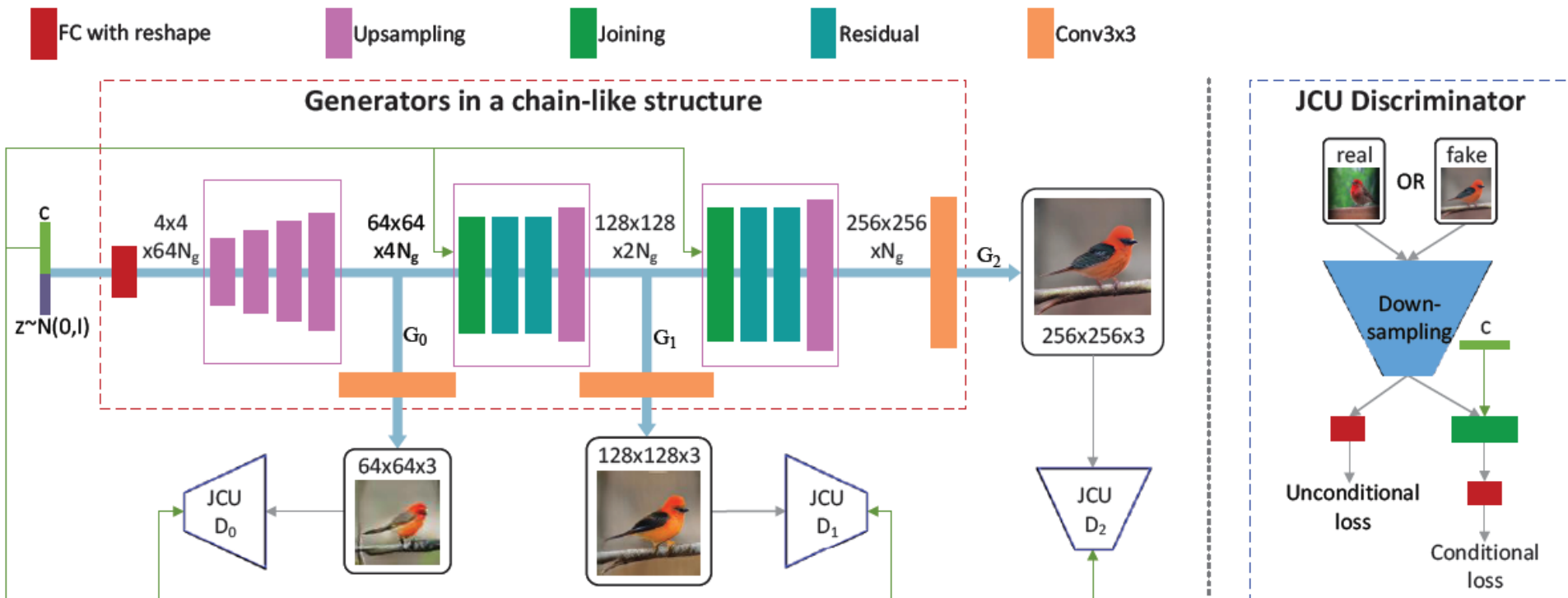


A group of people on skis stand in the snow



StackGAN-v2: Architecture

- Approximate multi-scale image distributions jointly
- Approximate conditional and unconditional image distributions jointly



StackGAN-v2: Results



256×256 samples by our StackGAN-v2 on LSUN bedroom dataset



256×256 samples by our StackGAN-v2 on LSUN church dataset



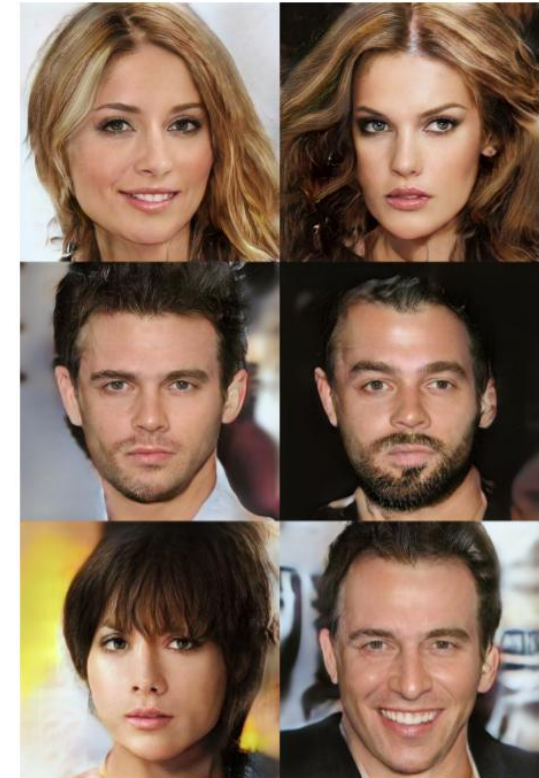
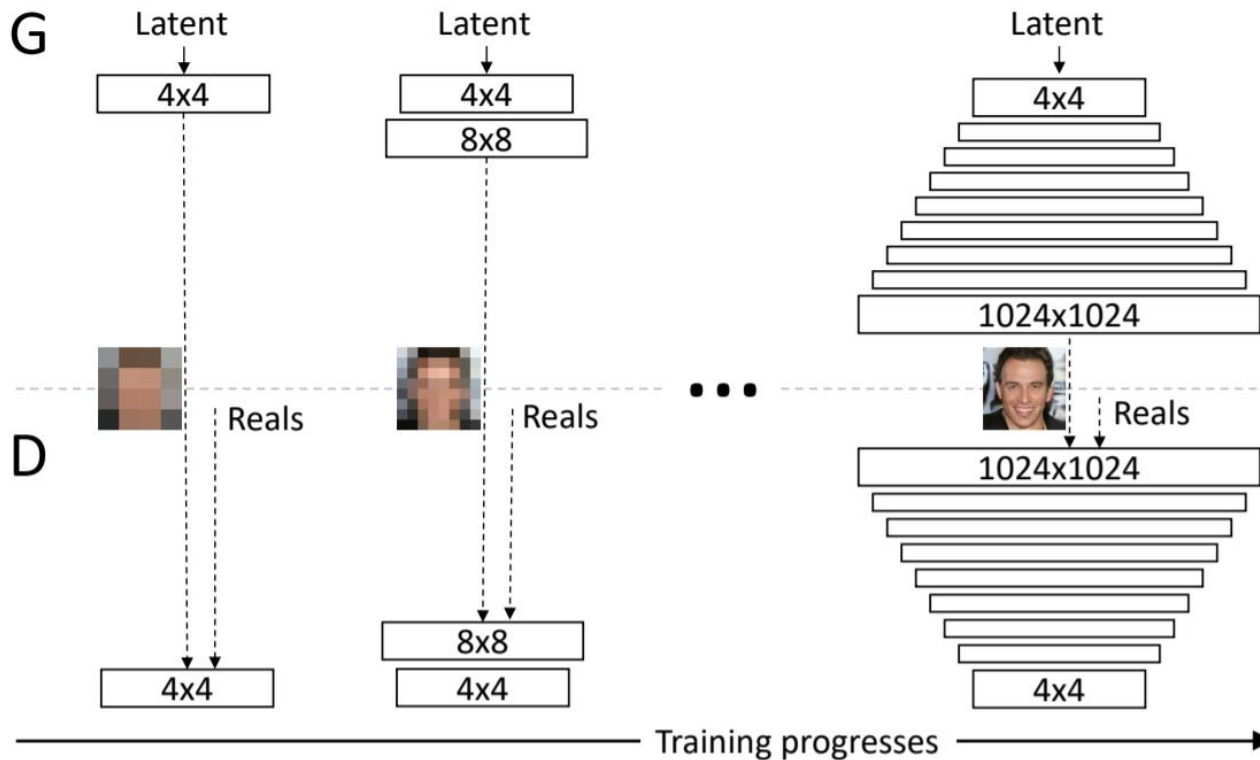
256×256 samples by StackGAN-v2 on ImageNet dog dataset



256×256 samples by StackGAN-v2 on ImageNet cat dataset

Progressive Growing of GAN

- Share the similar spirit with StackGAN-v1/-v2 but use a different training strategy



Progressive Growing of GAN

- Impressively realistic face images

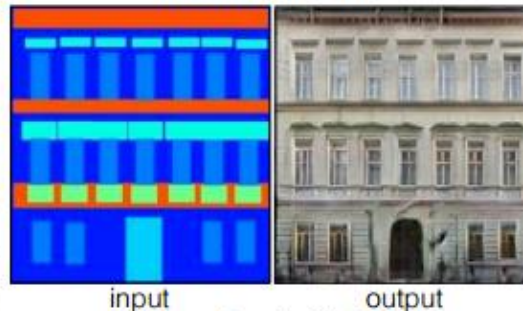


Image-to-Image Translation with Conditional GAN

Labels to Street Scene



Labels to Facade



BW to Color



Aerial to Map



Day to Night



Edges to Photo



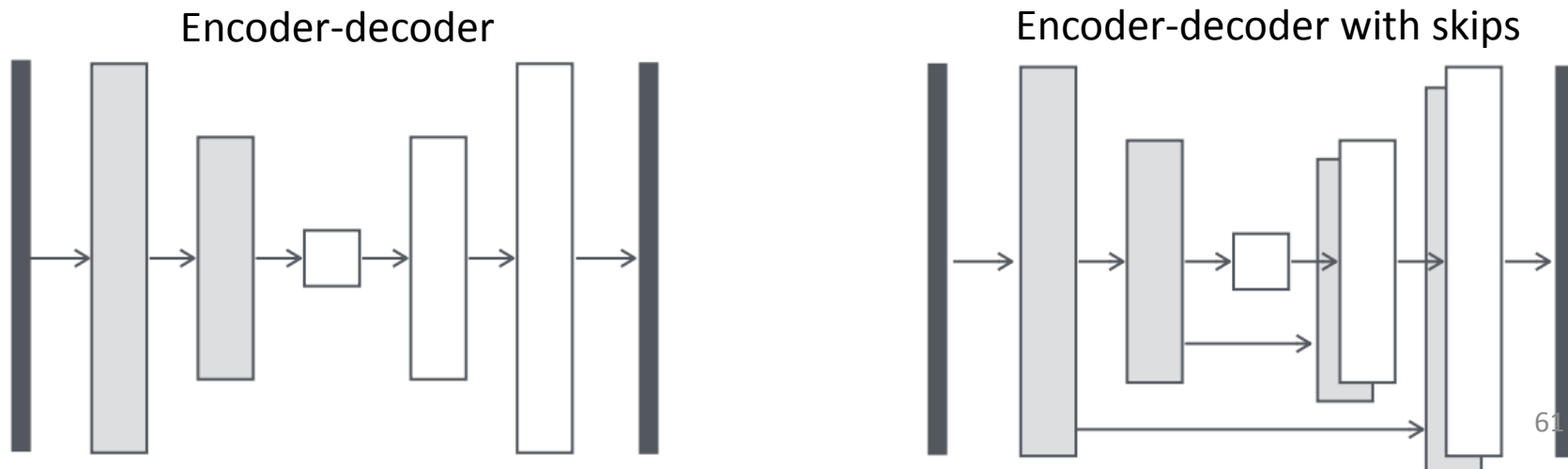
Image-to-Image Translation with Conditional GAN

- Incorporate L1 loss into the objective function

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y \sim p_{data}(x,y), z \sim p_z(z)} [\|y - G(x, z)\|_1].$$

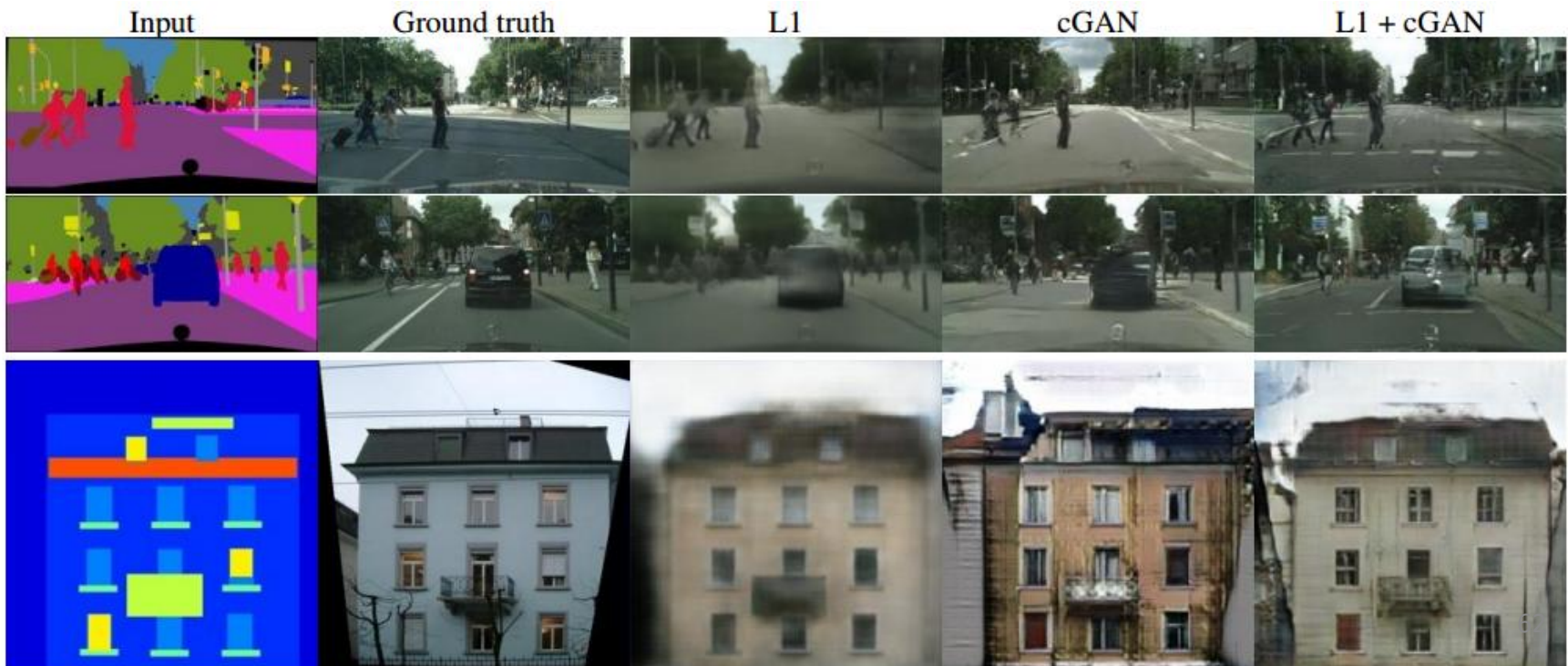
$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G).$$

- Adopt the U-net structure for the generator



Patch-based Discriminator

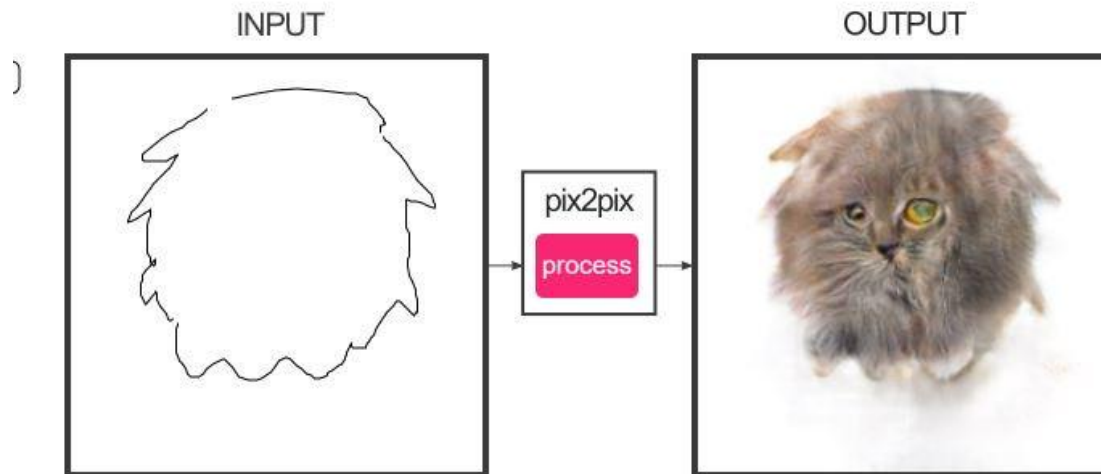
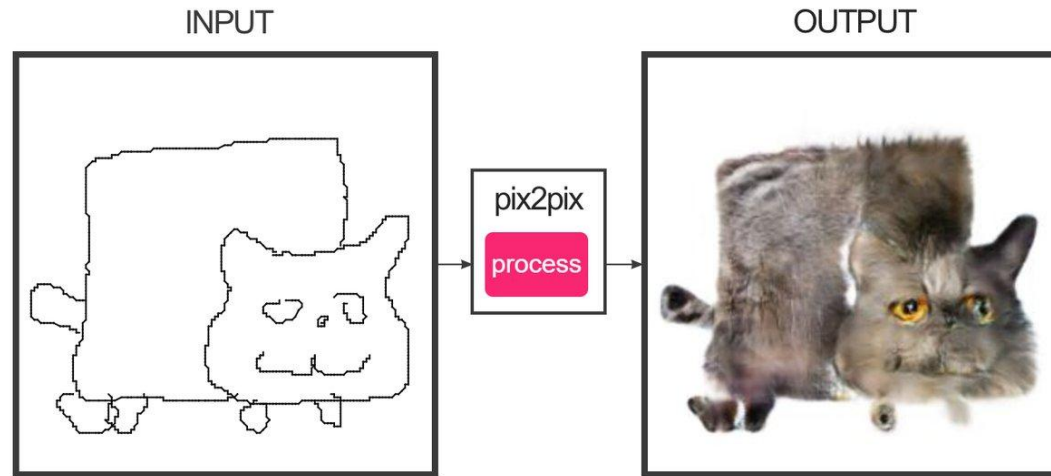
- Separate each image into $N \times N$ patches
- Instead of distinguish whether the whole image is real or fake, train a patch-based discriminator



More Results

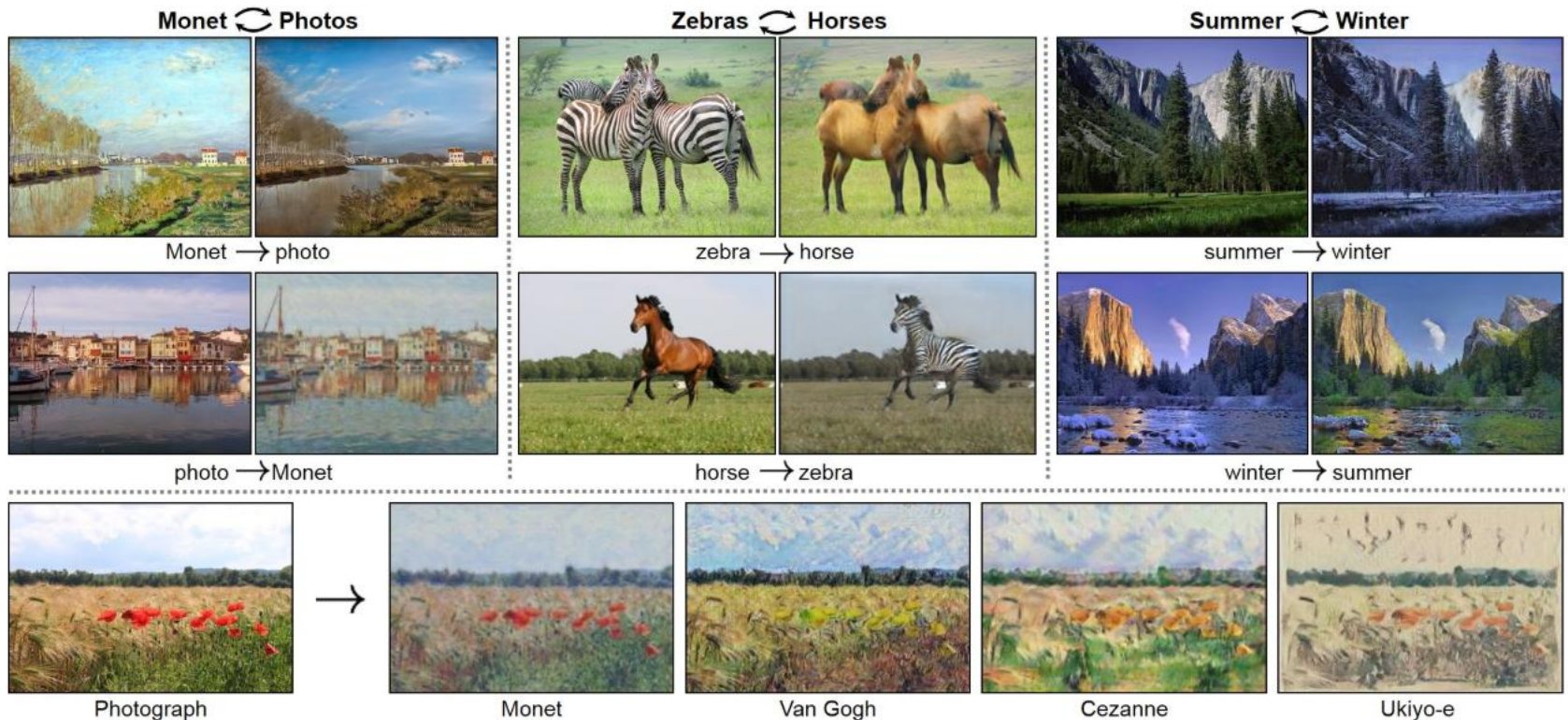


More Results



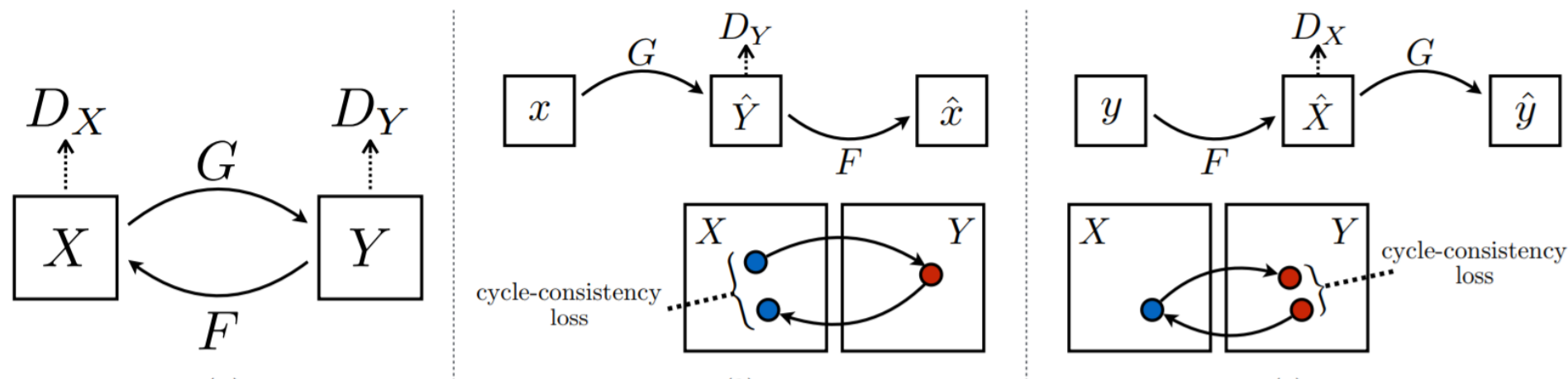
CycleGAN

- All previous methods require to have paired training data, i.e., exact input-output pairs, which can be extremely difficult to obtain in practice



CycleGAN

- The framework learns two mapping functions (generators) $G : X \rightarrow Y$ and $F : Y \rightarrow X$ with two domain discriminators D_X and D_Y



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

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

CycleGAN: Results



apple → orange



orange → apple

Input



Monet



Van Gogh



Cezanne



Ukiyo-e



CycleGAN: Results

Input



Output



Input

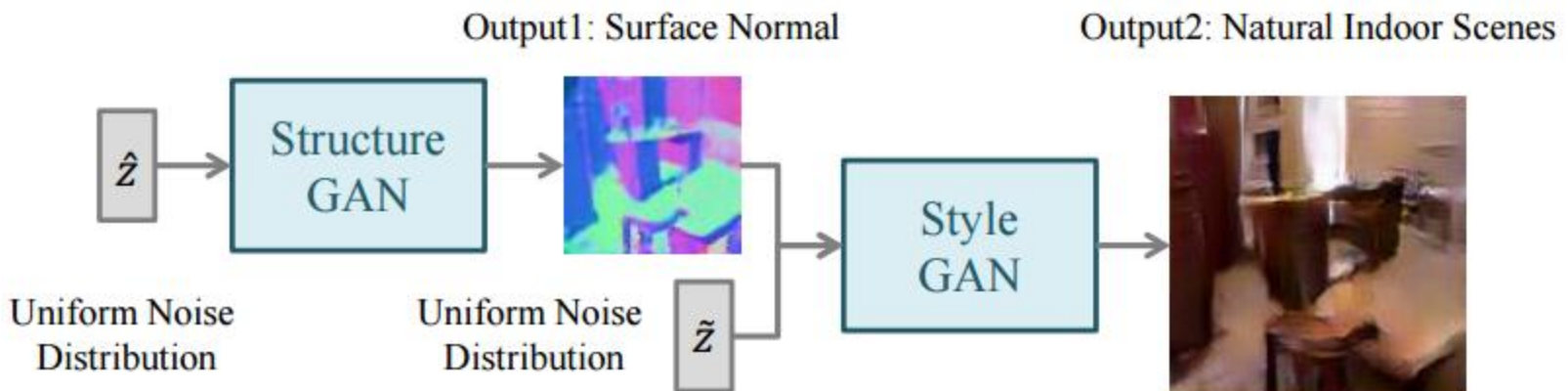


Output

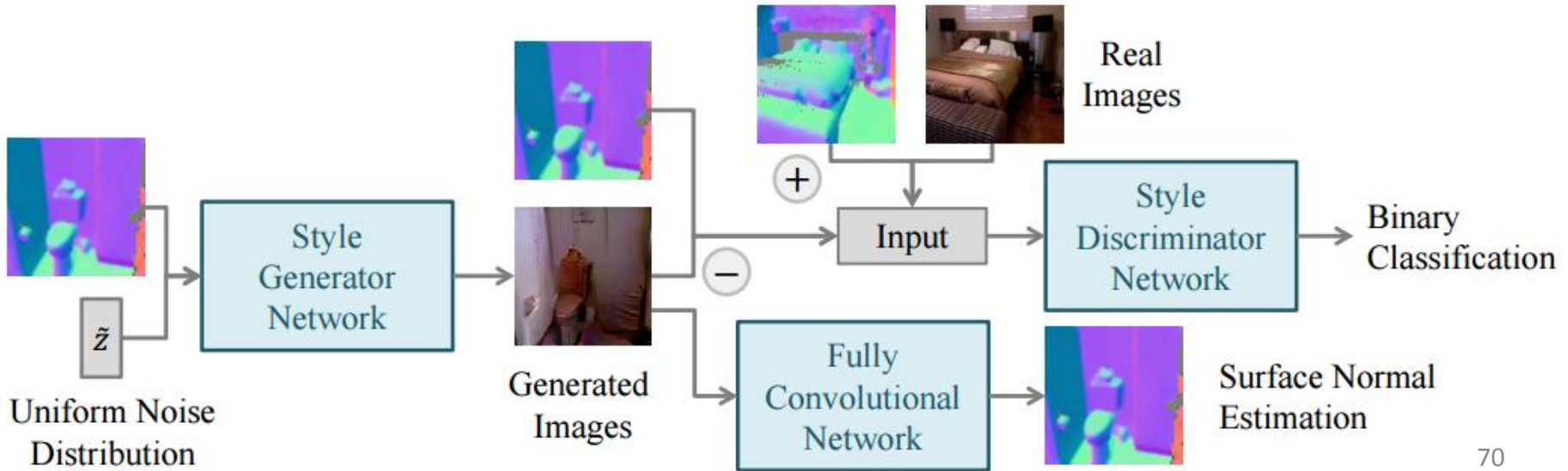
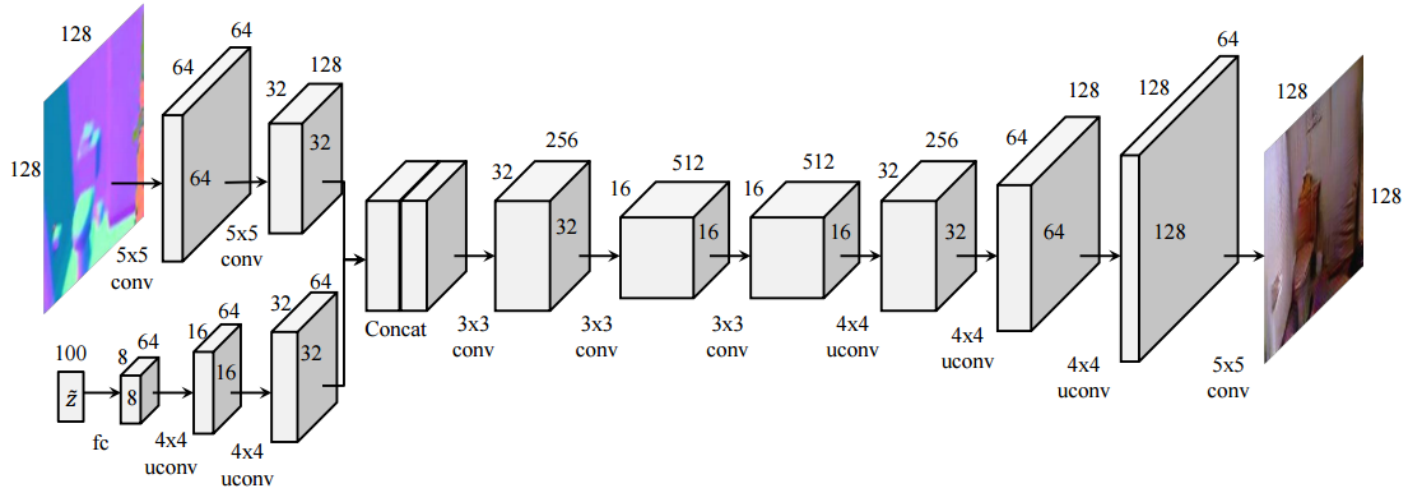


S²-GAN: Decomposing difficult problems into subproblems

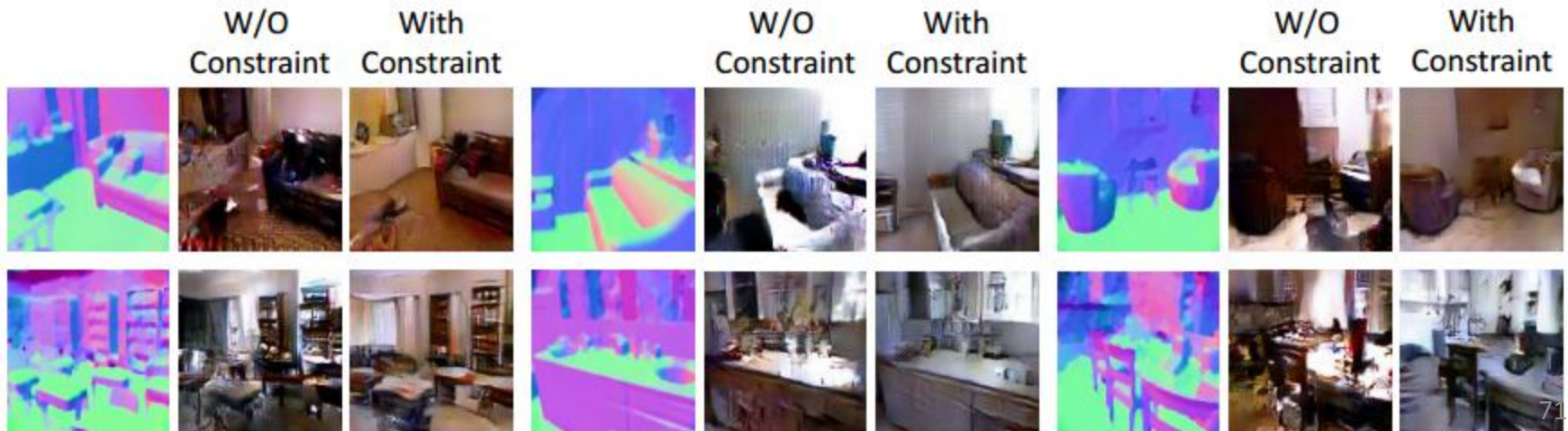
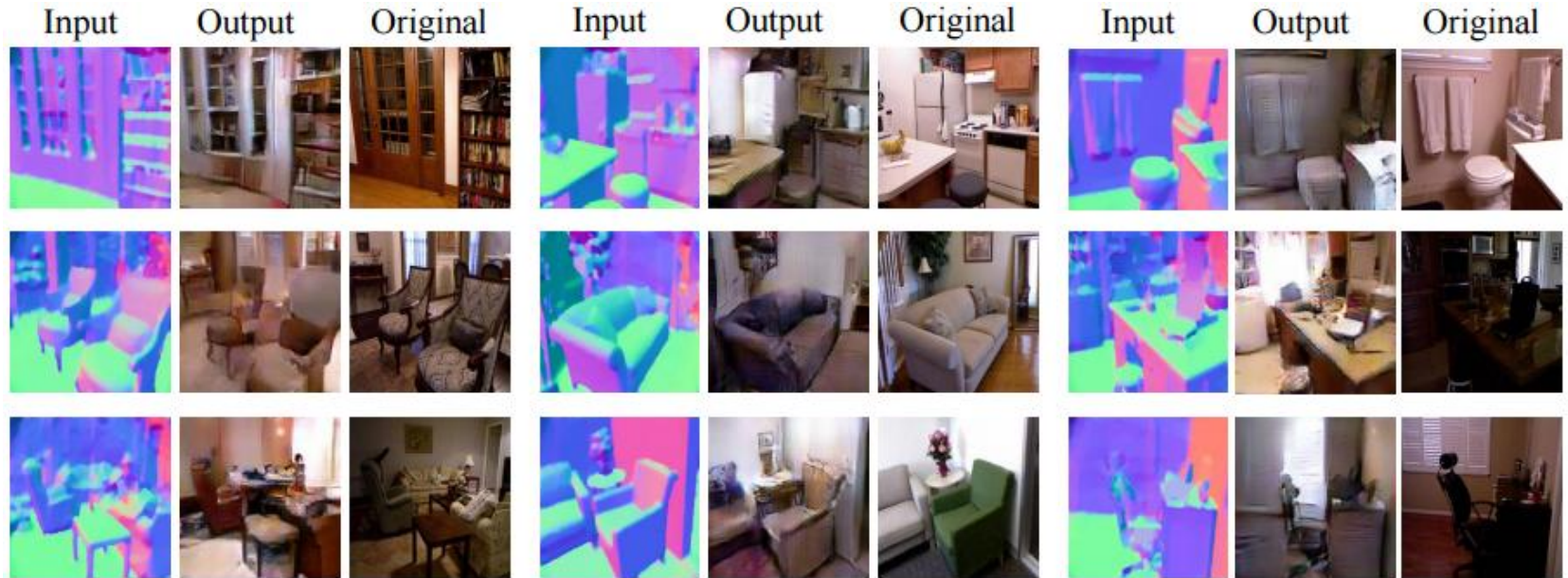
- Generating indoor images
- Generating surface normal map + surface style map



Style-GAN



S²-GAN Results



Insights

- Some insights
 - Decomposing the problems into easier problems
 - Spatially well-aligned conditioning variables are generally better

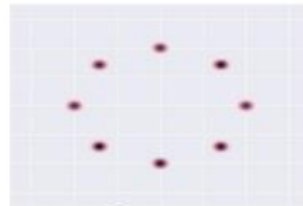
Non-convergence in GANs

- Finding equilibrium is a game of two players
- Exploiting convexity in function space, GAN training is theoretically guaranteed to converge if we can modify the density functions directly, but:
 - Instead, we modify G (sample generation function) and D (density ratio), not densities
 - We represent G and D as highly non-convex parametric functions
- “Oscillation”: can train for a very long time, generating very many different categories of samples, without clearly generating better samples
- Mode collapse: most severe form of non-convergence

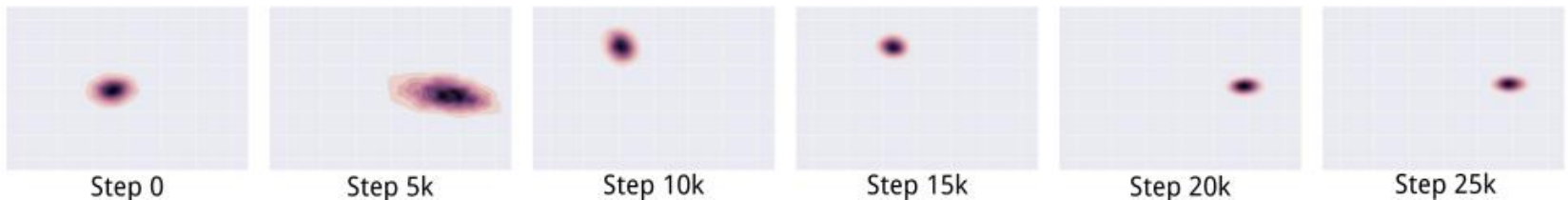
Mode Collapse

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

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



Target



Mode Collapse Causes Low Output Diversity

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



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



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



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



(Reed et al 2016)



(Reed et al, submitted to ICLR 2017)

Conditioning Augmentation

This small blue bird has a short pointy beak and brown on its wings



This bird is completely red with black wings and pointy beak



A small sized bird that has a cream belly and a short pointed bill



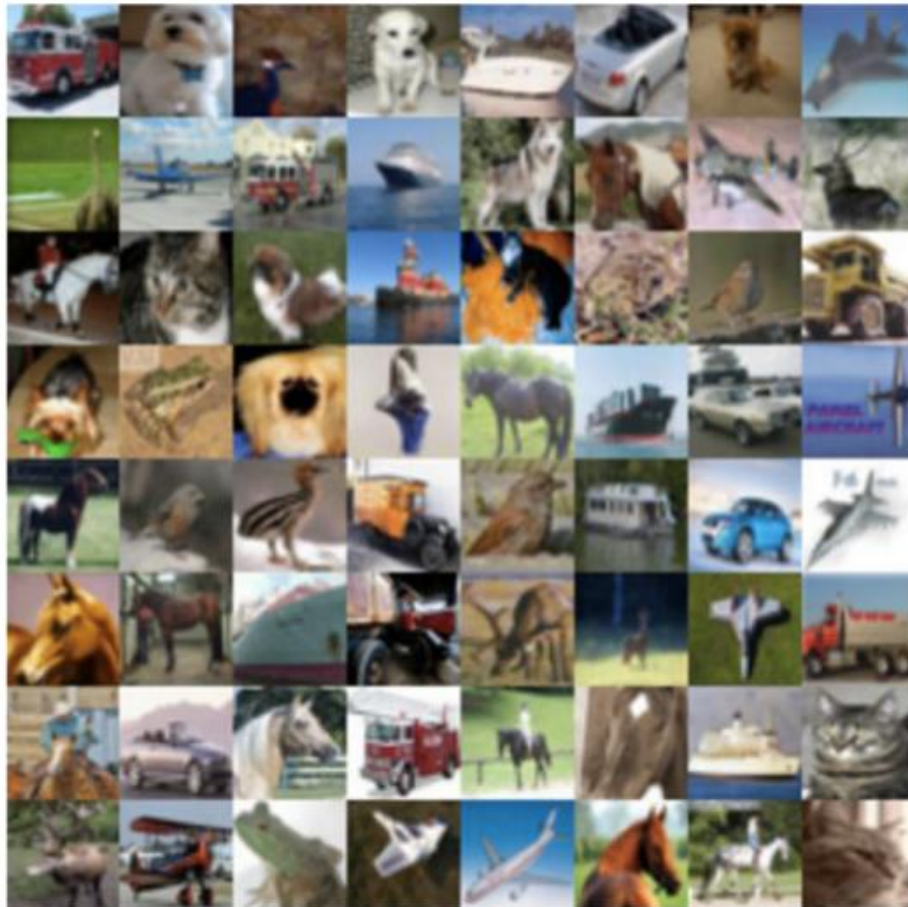
A small bird with a black head and wings and features grey wings



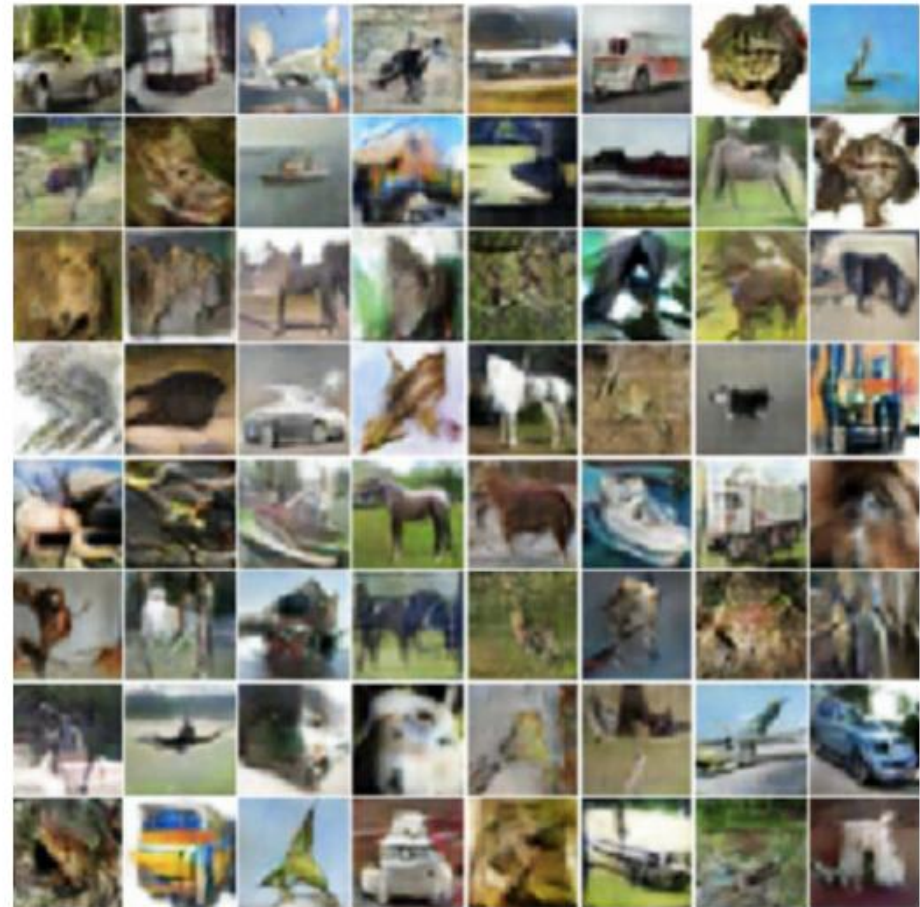
Minibatch Features

- Add minibatch features that classify each example by comparing it to other members of the minibatch (Salimans et al 2016)
- Nearest-neighbor style features detect if a minibatch contains samples that are too similar to each other

Minibatch GAN on CIFAR

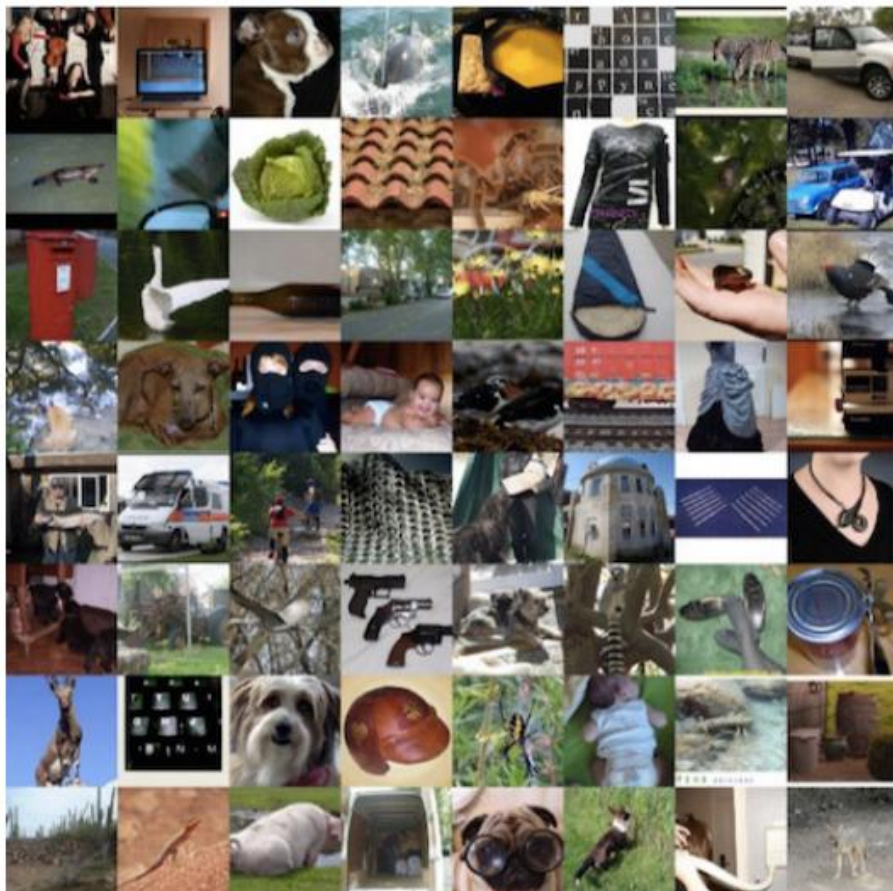


Training Data



Samples

Minibatch GAN on ImageNet



Cherry-picked Results

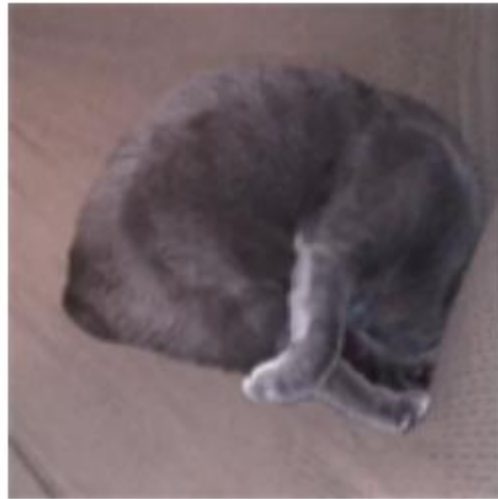


Problems with Counting



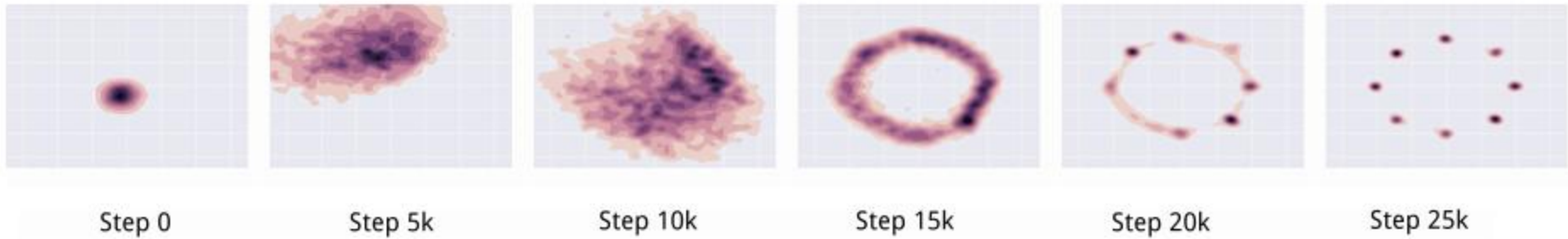
Problems with Perspective





Unrolled GANs

- A toy example

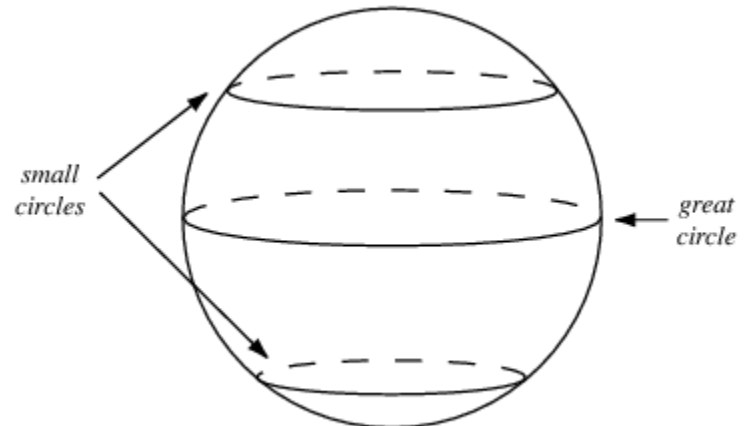
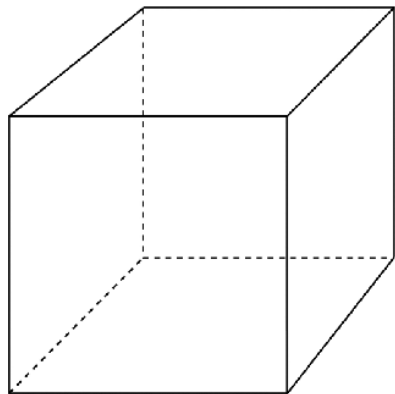


Tips and Tricks for training GAN

- **Normalize the inputs**
 - normalize the images between -1 and 1
 - Tanh as the last layer of the generator output
- **Modified loss function**
 - Because of the vanishing gradients (Goodfellow et al 2014). Use $J^{(G)} = -\frac{1}{2}\mathbb{E}_z \log D(G(z))$
 - Flip labels when training generator: real = fake, fake = real

Tips and Tricks for training GAN

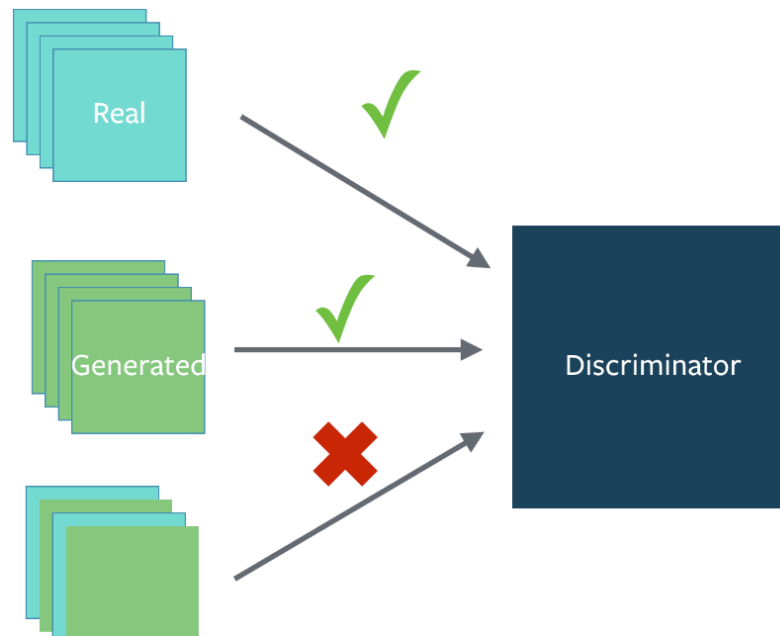
- **Use a spherical Z**
 - Don't sample from a Uniform distribution
 - Sample from a Gaussian distribution
 - When doing interpolations, do the interpolation via a great circle, rather than a straight line (White et al 2016)



Tips and Tricks for training GAN

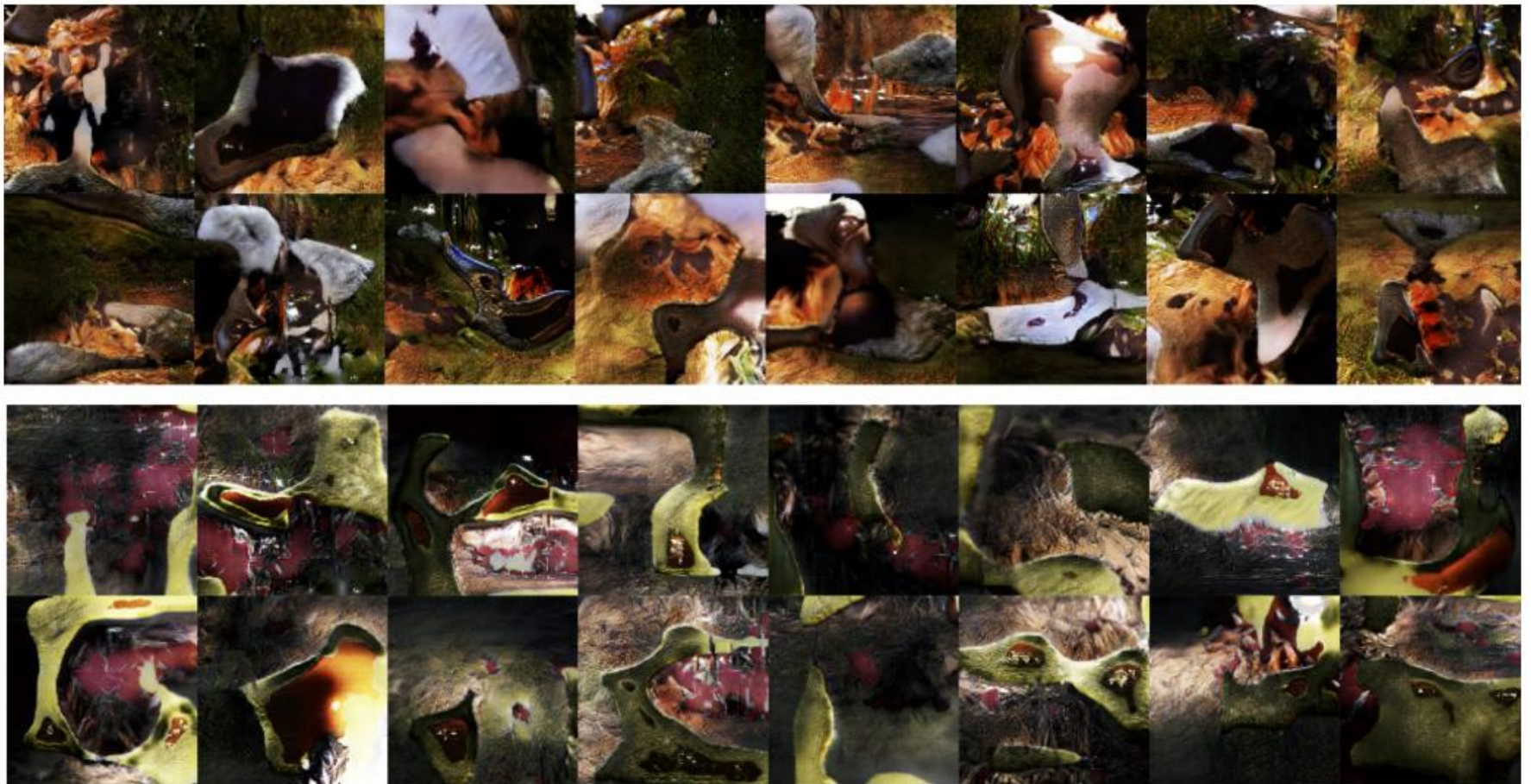
- **Batch normalization**

- Compute mean and standard deviation of features
- Normalize features (subtract mean, divide by standard deviation)



Tips and Tricks for training GAN

- **Batch normalization in G**



Tips and Tricks for training GAN

- **Reference Batch Normalization**

- Fix a reference batch $R = \{r^{(1)}, r^{(2)}, \dots, r^{(m)}\}$
- Given new inputs $X = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- Normalize the features of X using the mean and standard deviation from R
- Every $x^{(i)}$ is always treated the same, regardless of which other examples appear in the minibatch

Tips and Tricks for training GAN

- **Virtual Batch Normalization**

- Reference batch norm can overfit to the reference batch. A partial solution is virtual batch norm
- Fix a reference batch $R = \{r^{(1)}, r^{(2)}, \dots, r^{(m)}\}$
- Given new inputs $X = \{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- For each $x^{(i)}$
 - Construct a minibatch containing $x^{(i)}$ and all R
 - Compute mean and standard deviation of V
 - Normalize the features of $x^{(i)}$ using the mean and standard deviation

Tips and Tricks for training GAN

- **Use Adam optimizer**
 - Use SGD for discriminator & Adam for generator
- **Avoid Sparse Gradients: ReLU, MaxPool**
 - the stability of the GAN game suffers if you have sparse gradients
 - LeakyReLU = good (in both G and D)
 - For Downsampling, use: Average Pooling, Conv2d + stride
 - For Upsampling, use: Bilinear Interpolation, PixelShuffle

Tips and Tricks for training GAN

- **Use Soft and Noisy Labels**

- Default cost

- `cross_entropy(1., discriminator(data))`
`+ cross_entropy(0., discriminator(samples))`

- Label Smoothing (Salimans et al. 2016)

- For real ones (label=1), replace it with 0.9; For fake ones (label=0), keep it to 0.

- `cross_entropy(.9, discriminator(data))`
`+ cross_entropy(0., discriminator(samples))`

Tips and Tricks for training GAN

- **Use Soft and Noisy Labels**

```
cross_entropy(1.-alpha, discriminator(data))  
+ cross_entropy(beta, discriminator(samples))
```

$$D(\mathbf{x}) = \frac{(1 - \alpha)p_{\text{data}}(\mathbf{x}) + \beta p_{\text{model}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{model}}(\mathbf{x})}$$

- make the labels noisy for the discriminator:
occasionally flip the labels when training the discriminator

Tips and Tricks for training GAN

- **Track failures early**

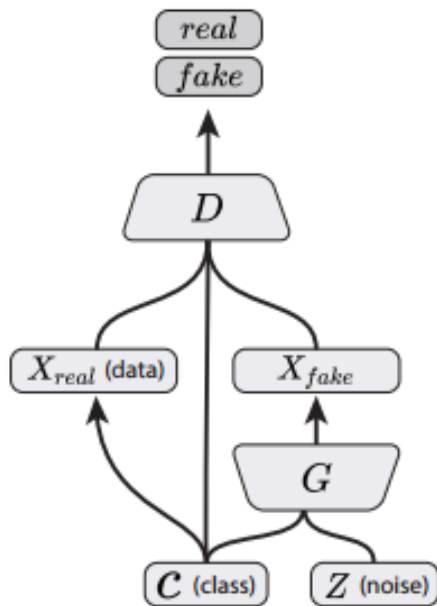
- D loss goes to 0: failure mode
- Check norms of gradients: > 100 is bad
- When training well, D loss has low variance and is going down. Otherwise, D loss is spiky
- If loss of G steadily decreases, then it's fooling D with garbage

Tips and Tricks for training GAN

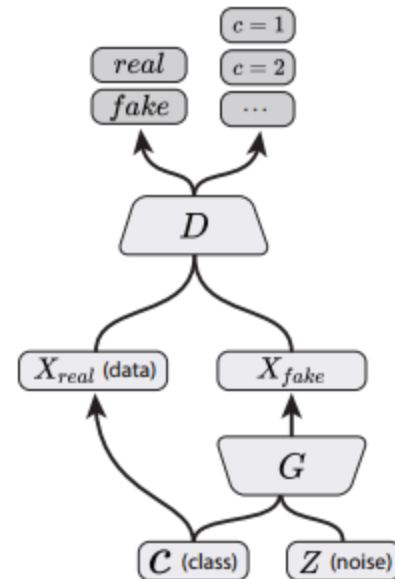
- **Discrete variables in Conditional GANs**
 - Use an Embedding layer
 - Add as additional channels to images
 - Keep embedding dimensionality low and upsample to match image channel size

Tips and Tricks for training GAN

- **Use label information when possible**
 - Used as conditioning variable
 - Auxiliary classifier GAN



Conditional GAN



AC-GAN

Tips and Tricks for training GAN

- **Balancing G and D**
 - Usually the discriminator “wins”
 - Good thing: theoretical justification are based on assuming D is perfect
 - Usually D is bigger and deeper than G
 - Sometimes run D more often than G . Mixed results
 - Do not try to limit D to avoid making it “too smart”
 - Use non-saturating cost
 - Use label smoothing

Research Directions

- Research direction of GANs
 - Better network structures
 - Better objective functions
 - Novel problem setups
 - Use of adversarial losses in other CV/ML applications
 - Theories on GAN