# EE655: Computer Vision & Deep Learning

## Lecture 16

Koteswar Rao Jerripothula, PhD
Department of Electrical Engineering
IIT Kanpur

# Overview

Neural Style Transfer

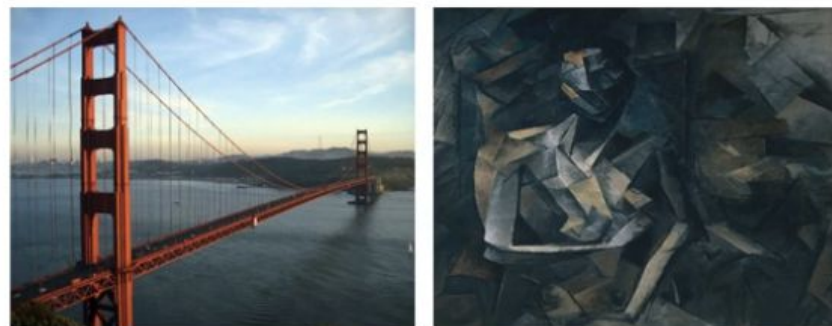Generative Adversarial Networks (GANs)

# Neural style transfer



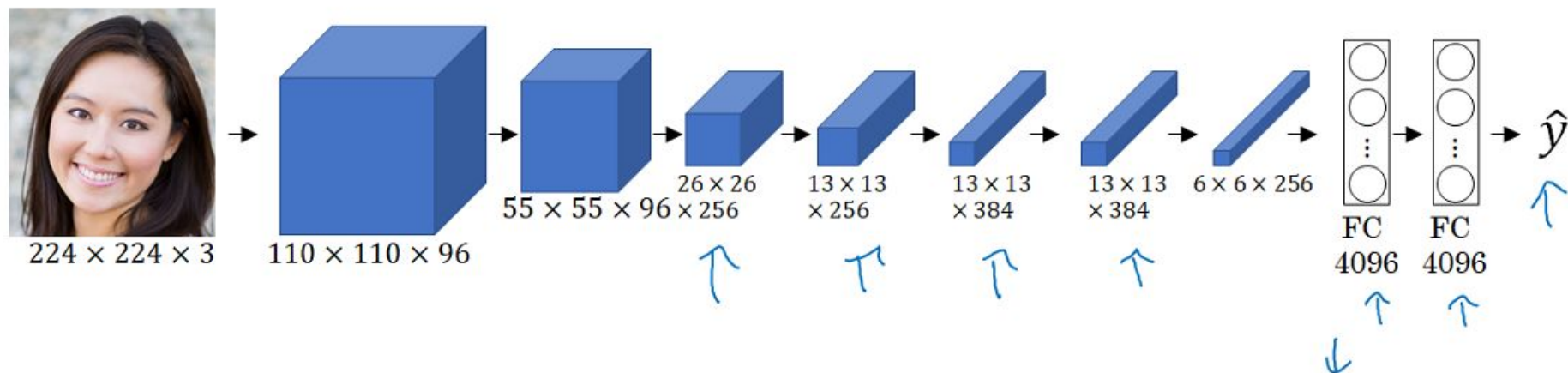Content (c)     Style (s)

Generated image (G)

Content (c)     Style (s)

Generated image (G)

# Visualizing what a deep network is learning



$224 \times 224 \times 3$    $110 \times 110 \times 96$    $55 \times 55 \times 96 \times 256$    $26 \times 26 \times 256$    $13 \times 13 \times 256$    $13 \times 13 \times 384$    $13 \times 13 \times 384$    $6 \times 6 \times 256$    FC 4096    FC 4096    $\hat{y}$

Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation.

Repeat for other units.

[Zeiler and Fergus., 2013, Visualizing and understanding convolutional networks]

# Visualizing deep layers: Layer 1



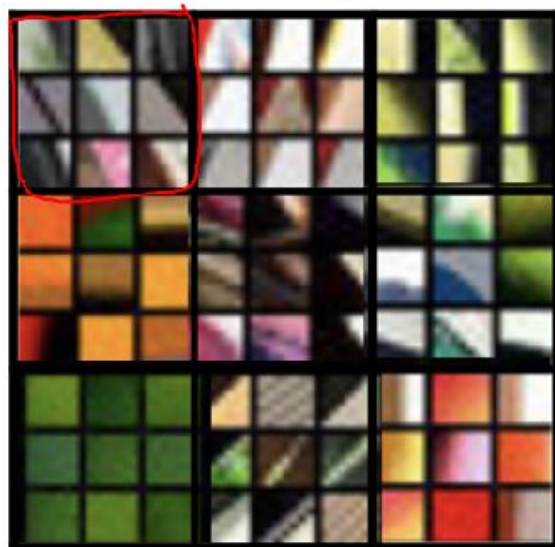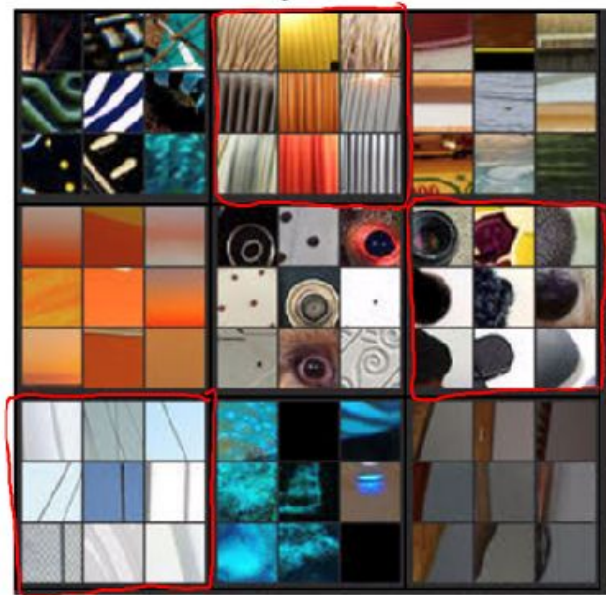Layer 1    Layer 2    Layer 3    Layer 4    Layer 5

# Visualizing deep layers: Layer 2



Layer 1    Layer 2    Layer 3    Layer 4    Layer 5

# Visualizing deep layers: Layer 3
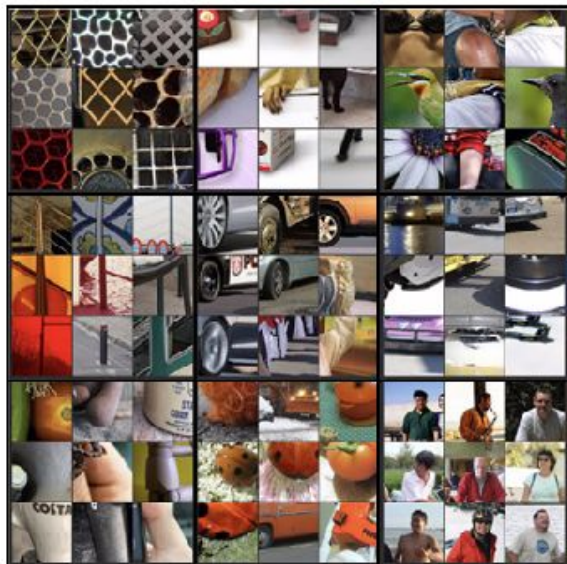


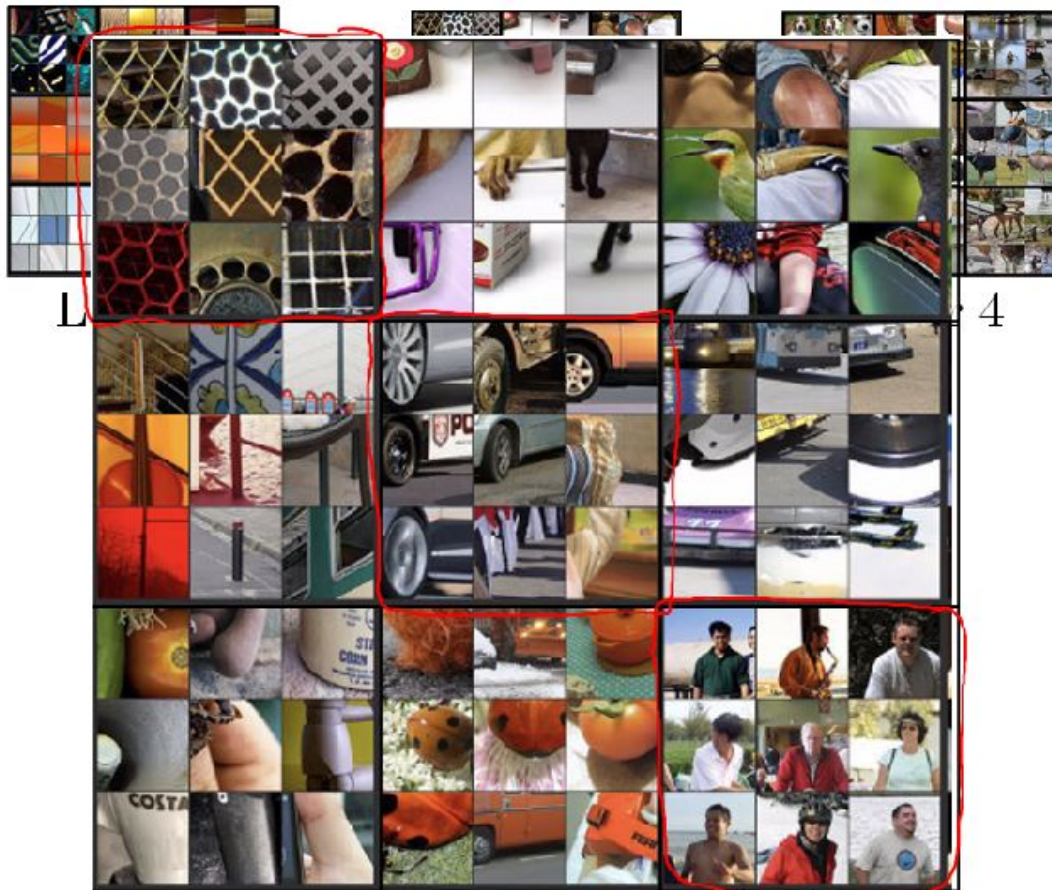Layer 1          Layer 2          Layer 3          Layer 4          Layer 5

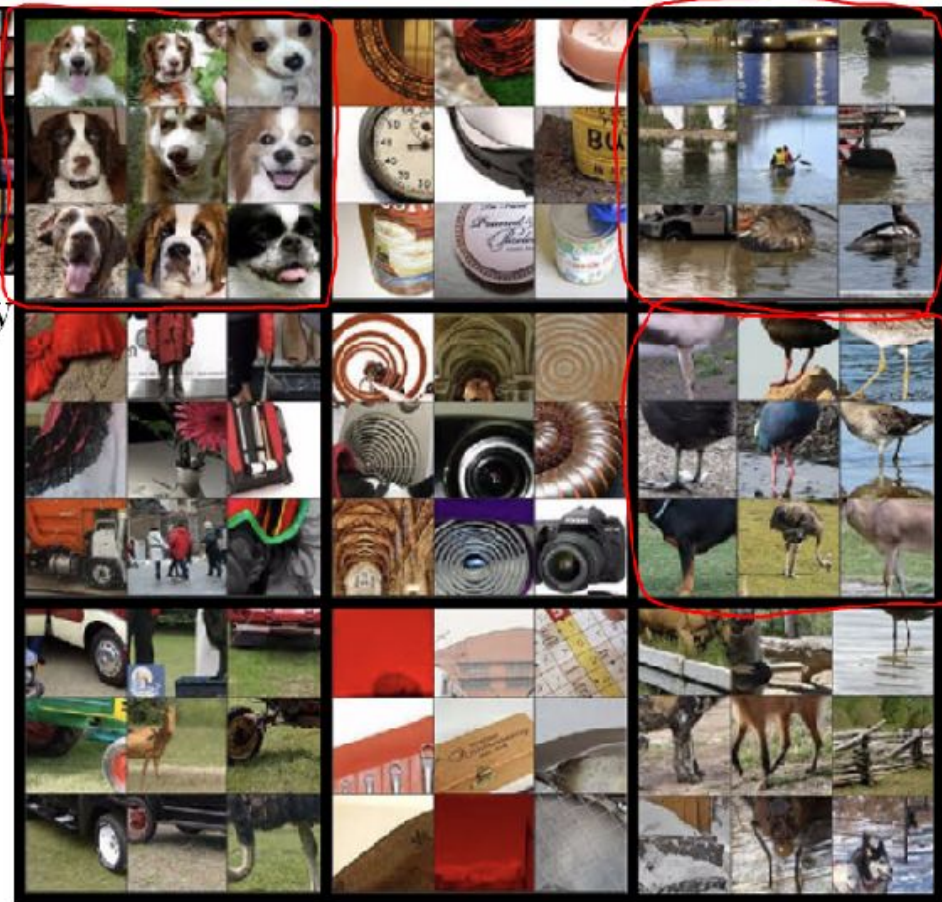# Visualizing deep layers: Layer 3



Layer 1    Layer 2    Layer 3    Layer 4    Layer 5

# Visualizing deep layers: Layer 4



Lay...

Layer 4

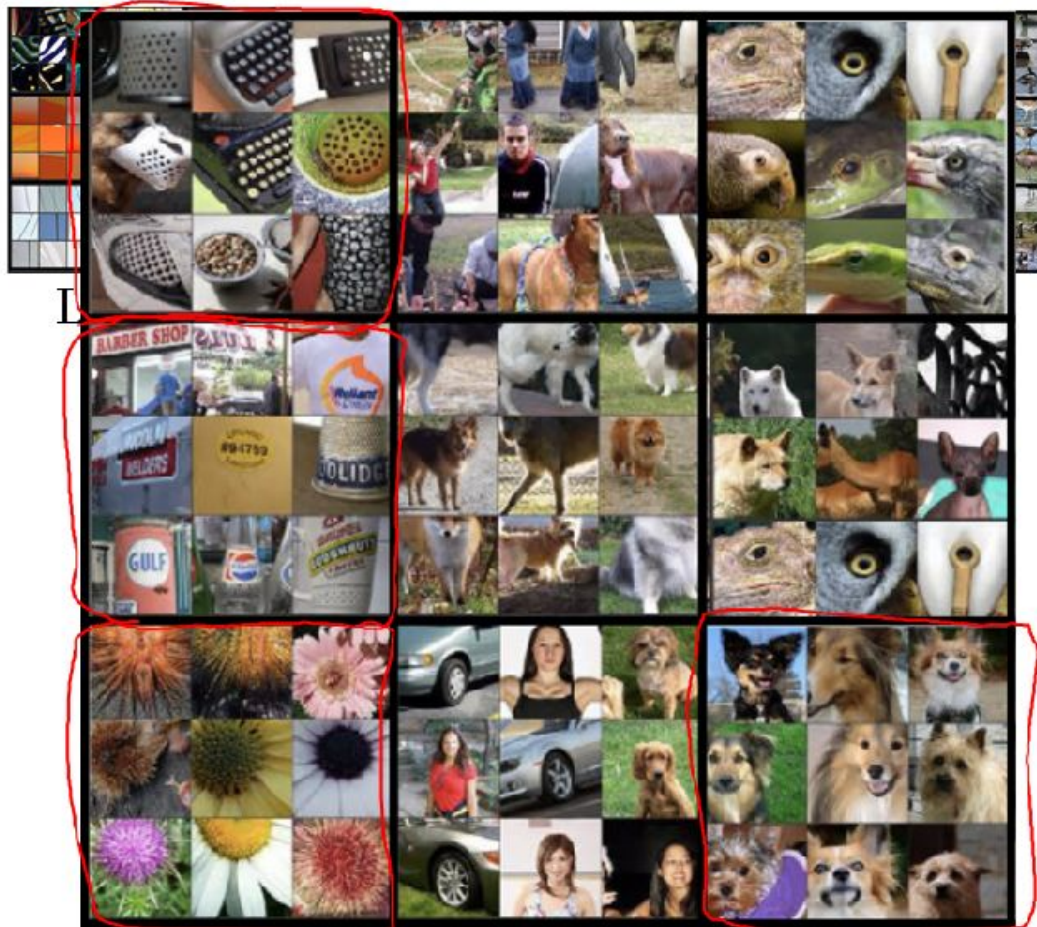Layer 5

# Visualizing deep layers: Layer 5



Layer 1

Layer 5

# Neural style transfer cost function



Content C  Style S

Generated image G ←

$$J(G) = \alpha \, J_{content}(C, G)$$

$$+ \beta \, J_{style}(S, G)$$

# Find the generated image G



1. Initiate G randomly

   G: $100 \times 100 \times 3$

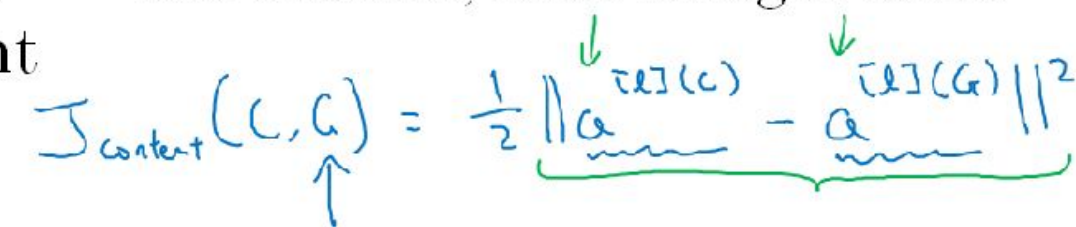   $\uparrow$
   RGB

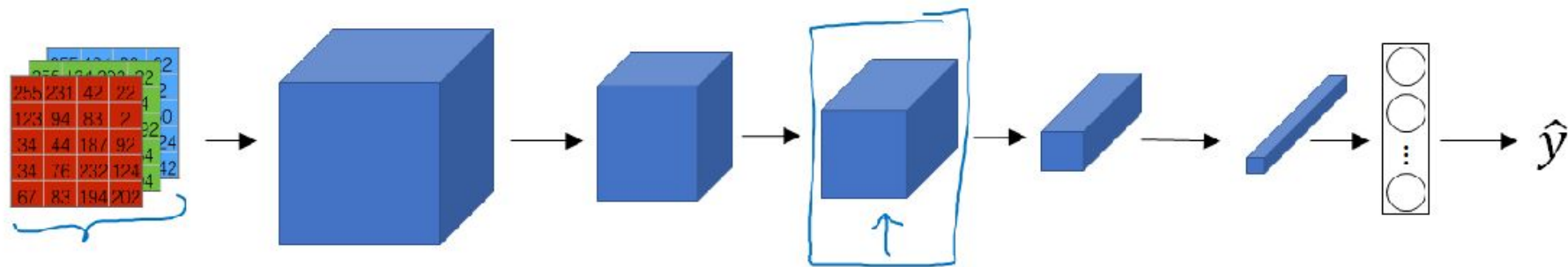2. Use gradient descent to minimize $J(G)$

$$G := G - \frac{d}{dG} J(G)$$



[Gatys et al., 2015. A neural algorithm of artistic style]

# Content cost function

$$J(G) = \alpha \, J_{content}(C, G) + \beta \, J_{style}(S, G)$$

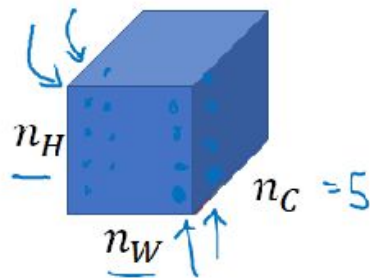- Say you use hidden layer $l$ to compute content cost.
- Use pre-trained ConvNet. (E.g., VGG network)
- Let $a^{[l](C)}$ and $a^{[l](G)}$ be the activation of layer $l$ on the images
- If $a^{[l](C)}$ and $a^{[l](G)}$ are similar, both images have similar content

$$J_{content}(C, G) = \frac{1}{2} \left\| a^{[l](C)} - a^{[l](G)} \right\|^2$$

[Gatys et al., 2015. A neural algorithm of artistic style]

# Meaning of the "style" of an image



Say you are using layer $l$'s activation to measure "style."
Define style as correlation between activations across channels.



How correlated are the activations across different channels?

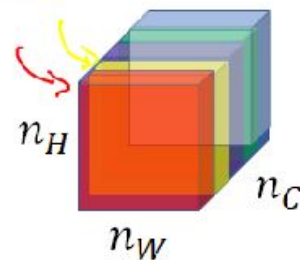[Gatys et al., 2015. A neural algorithm of artistic style]

# Intuition about style of an image

Style image



Generated Image



$n_H$

$n_C$

$n_W$

Correlated?

Uncorrelated

$n_H$

$n_C$

$n_W$

[Gatys et al., 2015. A neural algorithm of artistic style]

# Style matrix

$H$ $W$ $C$

Let $a^{[l]}_{i,j,k}$ = activation at $(i, j, k)$. $G^{[l]}$ is $n^{[l]}_c \times n^{[l]}_c$

$n_c$

$G^{[l]}_{kk'}$

$k = 1, \ldots, n^{[l]}_c$

$$G^{[l](S)}_{kk'} = \sum_{i=1}^{n^{[l]}_H} \sum_{j=1}^{n^{[l]}_W} a^{[l](S)}_{ijk} \, a^{[l](S)}_{ijk'}$$

$$G^{[l](G)}_{kk'} = \sum_{i=1}^{n^{[l]}_H} \sum_{j=1}^{n^{[l]}_W} a^{[l](G)}_{ijk} \, a^{[l](G)}_{ijk}$$

"Gram matrix"

$$J^{[l]}_{style}(S, G) = \frac{1}{(\cdots)} \left\| G^{[l](S)} - G^{[l](G)} \right\|^2_F$$

$\beta$

$$= \frac{1}{\left(2 n^{[l]}_H n^{[l]}_W n^{[l]}_c\right)^2} \sum_{k} \sum_{k'} \left( G^{[l](S)}_{kk'} - G^{[l](G)}_{kk'} \right)^2$$

[Gatys et al., 2015. A neural algorithm of artistic style]

# Style cost function

$$\left\| G^{[l](S)} - G^{[l](G)} \right\|_F^2$$

$$J_{style}^{[l]}(S,G) = \frac{1}{\left(2 n_H^{[l]} n_W^{[l]} n_C^{[l]}\right)^2} \sum_k \sum_{k'} \left(G_{kk'}^{[l](S)} - G_{kk'}^{[l](G)}\right)^2$$

$$J_{style}(S,G) = \sum_l \lambda^{[l]} J_{style}^{[l]}(S,G)$$

$$J(G) = \alpha\, J_{content}(C,G) + \beta\, J_{style}(S,G)$$

$$G$$

[Gatys et al., 2015. A neural algorithm of artistic style]

# Generative

# Adversarial



Generates data
(Creates fake data)



**Generator** and **discriminator**, each
competing to win.

Generator trying to fake
and
Discriminator, trying not to be fooled.

**Z**

| 0.1 |
|-----|
| -0.3 |
| 0.6 |
| ... |
| ... |
| ... |
| -0.7 |

Random Noise
(Latent vector)

**Real Images**

Training goal for the generator is to maximize the probability of the discriminator making a mistake.
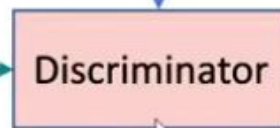
50/50

**Generator**

**Fake Images**

**Discriminator**
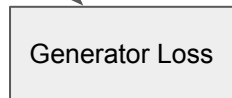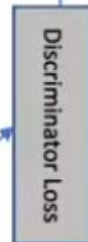
Training goal for the discriminator is to maximize the probability of identifying real vs. fake images correctly.

Real

Fake

Discriminator Loss

Generator Loss

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, $k$, is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

---

**for** number of training iterations **do**

    **for** $k$ steps **do**

        • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.

        • Sample minibatch of $m$ examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.

        • Update the discriminator by ascending its stochastic gradient:

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

    **end for**

    • Sample minibatch of $m$ noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.

    • Update the generator by descending its stochastic gradient:

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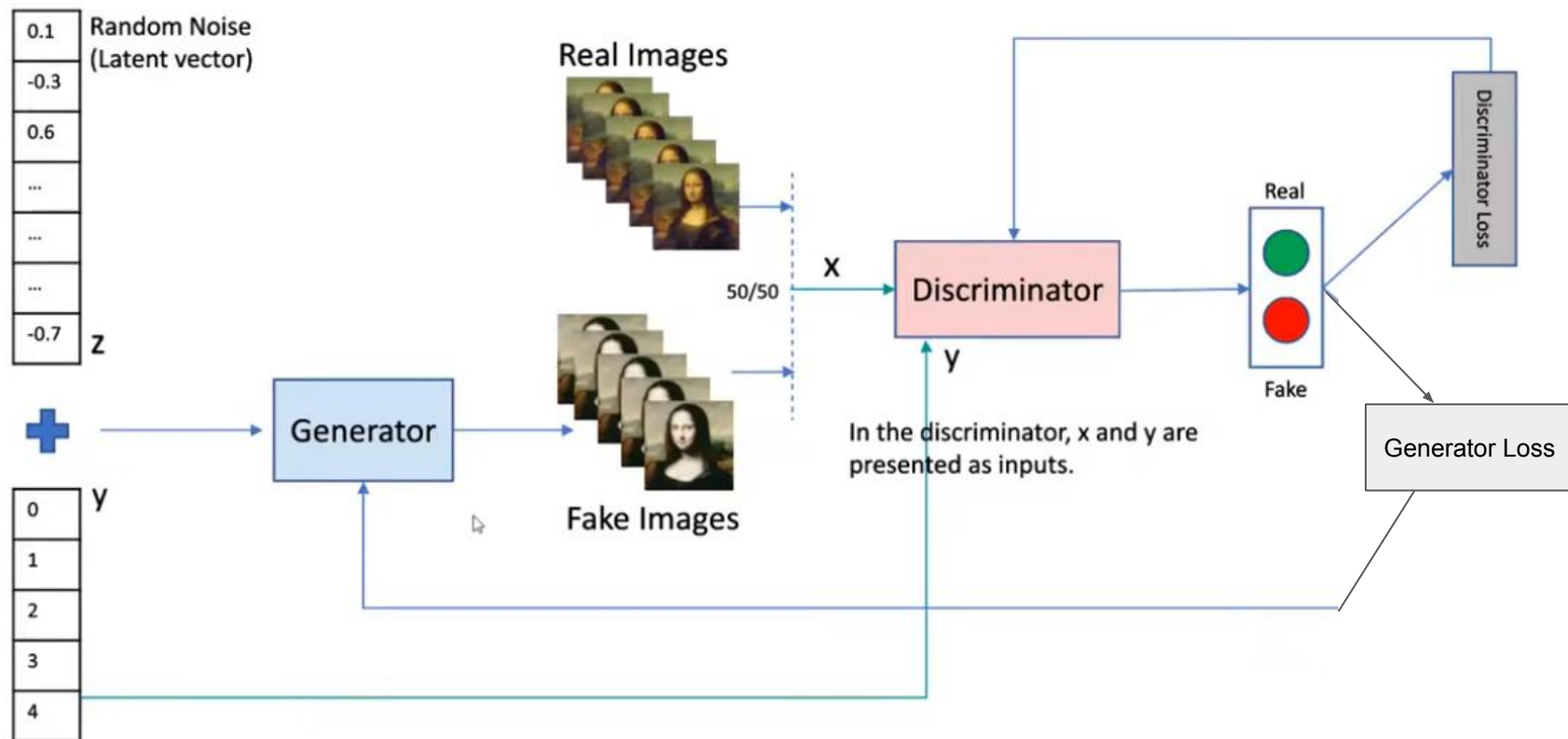
**end for**

**5 steps to training a GAN**

1. Define GAN architecture based on the application

2. Train the discriminator to distinguish between real vs fake data

3. Train the generator to fake data that can fool the discriminator

4. Continue discriminator and generator training for multiple epochs

5. Save generator model to create new, realistic fake data

NOTE: When training the discriminator, hold the generator values constant; and when training the generator, hold the discriminator values constant. Each should train against a static adversary.

# Conditional Generative Adversarial Network



0.1
-0.3
0.6
...
...
...
-0.7 z

Random Noise
(Latent vector)

Real Images

50/50

X

Discriminator

y

In the discriminator, x and y are presented as inputs.

Real

Fake

Discriminator Loss

Generator Loss

Generator

Fake Images

0
1
2
3
4

y

Conditional data (y): can be class labels
or data from other modalities.

# Applications of Conditional GANs

## Image-to-Image Translation: Pix2Pix GAN

- Takes an image as input and maps it to a generated output image with different properties.

- Example: Train an image-to-image GAN to take sketches of handbags and turn them into photorealistic images of handbags.

- The system requires pairwise correspondences between images during training.



Labels to Street Scene

input    output

Aerial to Map

Day to Night

input    output

# Applications of Conditional GANs

## Super-resolution

Increase the resolution of images, adding detail where necessary to fill in blurry areas.



bicubic (21.59dB/0.6423)    SRResNet (23.53dB/0.7832)    SRGAN (21.15dB/0.6868)    original

The GAN-generated image looks very similar to the original image, but if you look closely at the headband...
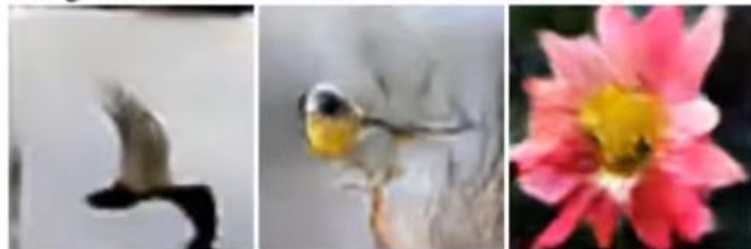
**Applications of Conditional GANs**

**Text-to-Image Synthesis**

Take text as input and produce images as described by the text.

## CycleGAN Architecture

Generator 1 → Generator 2

Discriminator → Real / Fake Zebra?
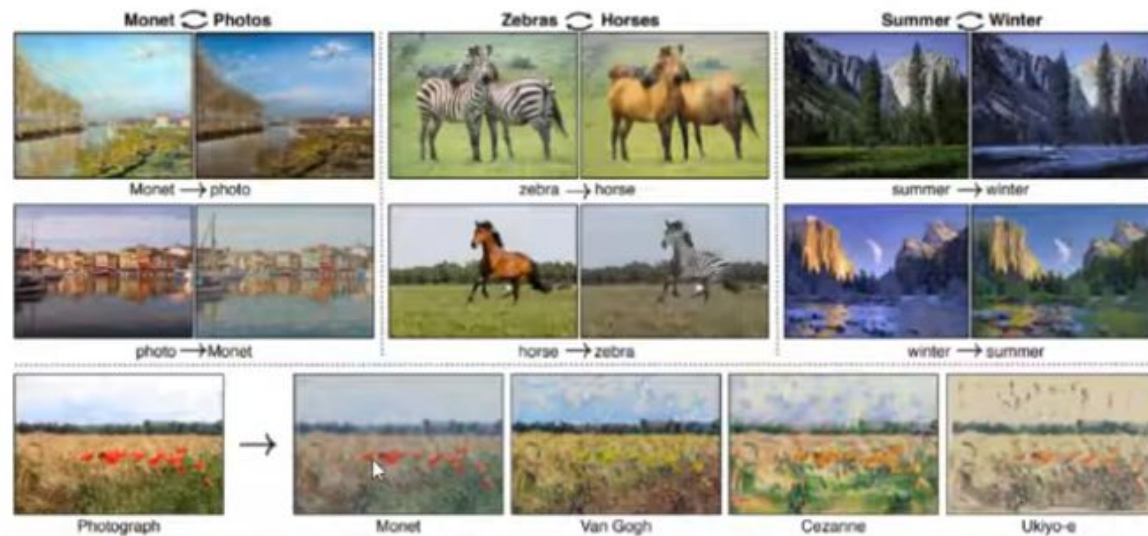
## CycleGAN Generator Objectives

1. Ensure the translated image looks like a zebra.
   a. This is trained using the GAN objective with the discriminator.
2. Ensure the translated image still looks mostly like the original.
   a. This is trained using a reconstruction objective with the second generator.
   b. This is the novel cycle-consistency loss.

# Applications of Conditional GANs

## CycleGAN

- Transform images from one set into images that could belong to another set.

- Example: Convert an image of a horse into an image of a zebra.

- The training data for the CycleGAN is simply two sets of images (e.g., a set of horse images and a set of zebra images).

- The system requires no labels or pairwise correspondences between images for training.

# References

https://www.youtube.com/watch?v=R39tWYYKNcI&list=PLkDaE6sCZn6Gl29AoE31iwdVwSG-KnDzF&index=37

https://www.youtube.com/watch?v=ChoV5h7tw5A&list=PLkDaE6sCZn6Gl29AoE31iwdVwSG-KnDzF&index=38

https://www.youtube.com/watch?v=xY-DMAJpIP4&list=PLkDaE6sCZn6Gl29AoE31iwdVwSG-KnDzF&index=39

https://www.youtube.com/watch?v=b1I5X3UfEYI&list=PLkDaE6sCZn6Gl29AoE31iwdVwSG-KnDzF&index=40

https://www.youtube.com/watch?v=W5NPlZzebO0

https://www.youtube.com/watch?v=-8hfnlxEPn4