

Deep Generative Modeling

Saqib Shamsi

Unsupervised Learning

- No labels are provided during training
- General objective: inferring a function to describe hidden structure from unlabeled data
 - Clustering (discrete labels)
 - Representation/feature learning (continuous vectors)
 - Dimensionality reduction (lower-dimensional representation)
 - Density estimation (continuous probability)

Unconditional Data Generation

- **Problem:** Given some high dimensional data, generate new samples from the same distribution.



Training data $\sim p_{\text{data}}(x)$

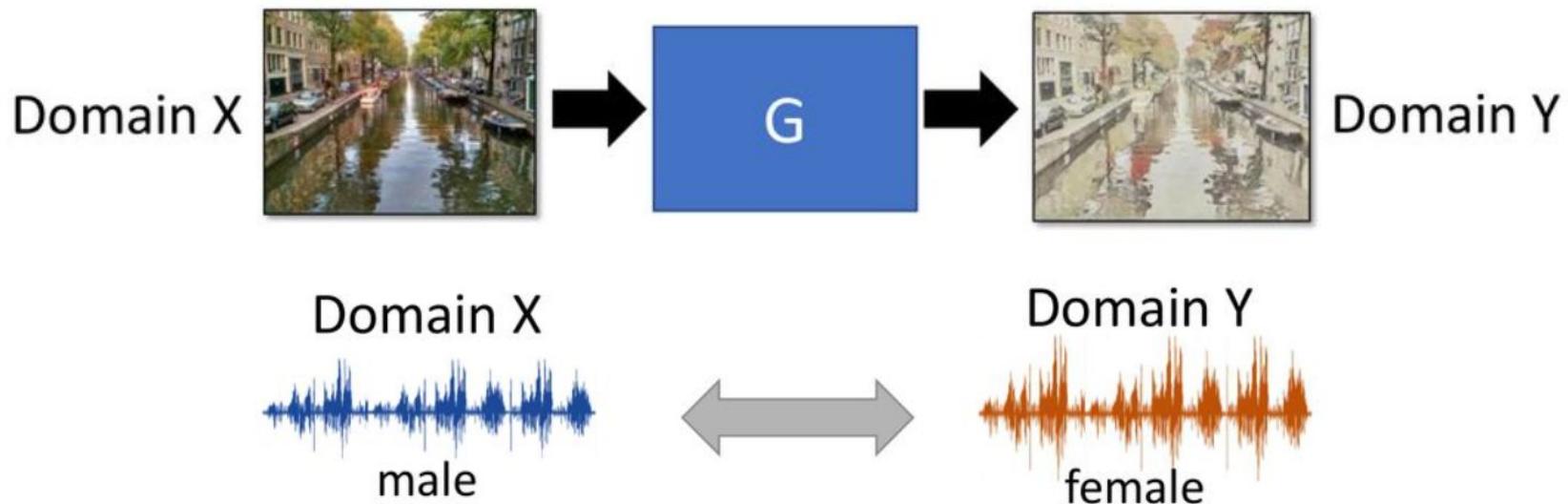


Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

Conditional Data Generation

- **Problem:** Given some high dimensional data, generate new data subject to some constraint.



[Source](#)

Unconditional Image Generation

What About This Solution?

```
import glob

def augment(image):
    """
    Applies random rotations, brightness, contrast,
    flips, shearing, colour-shift and noise
    """
    # Code for augmenting `image`
    return augmented_image

def generate(given_data_path):
    images = glob.glob(given_data_path + '/*.jpg')
    for image in images:
        yield augment(image)
```

What About This One?

```
def generate_all(num_pixels, intensity_levels):
    """
    Generates all possible images with the number
    of pixels equal to `num_pixels` and of given
    `intensity_levels`
    """
    # Code to generate all possible images
    return all_images

def main():
    DESIRED_SIZE = 100 * 100 * 3
    INTENSITIES = 256
    all_images = generate_all(DESIRED_SIZE, INTENSITIES)
    # Pick whatever images you want from `all_images`
```

What About This One?

```
def generate_all(num_pixels, intensity_levels):
    """
    Generates all possible images with the number
    of pixels equal to `num_pixels` and of given
    `intensity_levels`
    """
    # Code to generate all possible images
    return all_images
```

```
def main():
    DESIRED_SIZE = 100 * 100 * 3
    INTENSITIES = 256
    all_images = generate_all(DESIRED_SIZE, INTENSITIES)
    # Pick whatever images you want from `all_images`
```

$$1.58 \times 10^{72245}$$

What About This One?

```
def generate_all(num_pixels, intensity_levels):
    """
    Generates all possible images with the number
    of pixels equal to `num_pixels` and of given
    `intensity_levels`
    """
    # Code to generate all possible images
    return all_images
```

```
def main():
    DESIRED_SIZE = 100 * 100 * 3
    INTENSITIES = 256
    all_images = generate_all(DESIRED_SIZE, INTENSITIES)
    # Pick whatever images you want from `all_images`
```

Estimated #protons + #electrons
in the observable universe: 10^{80}

1.58×10^{72247}

What About This One?

```
def generate_all(num_pixels, intensity_levels):
    """
    Generates all possible images with the number
    of pixels equal to `num_pixels` and of given
    `intensity_levels`
    """
    # Code to generate all possible images
    return all_images
```

Moreover,
Undesired Images

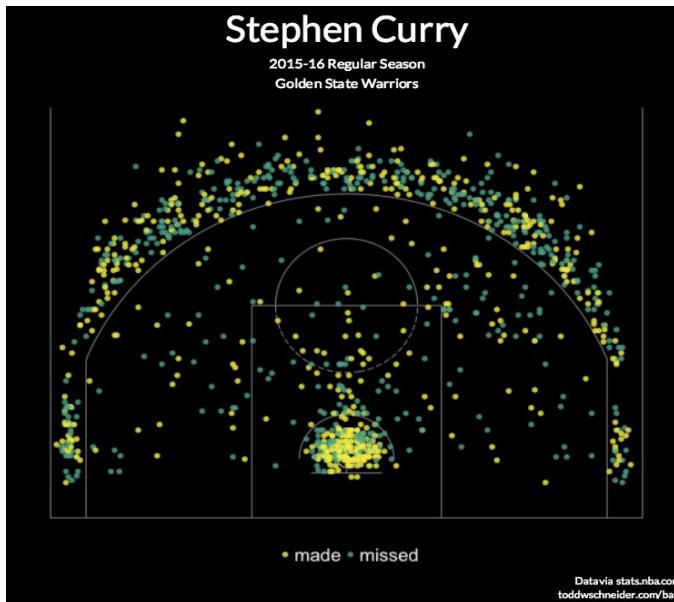
```
def main():
    DESIRED_SIZE = 100 * 100 * 3
    INTENSITIES = 256
    all_images = generate_all(DESIRED_SIZE, INTENSITIES)
    # Pick whatever images you want from `all_images`
```

What We Want

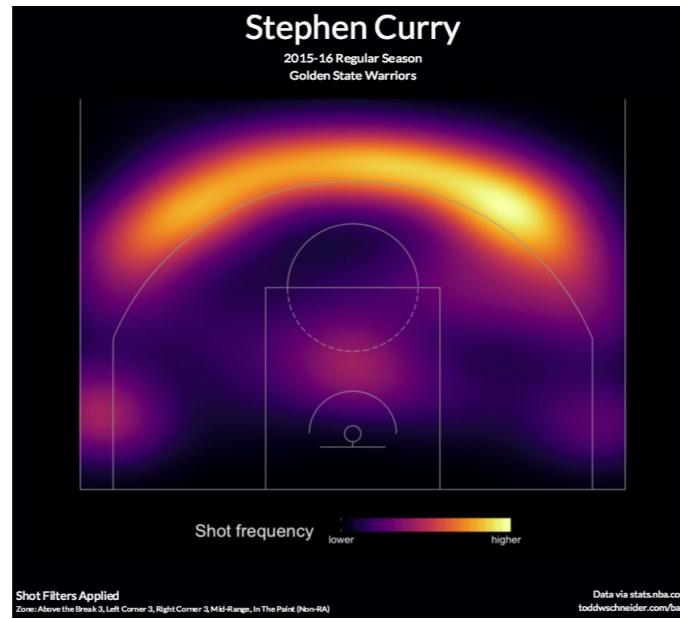
- Meaningfully different samples. Not the same data in augmented form
 - Samples that have the same underlying pattern as the training data, but should be different from the training data
 - Eg, digits in a different handwriting, cats in new poses, paintings with different scenery and brushstrokes
- The generation process must be tractable

Density Estimation

- Goal: Estimate the probability density function for given data.



[Image Source](#)

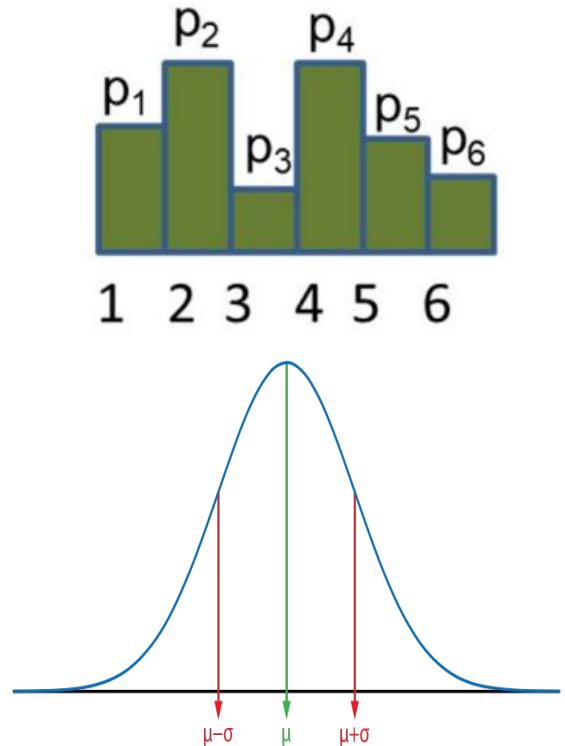


[Image Source](#)

Generative Model

- Simple Generative Models:
 - **Multinomial PMF:** For discrete data. Can be modeled as a probability table.
 - Probability of observing a data point can be obtained directly from the table
 - **Gaussian PDF:** For continuous data. Distribution is parameterized by mean μ and standard deviation σ
 - Probability of observing a data point x can be obtained by

$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$



[Image Source](#)

[Source](#)

Sampling From a Distribution

- Since our goal was to generate new data, we have to sample from the density function.
- Sampling from a multinomial PMF:
 - From the model probabilities p_1, \dots, p_k , compute the cumulative distribution
$$F_i = p_1 + \dots + p_i \quad \text{for all } i \in \{1, \dots, k\}$$
 - Draw a uniform random number $u \sim [0, 1]$
 - Return the smallest i such that $u \leq F_i$
- We have algorithms to sample from a normal distribution given μ and σ
- We have algorithms to sample from uniform and other simple distributions

Generative Models for ‘Simple’ Data

- We are given some data $X = \{x_i\}_{i=1}^N$
- We choose a model $P(x;w)$ (that we can sample from later) to model the distribution of X
- We estimate ‘ w ’ using some method such that $P(x;w)$ best fits the observations X , hoping that it will account for data not in X
- Finally, we sample from P to generate new data such that the new points generated follow the same pattern as X

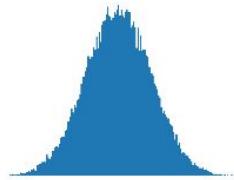
What About High Dimensional Data?

- We are interested in generating high dimensional data like images, audio etc
- We can sample from simple distributions but they do not capture the distribution of high dimensional data well
- Neural networks have the capacity to model the distribution of high dimensional data, but how do we sample from them?
- Can we combine the two above?

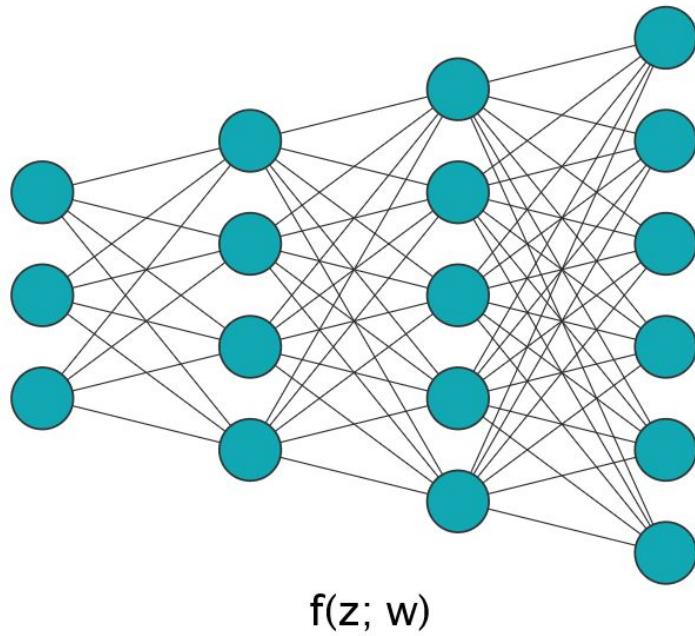
Mapping Distributions

- We create a function `f` such that it maps a simple distribution $Z=\{z\}$ to a high dimensional distribution $X=\{x\}$. $f: z \rightarrow x$
- `f` would be a neural network which has the capacity to model the distribution of `X`
- When we want to sample new data, we sample z from Z , provide it as input to f , which gives a data point x , hoping the generated x is different from the ones in X

Putting it all Together



$$z \sim p(z)$$



$$x \sim p(x)$$

Learning Network Parameters

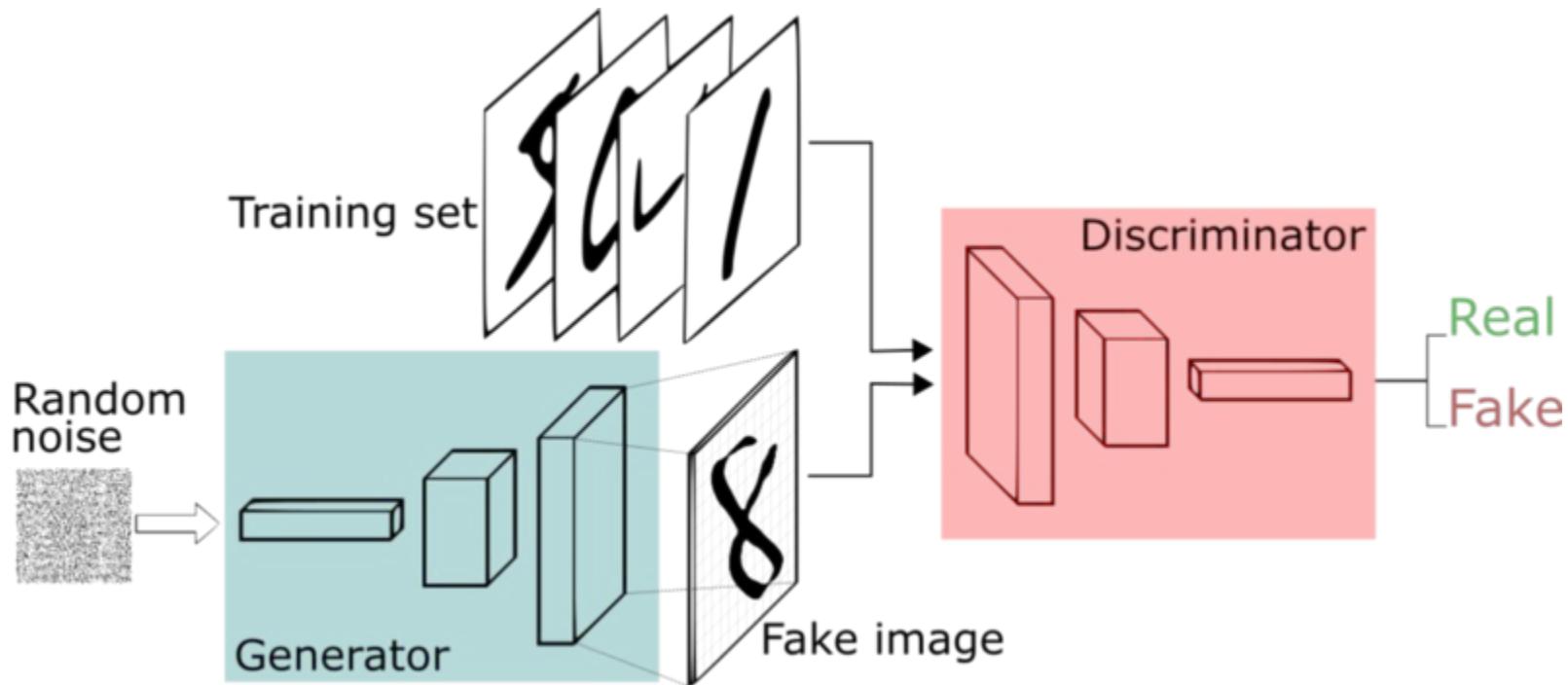
- Frameworks to learn the parameters of the generative network 'f':
 - Generative Adversarial Networks (GAN)
 - Variational Autoencoders (VAE)

Learning Network Parameters

Generative Adversarial Networks

- [Goodfellow et al](#), 2014
- Goal is to model the data distribution
- Trained using a pair of networks playing a minmax game
- They are termed as “adversaries” because they have conflicting loss functions
- New data are sampled by “seeding” the trained generator

GAN Architecture



[Image Source](#)

Vanilla GAN Loss

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

- Minmax game being played between a generator (G) and a discriminator (D)
- D tries to maximize the log-likelihood for the binary classification problem of real (1) vs fake (0)
- G tries to minimize the log-probability of its samples being classified as fake (0) by D (“fooling” the discriminator)

Vanilla GAN Training

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

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

end for

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

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

end for

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

Problems with Training Vanilla GANs

- Non-convergence
- A balancing act
 - If the discriminator gets too strong too fast, generator does not learn
 - If the discriminator is weak generator does not receive reliable feedback
- Too sensitive to the choice of hyperparameters. Require “hacks” to behave well
- Mode Collapse: generated samples lack diversity

Non-Convergence

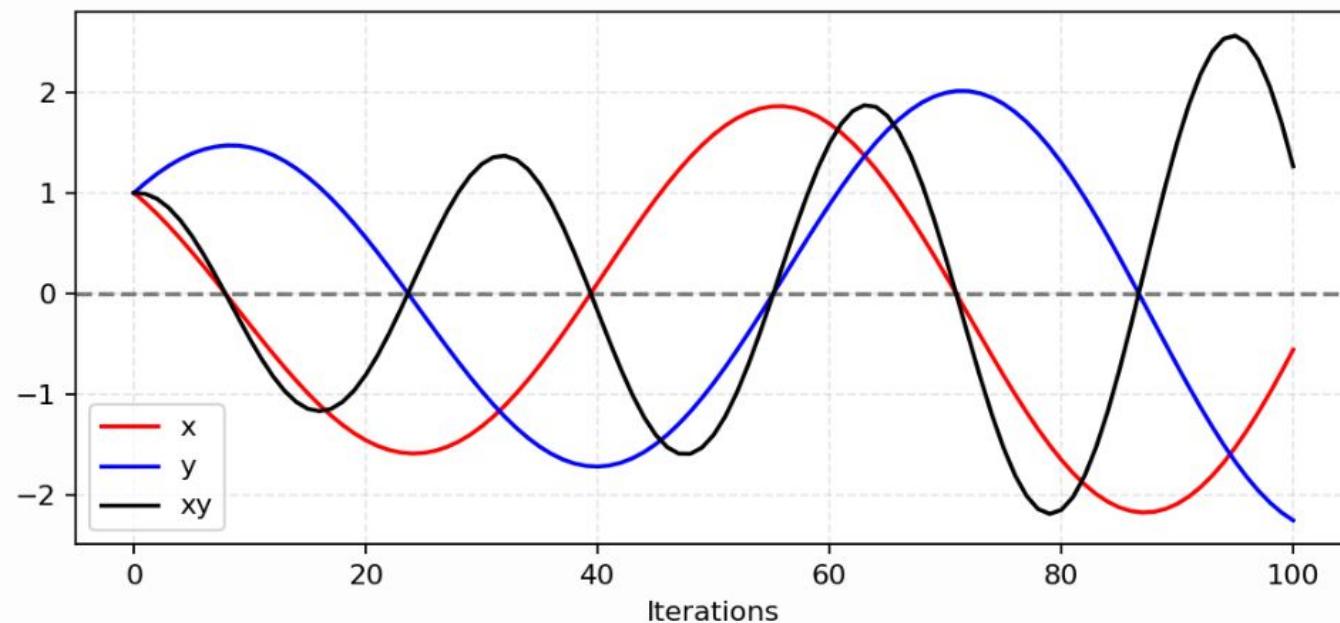
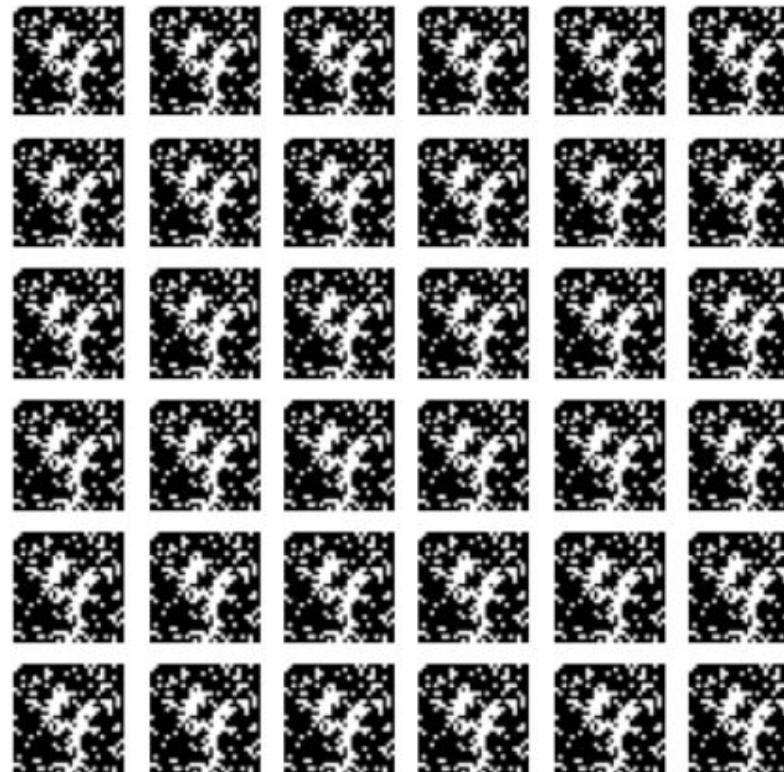


Fig. 3. A simulation of our example for updating x to minimize xy and updating y to minimize $-xy$. The learning rate $\eta = 0.1$. With more iterations, the oscillation grows more and more unstable.

Mode Collapse

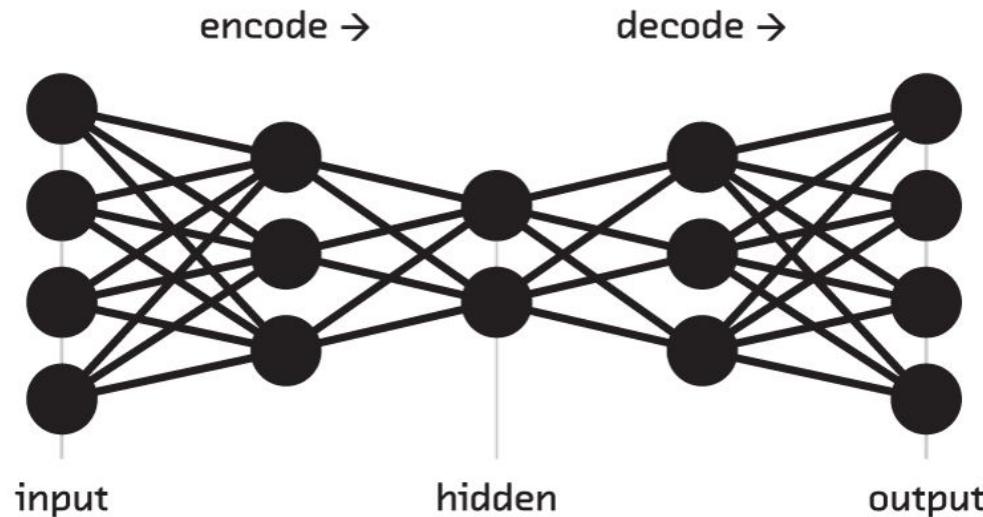


Learning Network Parameters

Variational Autoencoders

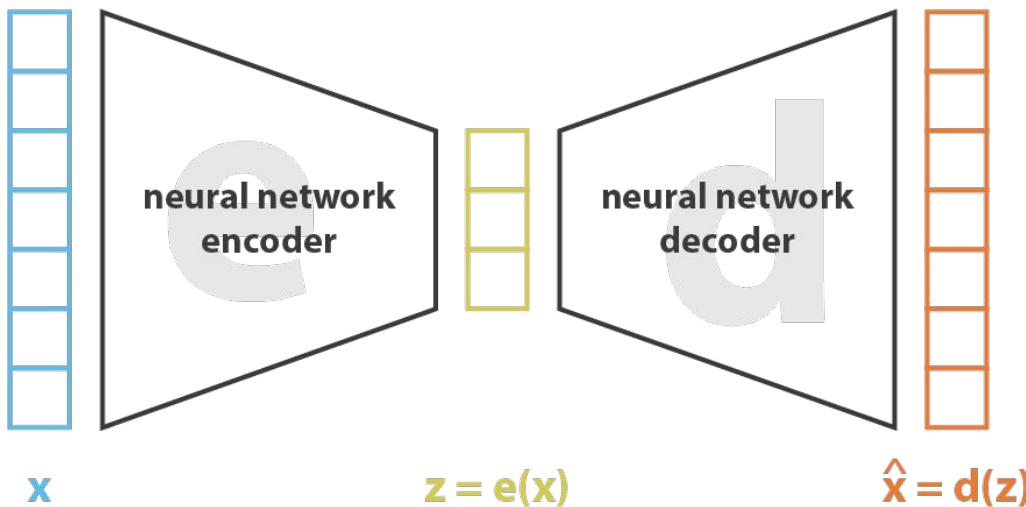
- [Kingma and Welling, 2013](#)
- Goal is to compress dataset information to a latent space
- Autoencoder whose encodings distribution is regularised
- Trained using gradient based methods
- New data are sampled by “seeding” the latent space

Autoencoder



[Image Source](#)

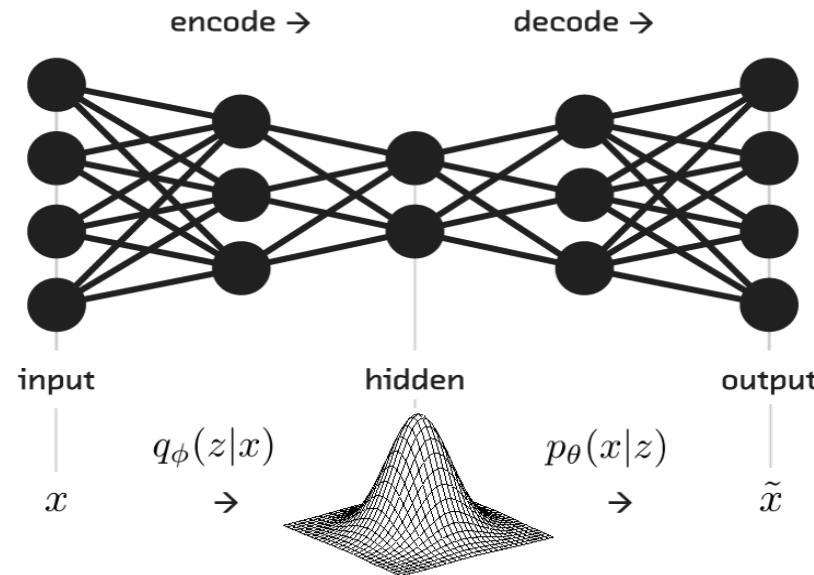
Autoencoder Loss



$$\text{loss} = \| x - \hat{x} \|^2 = \| x - d(z) \|^2 = \| x - d(e(x)) \|^2$$

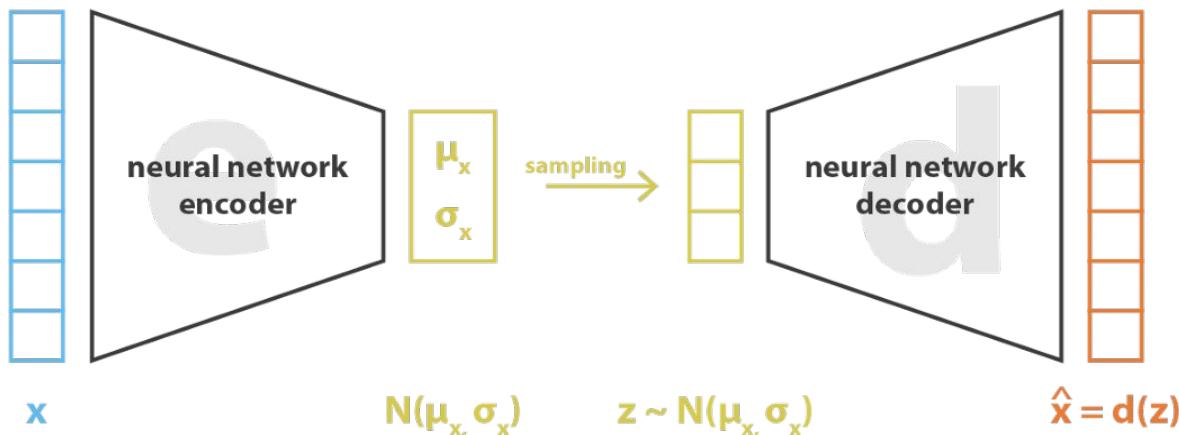
[Image Source](#)

Variational Autoencoder



[Image Source](#)

Variational Autoencoder Loss



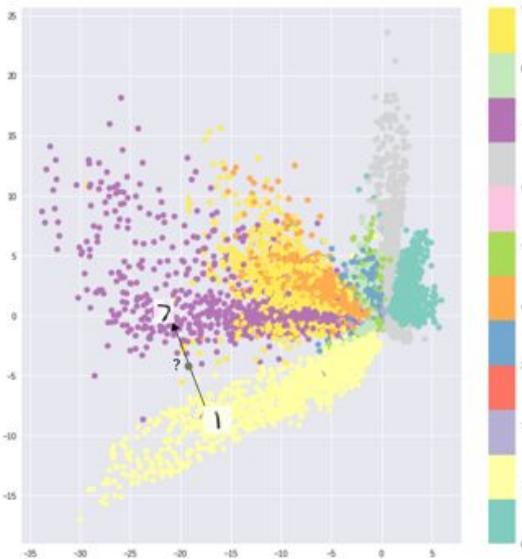
$$\text{loss} = \|x - \hat{x}\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)] = \|x - d(z)\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$

[Image Source](#)

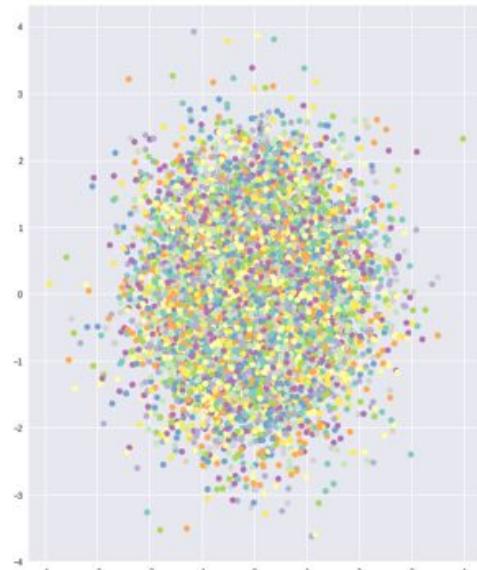
$$-\frac{1}{2} \sum (1 + \log(\sigma_x^2) - \mu_x^2 - \sigma_x^2)$$

The Need to Regularize the Latent Space

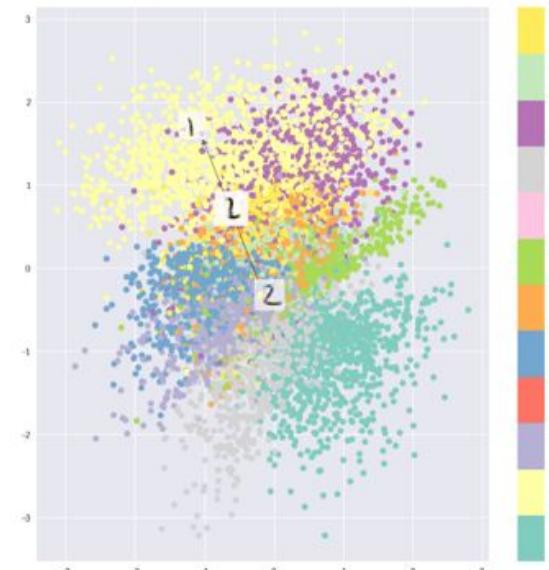
Only reconstruction loss



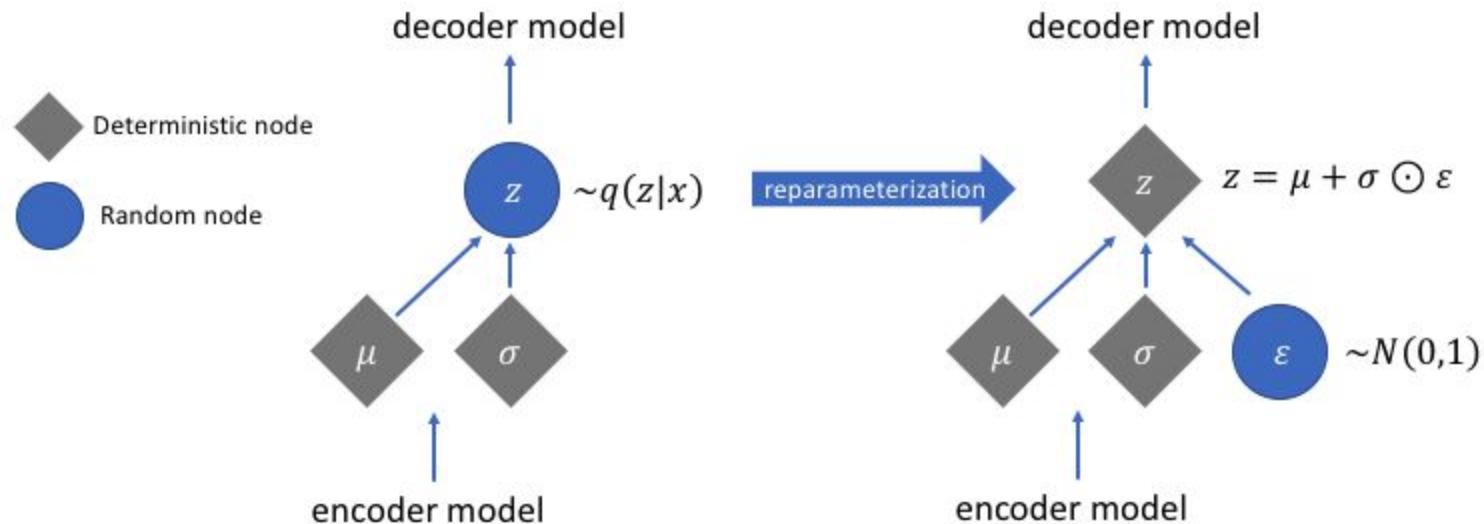
Only KL divergence



Combination



Reparameterization Trick



[Image Source](#)

Data Generation



(a) 2-D latent space

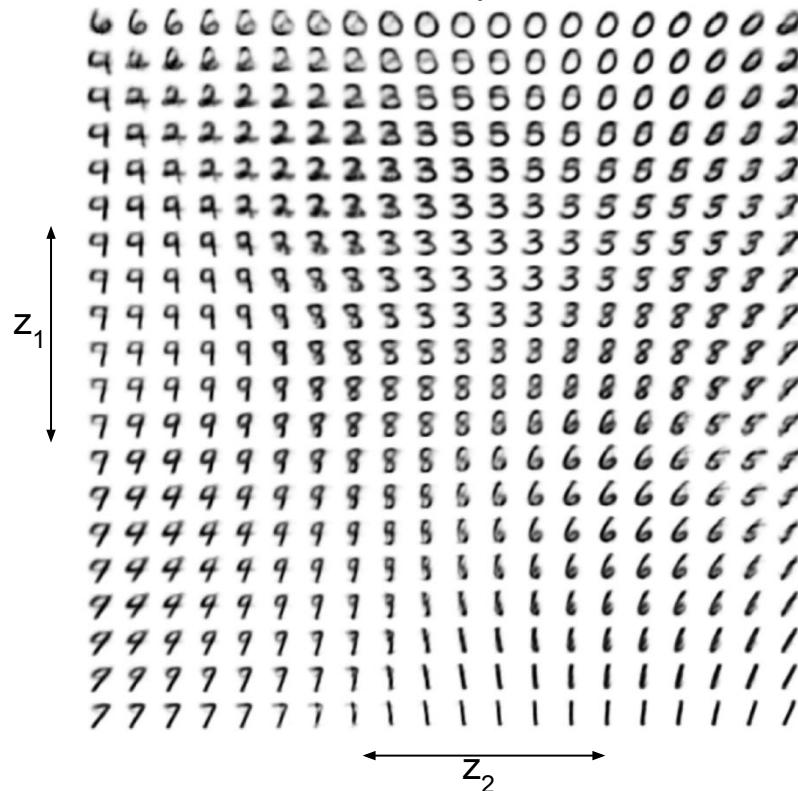
(b) 5-D latent space

(c) 10-D latent space

(d) 20-D latent space

Data Generation

2d latent space z



Applications

Data Generation

GAN Progress on Face Generation



Ian Goodfellow
@goodfellow_ian



4.5 years of GAN progress on face generation.

arxiv.org/abs/1406.2661 arxiv.org/abs/1511.06434

arxiv.org/abs/1606.07536 arxiv.org/abs/1710.10196

arxiv.org/abs/1812.04948



6:10 AM · Jan 15, 2019 · Twitter Web Client

[Link](#)

Generated Face Images



Generated Car Images



Generated Cat and Church Images

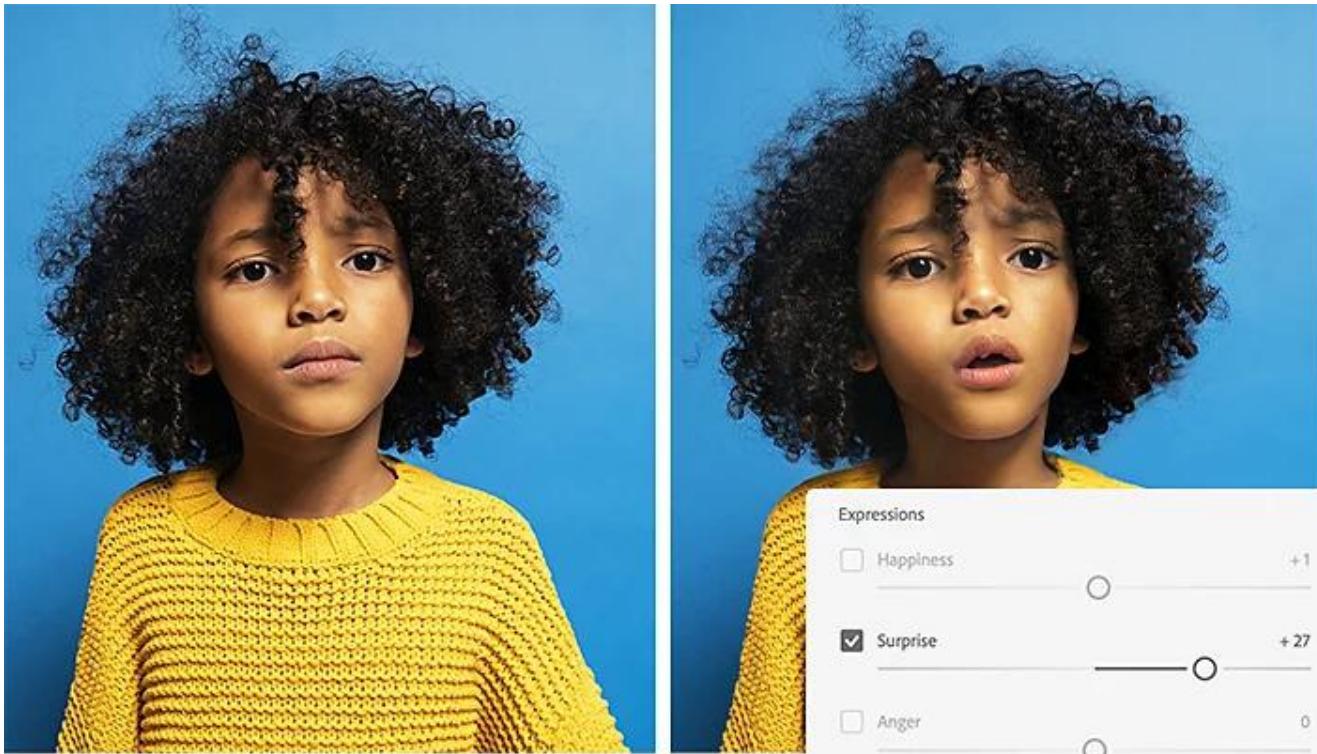


rob haroldone kanleor



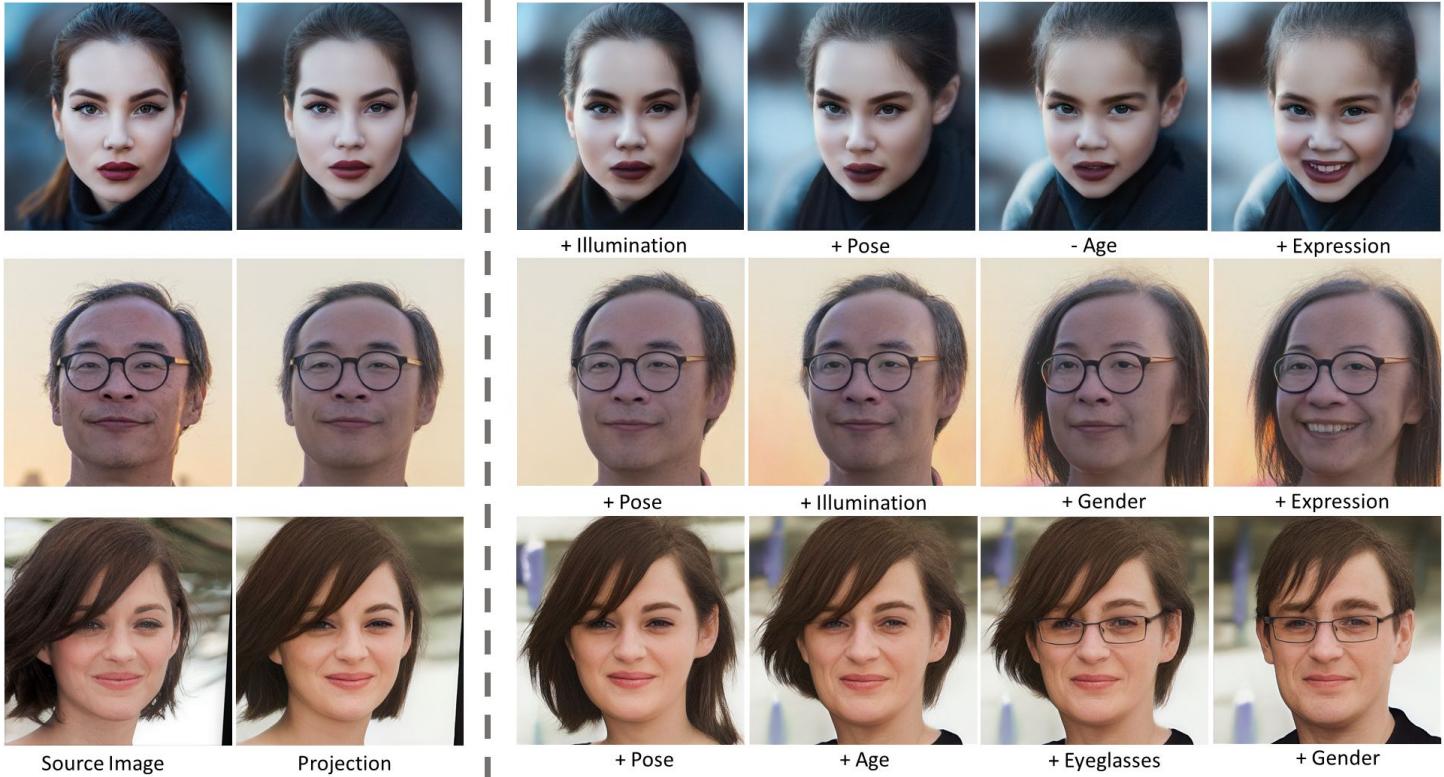
Image Editing

Image Editing



[Adobe Neural Filters](#)

Image Editing





Gender: ♂

Glasses: ✘

Year: 2010

Width: 400px

Bald: ✘

Beard: ✘

Age: +4

Dimensions: 100%

Mapping:

Left → Right: 0

Right → Left: 100%

Down → Up: 0

Up → Down: 0

Eye Width: 0

Pupil Width: 0



Image to Embedding

Algorithm 1: Latent Space Embedding for GANs

Input: An image $I \in \mathbb{R}^{n \times m \times 3}$ to embed; a pre-trained generator $G(\cdot)$.

Output: The embedded latent code w^* and the embedded image $G(w^*)$ optimized via F' .

```
1 Initialize latent code  $w^* = w$ ;  
2 while not converged do  
3    $L \leftarrow L_{percept}(G(w^*), I) + \frac{\lambda}{N} \|G(w^*) - I\|_2^2$  ;  
4    $w^* \leftarrow w^* - \eta F'(\nabla_{w^*} L)$ ;  
5 end
```

Image to Embedding

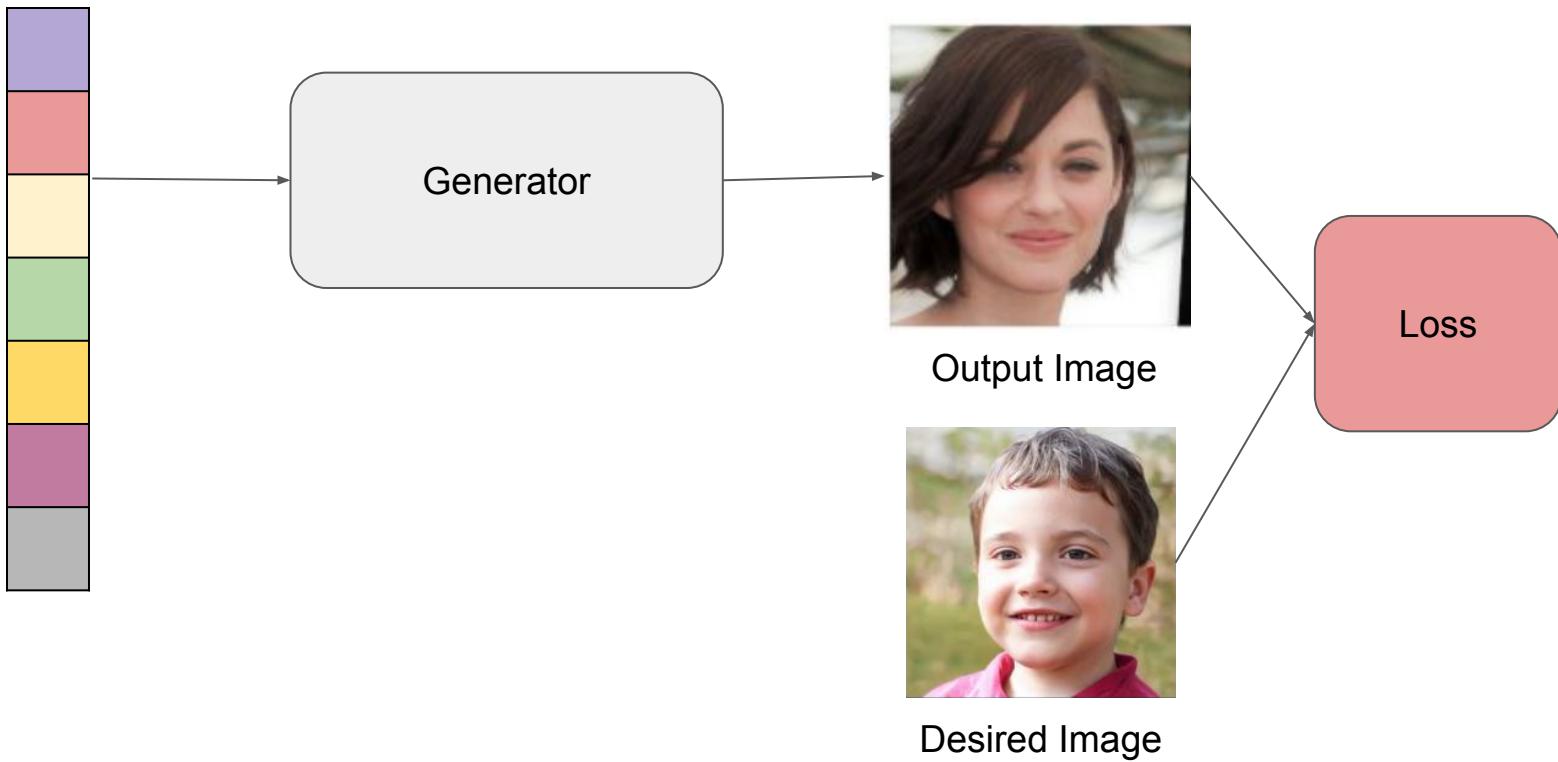
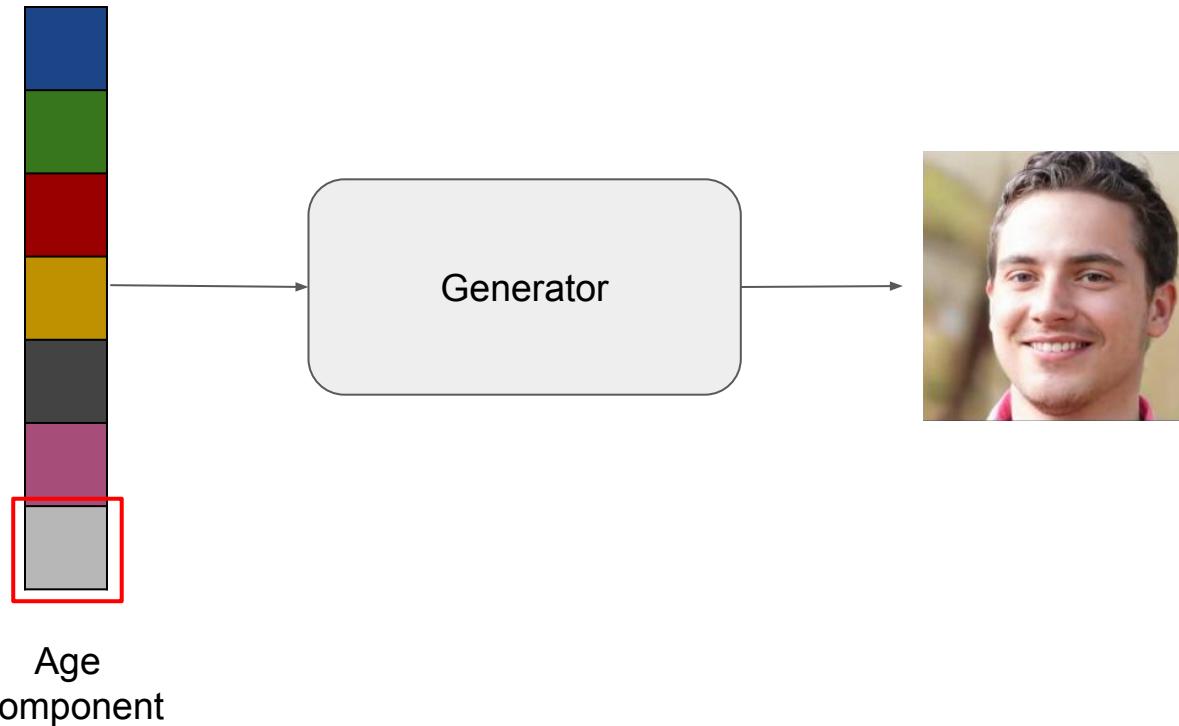


Image Transformation

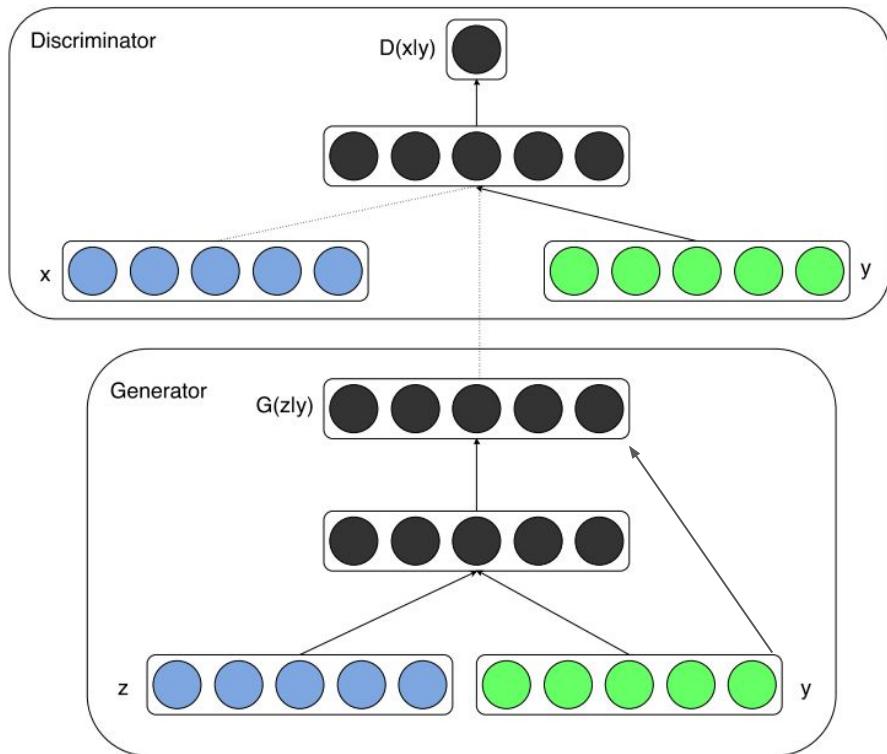


Conditional Image Generation

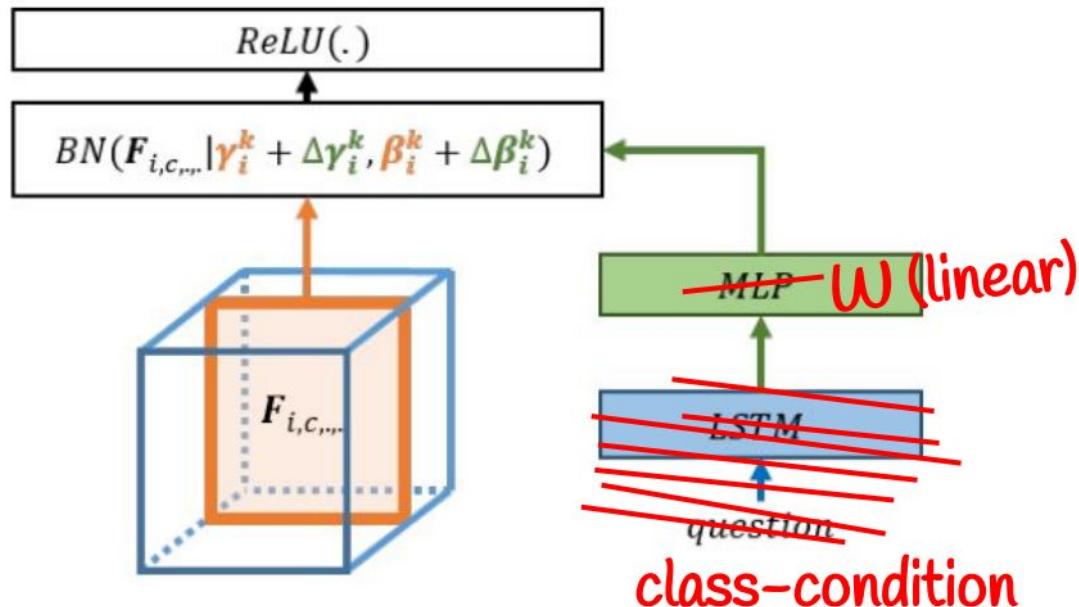
Conditional Generation

- In an unconditional model, there is no control over the characteristics of the data being generated
- A conditional model generates fake samples that satisfy some condition (a class label, base image, tags, text description, etc), rather than generating generic samples

Vanilla cGANs



BigGAN



- Among other things, use class-conditional batchnorm layers

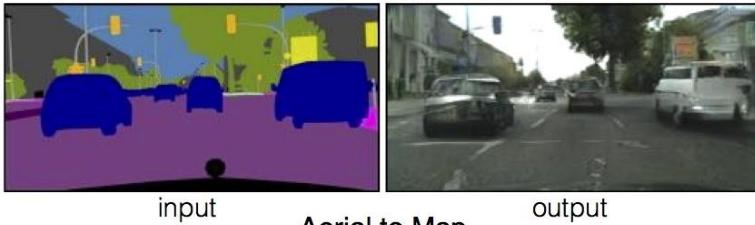
[Image Source](#)

BigGAN



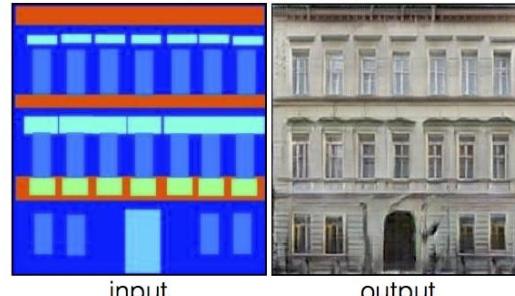
Paired Image-to-Image Translation

Labels to Street Scene



input

Labels to Facade



input

BW to Color



input

output

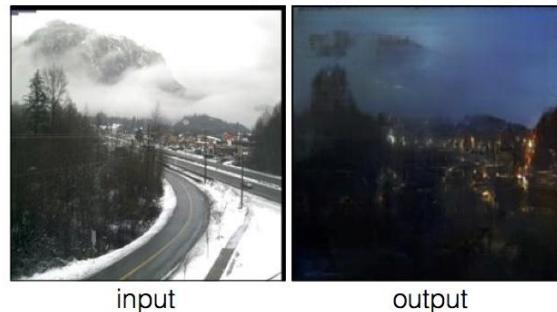
Aerial to Map



input

output

Day to Night



input

output

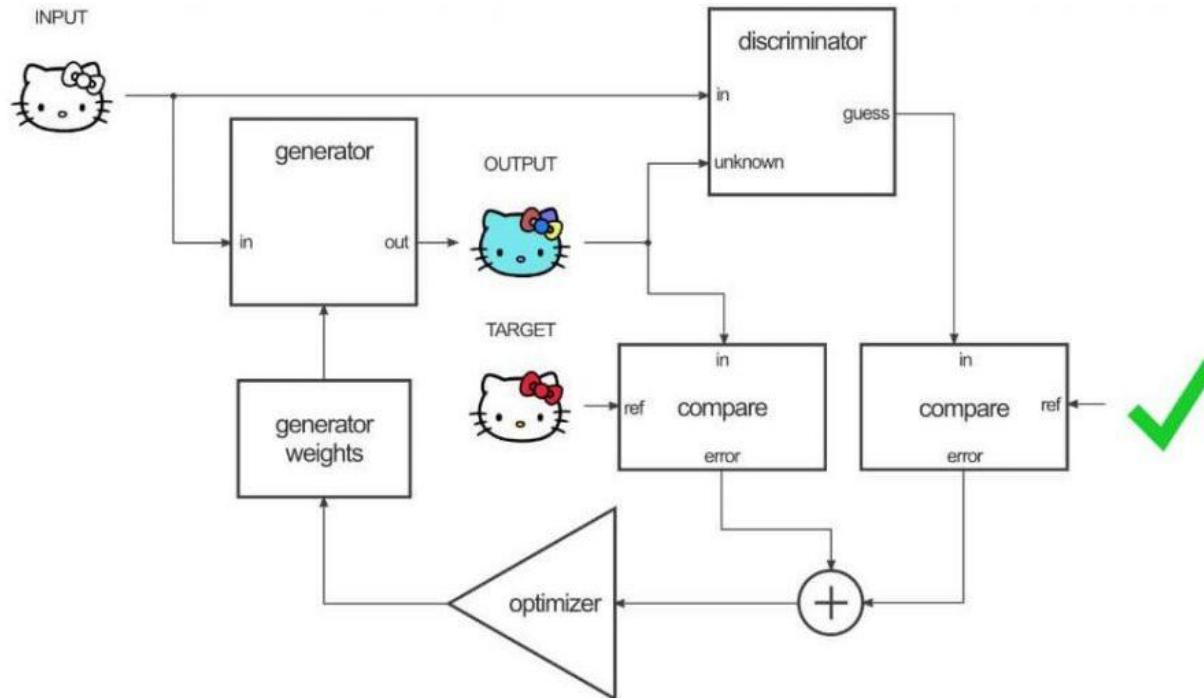
Edges to Photo



input

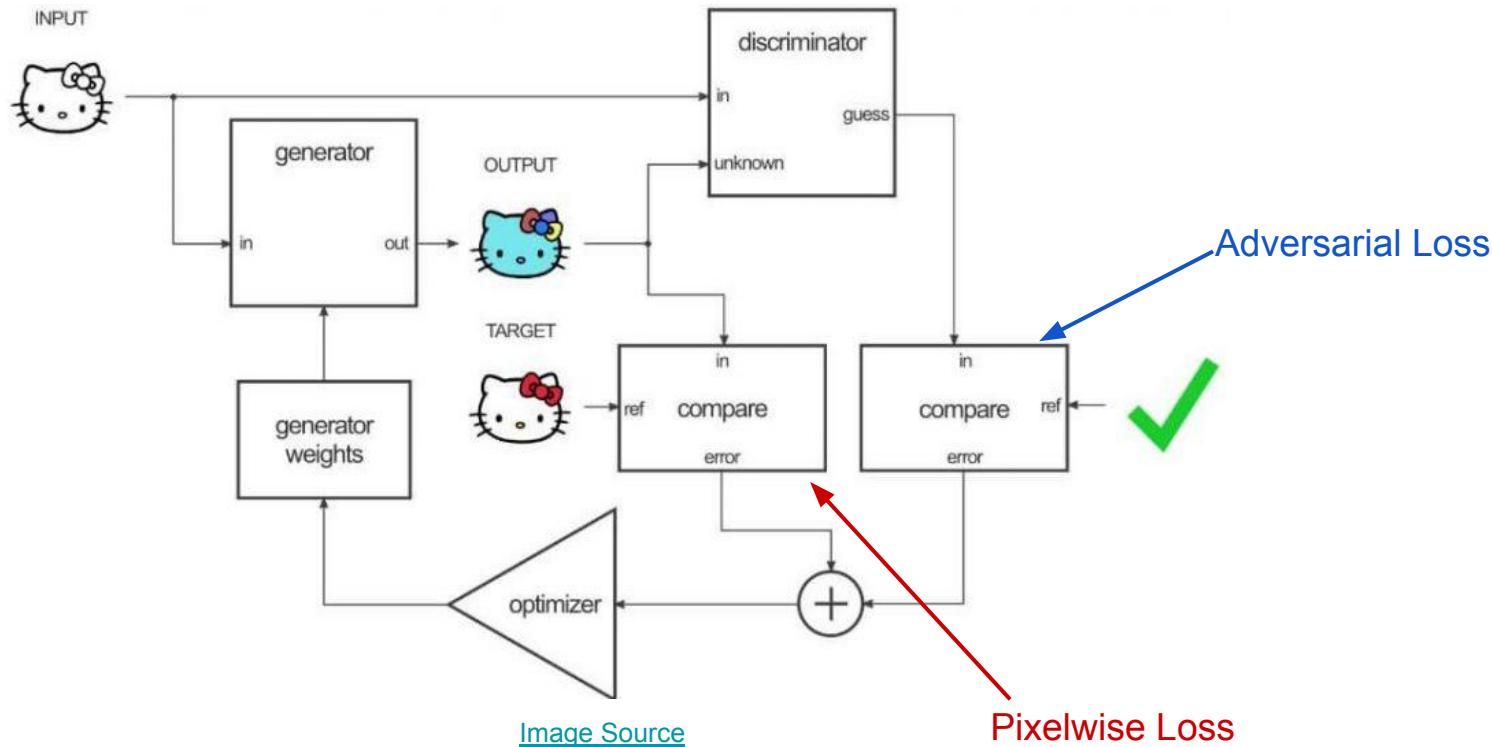
output

Paired Image-to-Image Translation



[Image Source](#)

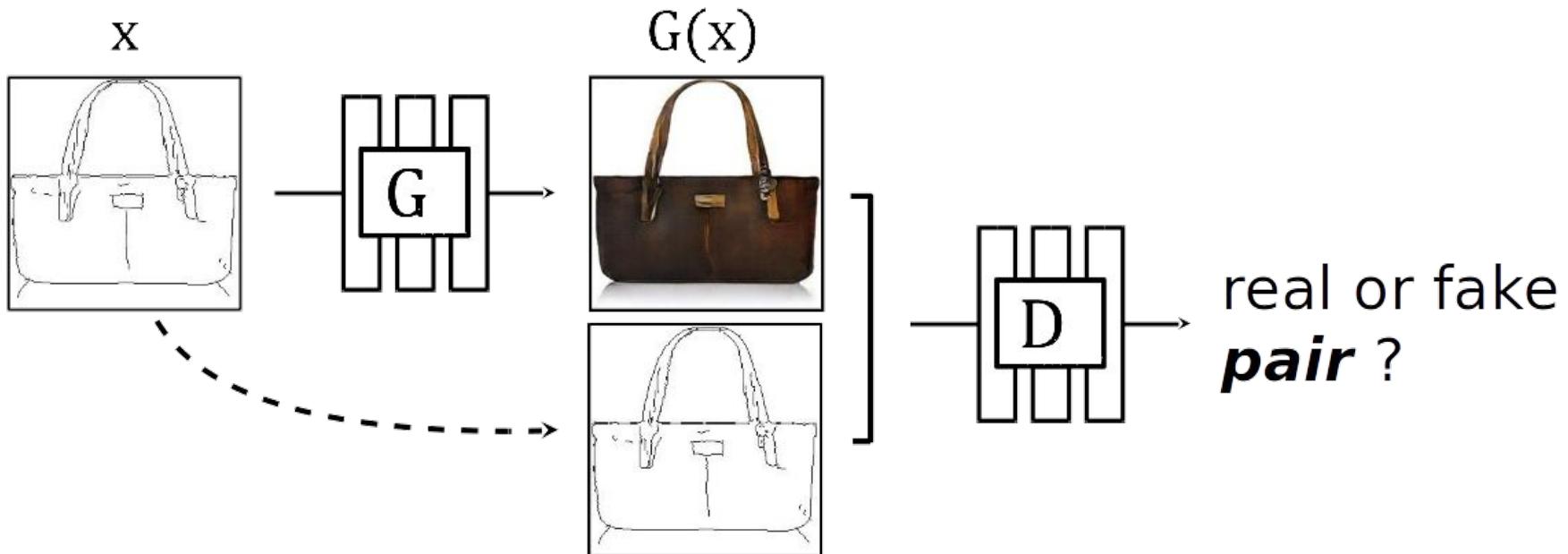
Paired Image-to-Image Translation



Why cGAN?



Why cGAN?



Effect of Different Losses



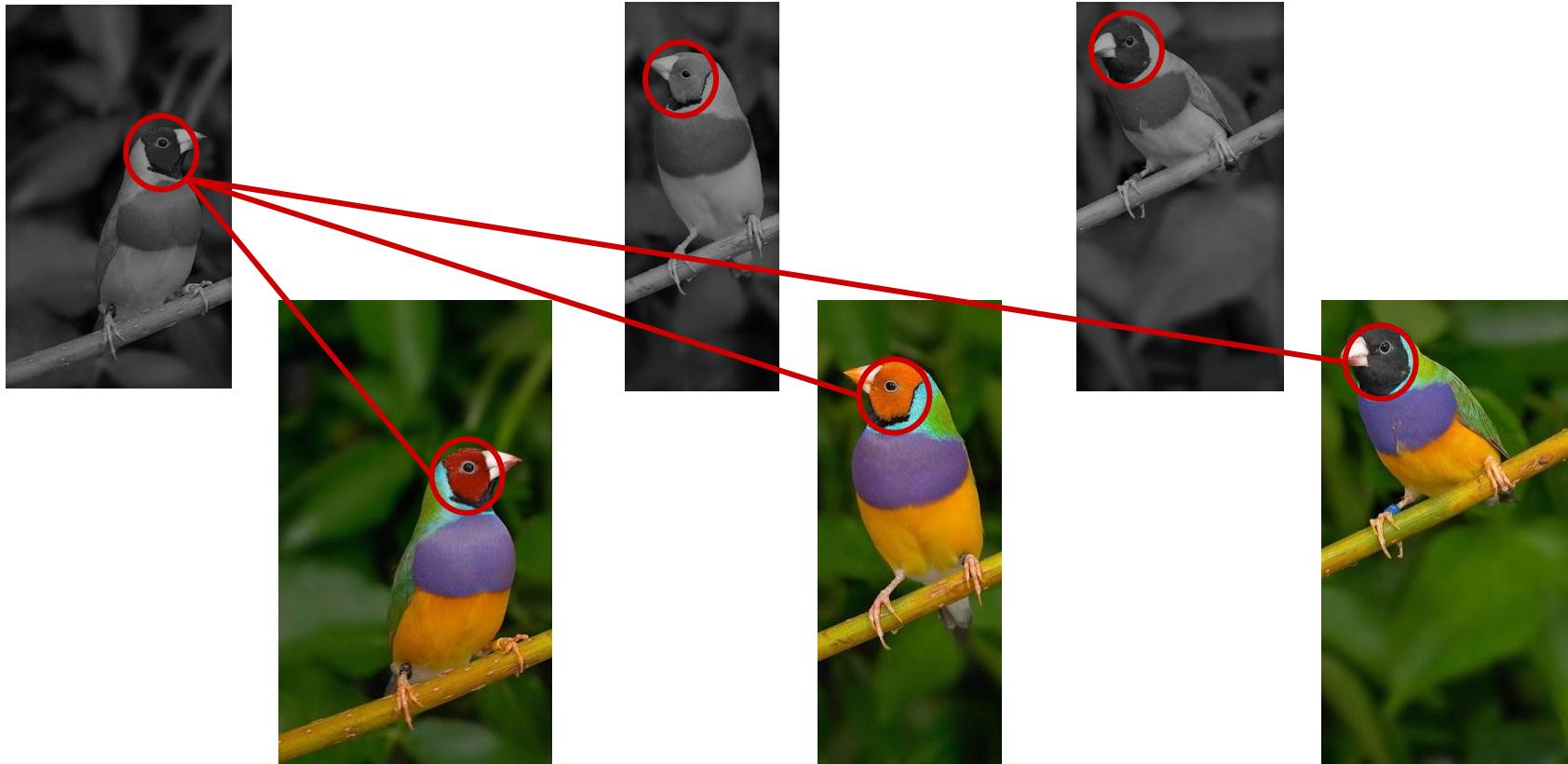
Notion of Adversarial Loss

Understanding Blurry Results Due to L1/L2



[Image Source](#)

Understanding Blurry Results Due to L1/L2



Understanding Blurry Results Due to L1/L2

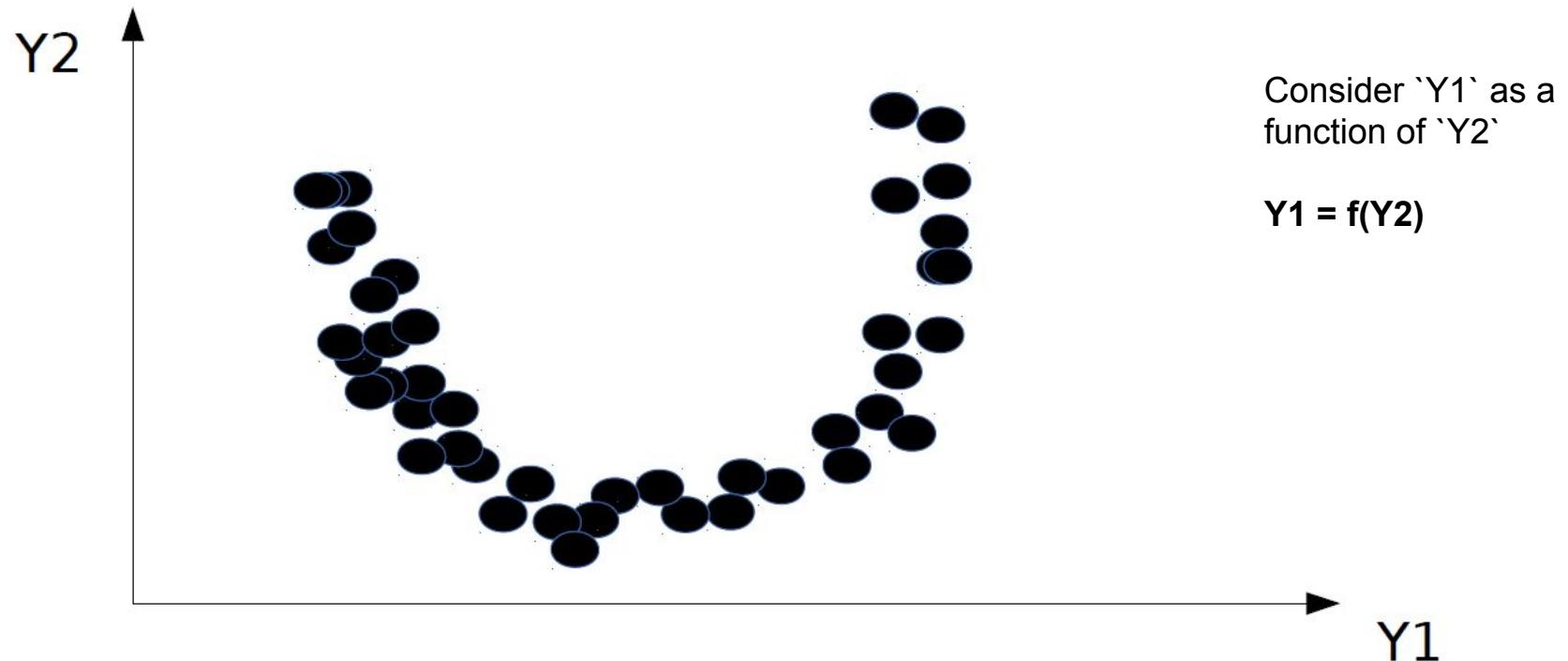
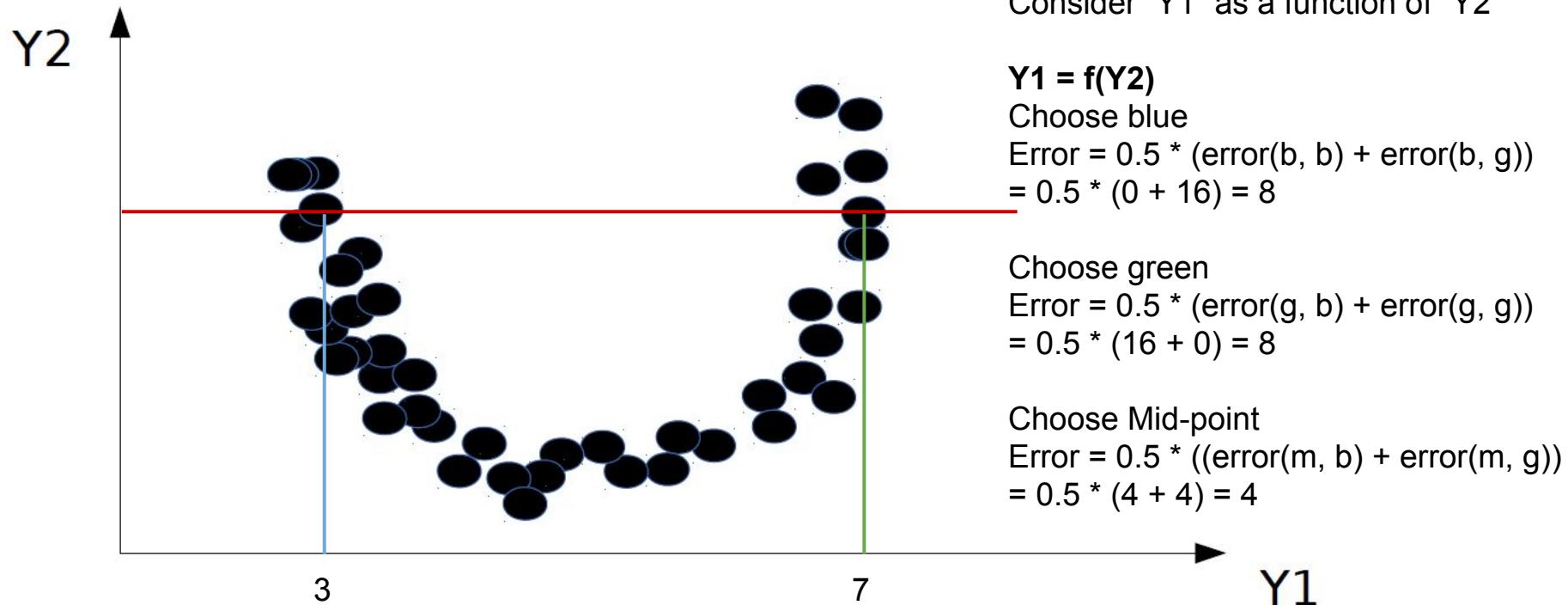
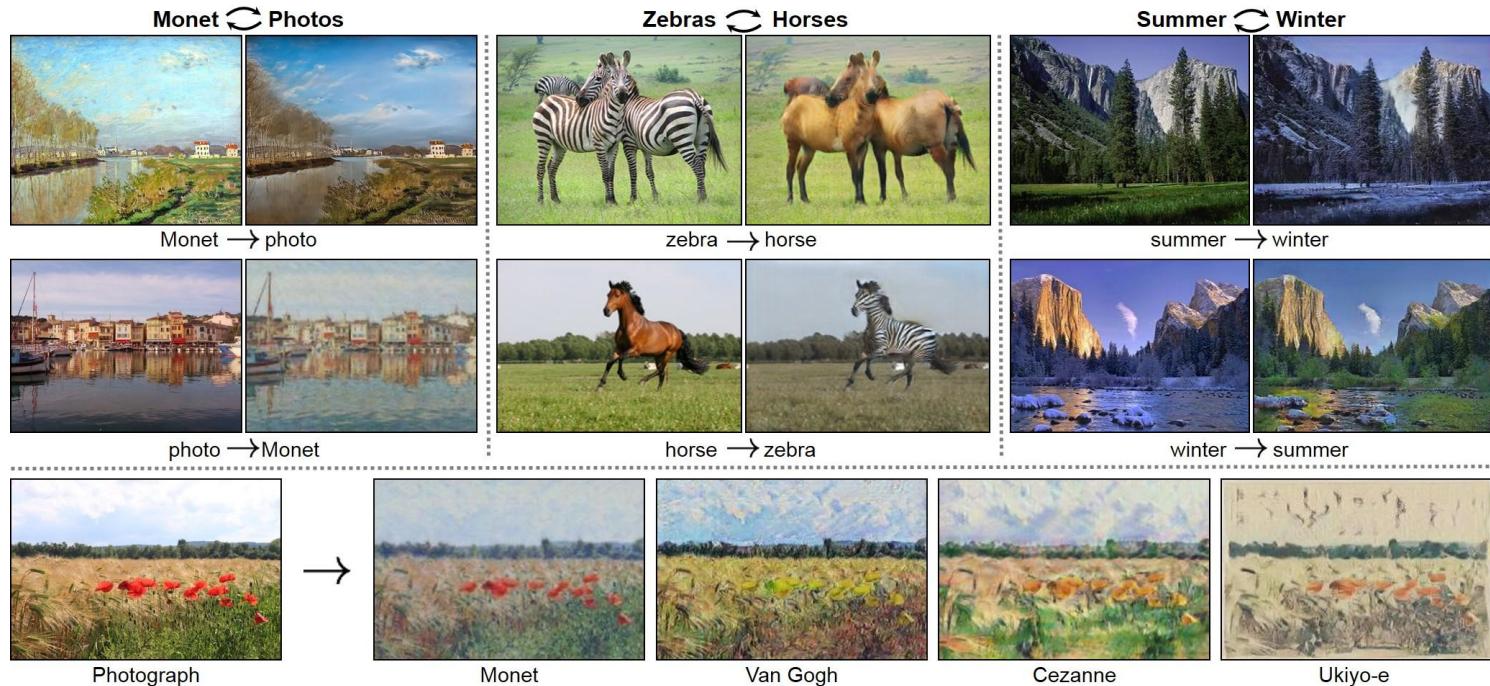


Image Credits: Prof Yann LeCun

Understanding Blurry Results Due to L1/L2

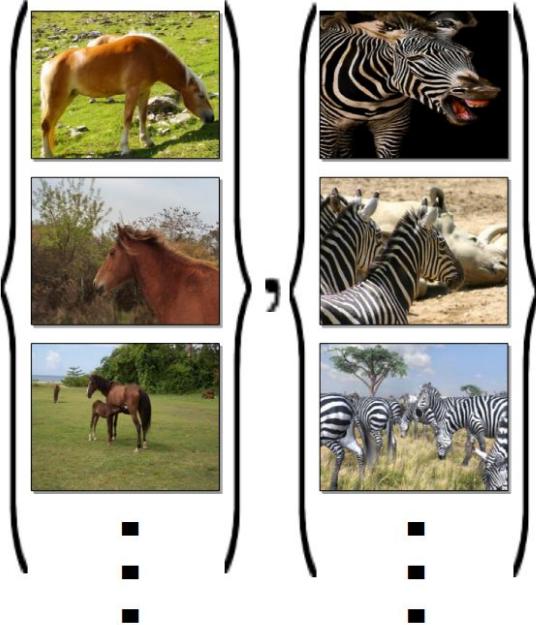


Unpaired Image-to-Image Translation



Snapshot of Data

Unpaired
 X Y



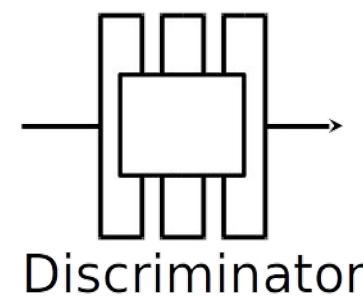
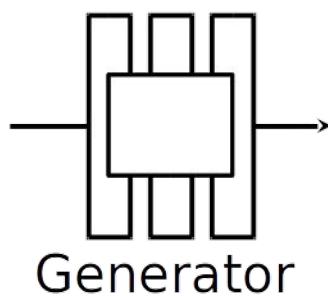
Unpaired
 X Y



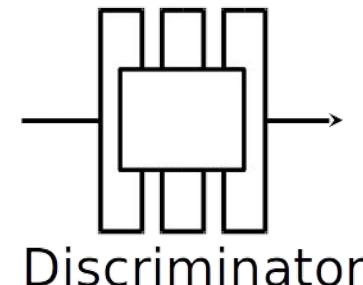
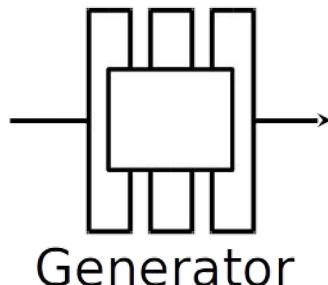
What About This Method?

- Train the generator to generate something from horses
- Sample a generated batch, and a batch of zebras
- Ask the discriminator to differentiate between them
- Train like a normal GAN

Problem



Real!



Real too!

Moreover



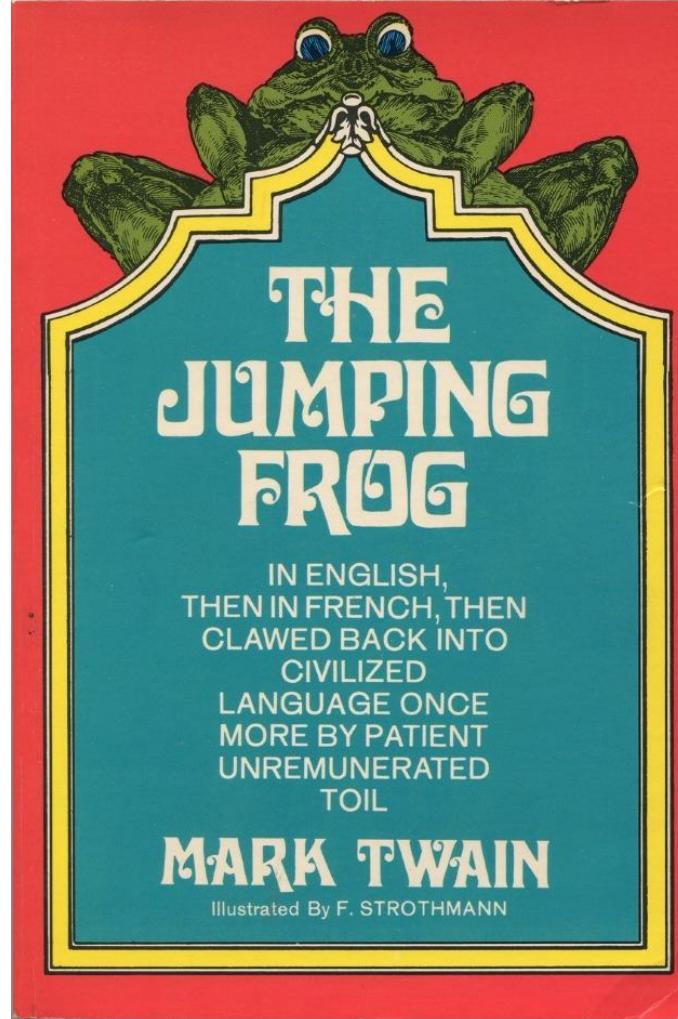
mode collapse!

Cycle Consistency

Original: "There was a feller here once by the name of Jim Smiley, in the winter of '49 or maybe it was the spring of '50 I don't recollect exactly"

Back Translation: "It there was one time here an individual known under the name of Jim Smiley; it was in the winter '49, possibly well at the spring of '50, I no me recollect not exactly."

— Mark Twain



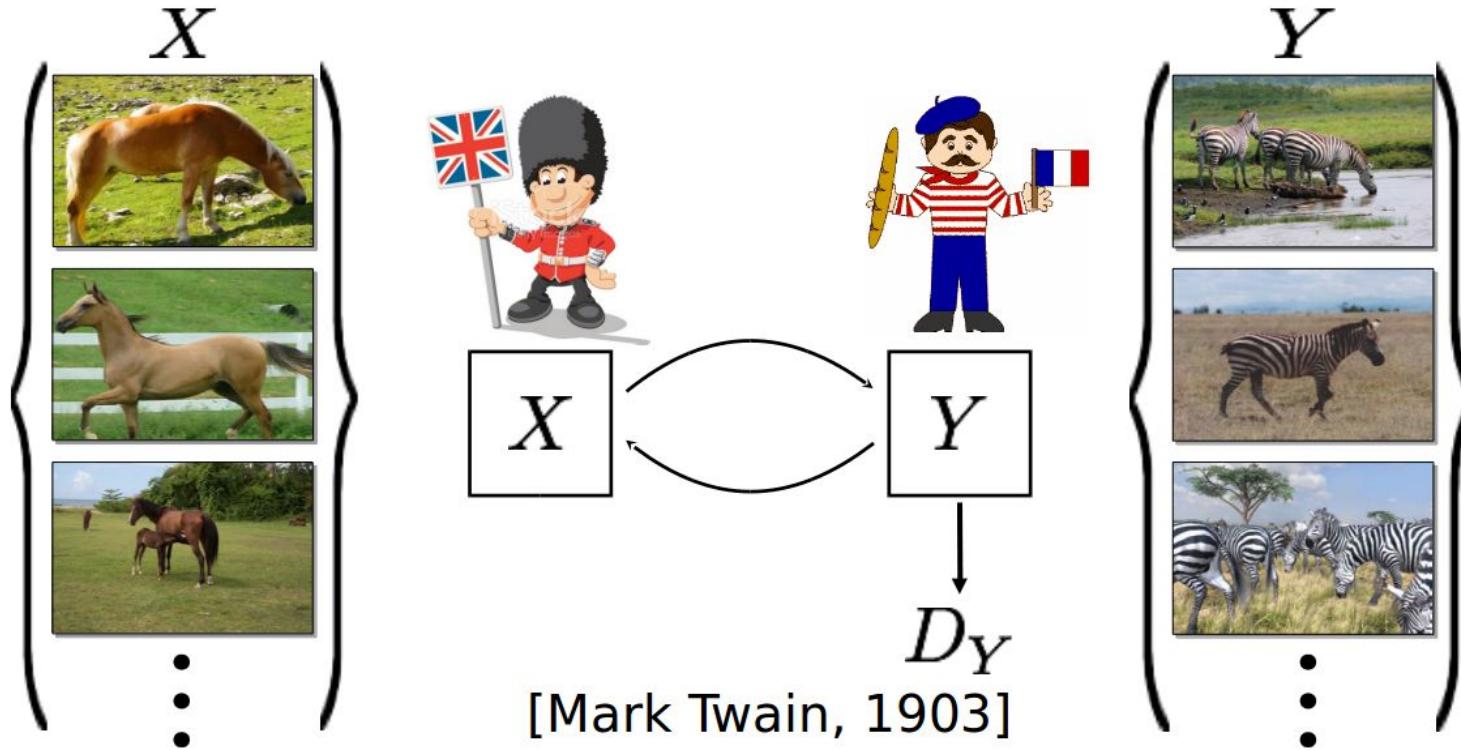
Cycle Consistency

- $G(F(x)) = x$, and
- $F(G(y)) = y$
- F and G are inverses of each other
- F and G are neural networks
- We will learn a function and its inverse together

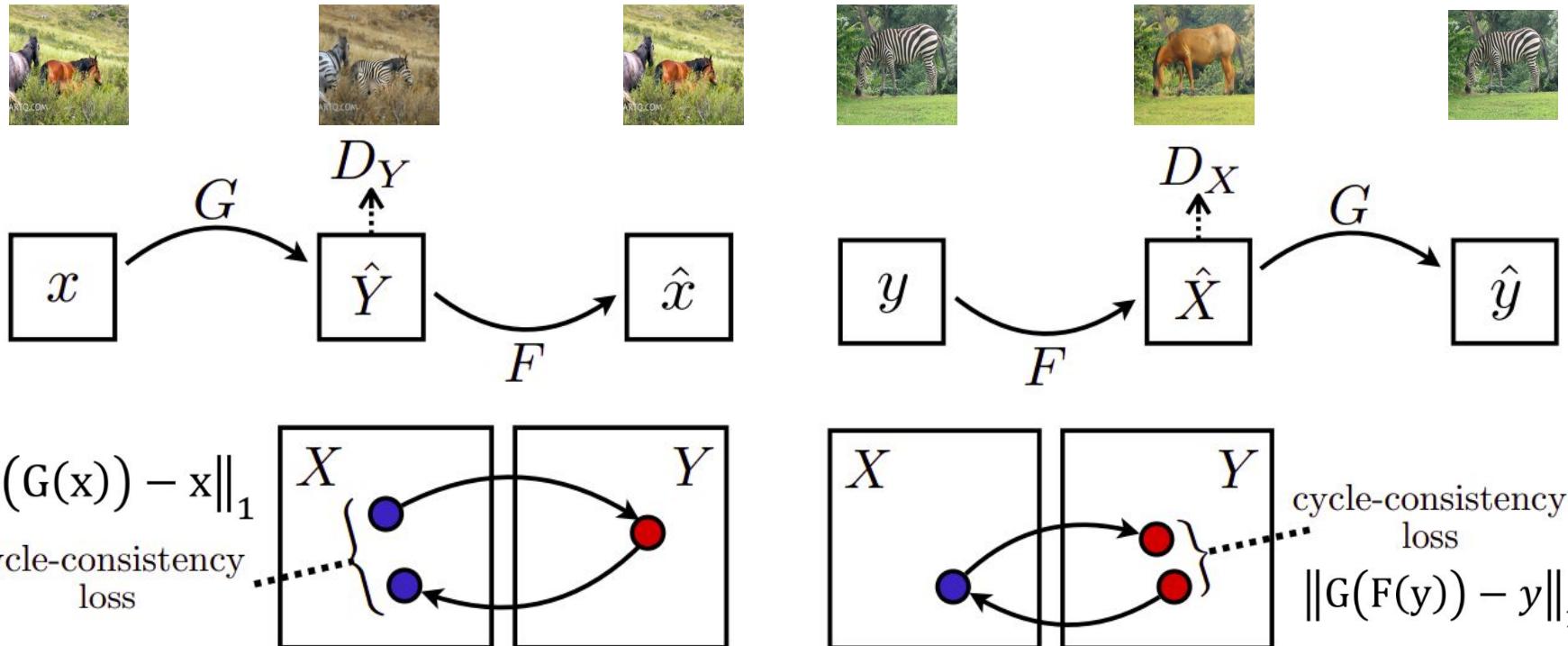


[Image Source](#)

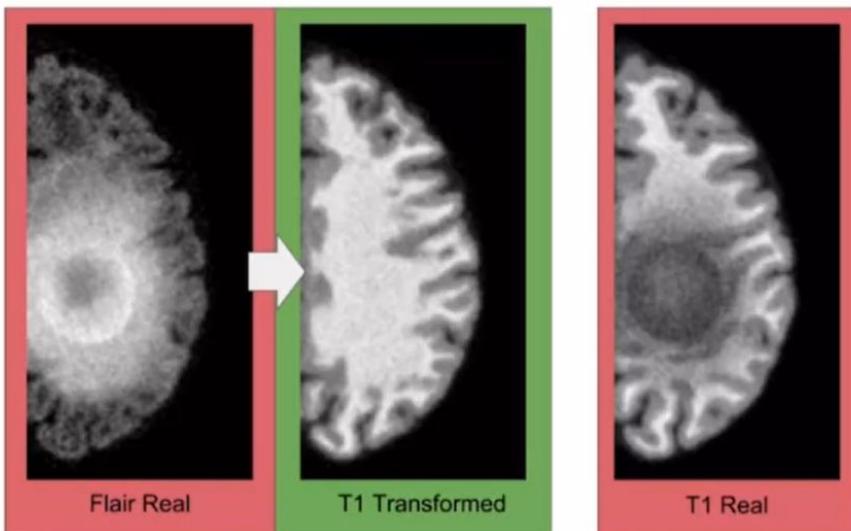
Cycle Consistency



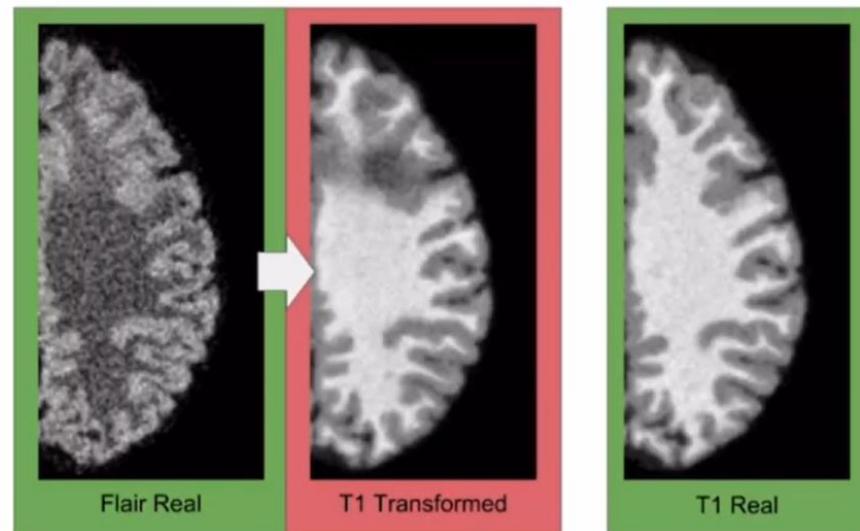
Cycle Consistency Loss



Applications



(a) A translation removing tumors



(b) A translation adding tumors

Applications

- Deepfakes: Cadbury Ad Campaign
- Superresolution (improving video game quality, MRIs)
- Image Restoration
- Background Matting
- Image Compression
- JPEG Artifact Removal
- Fashion Models Generation

Applications

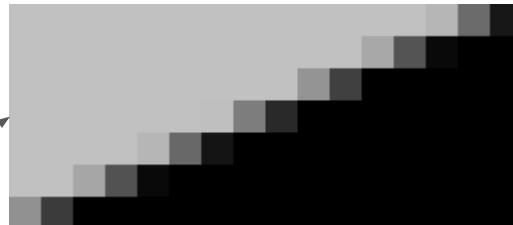
Super-resolution

Handcrafted Upsampling Algorithms



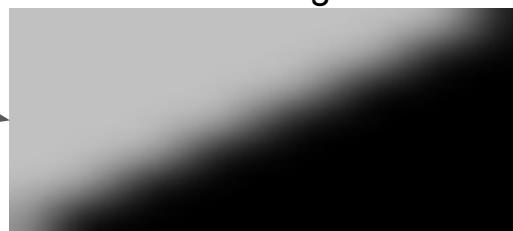
Test Image

Nearest
Neighbour



Aliasing

Bilinear with
Blur



Blurring

Bicubic



Edge Halos

[Source](#)

Super-resolution

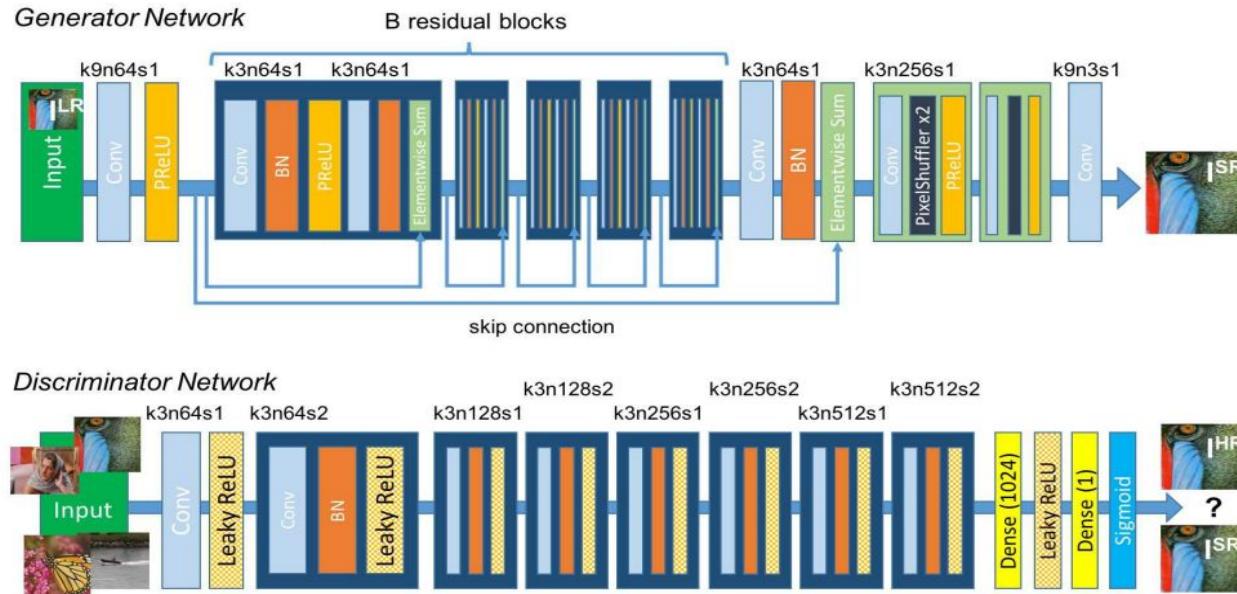
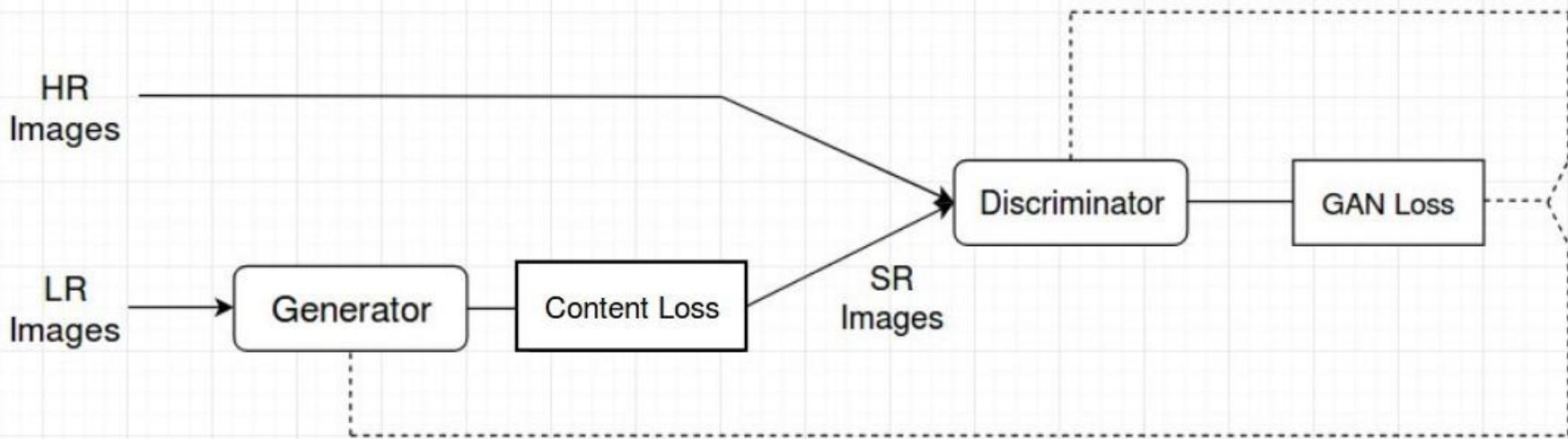


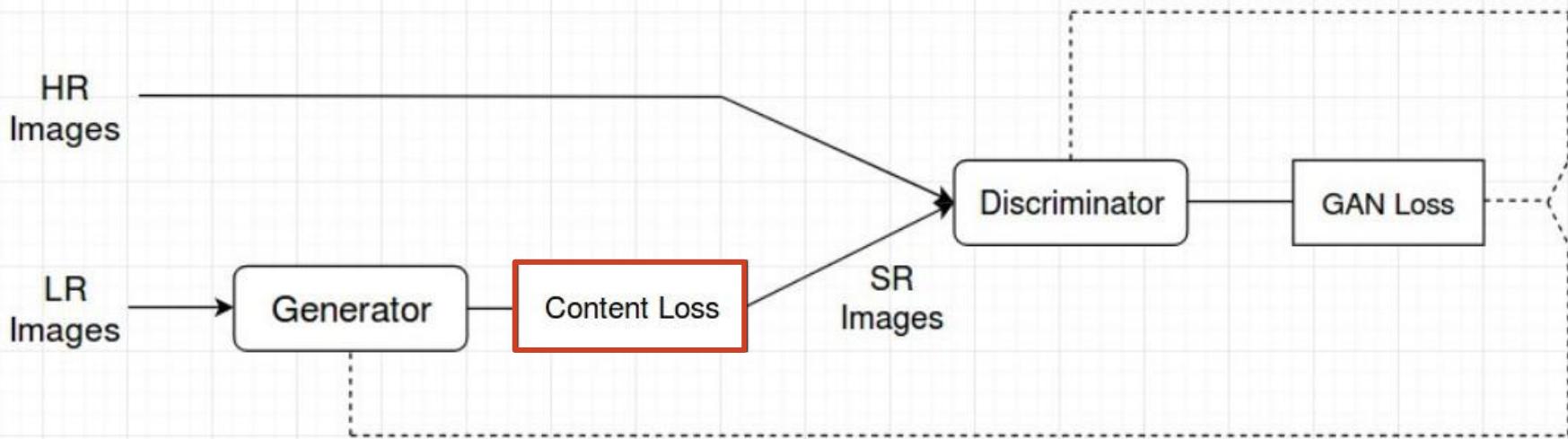
Figure 4: Architecture of Generator and Discriminator Network with corresponding kernel size (k), number of feature maps (n) and stride (s) indicated for each convolutional layer.

SRGAN Architecture



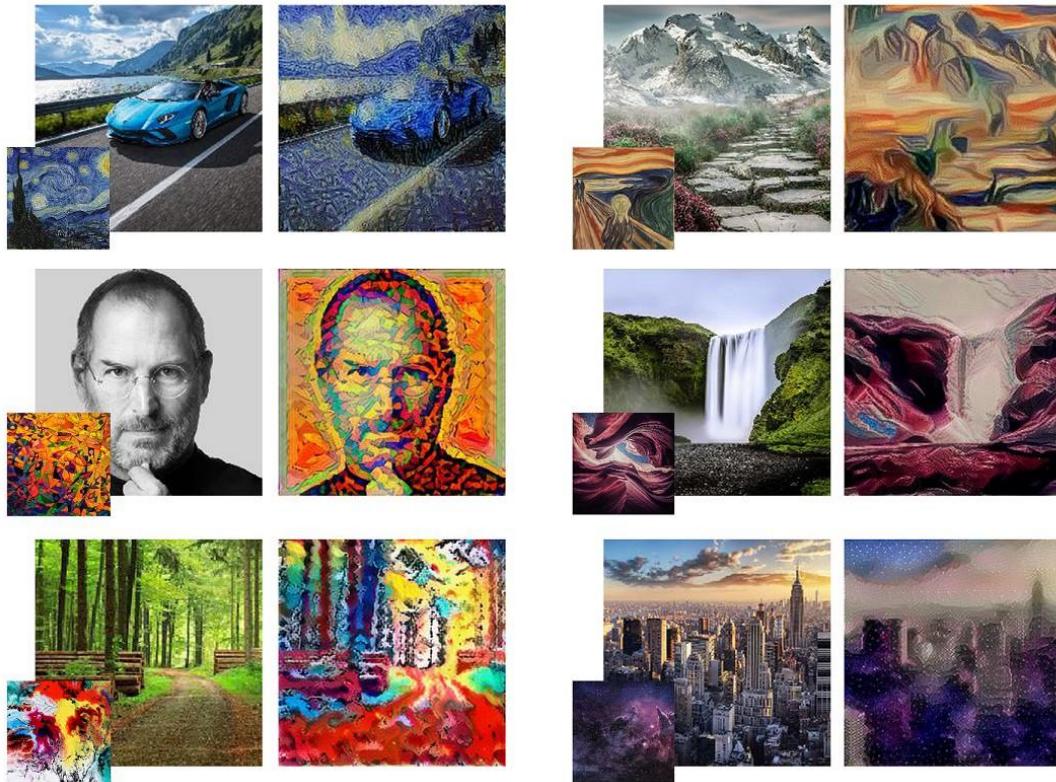
[Image Source](#)

SRGAN Architecture



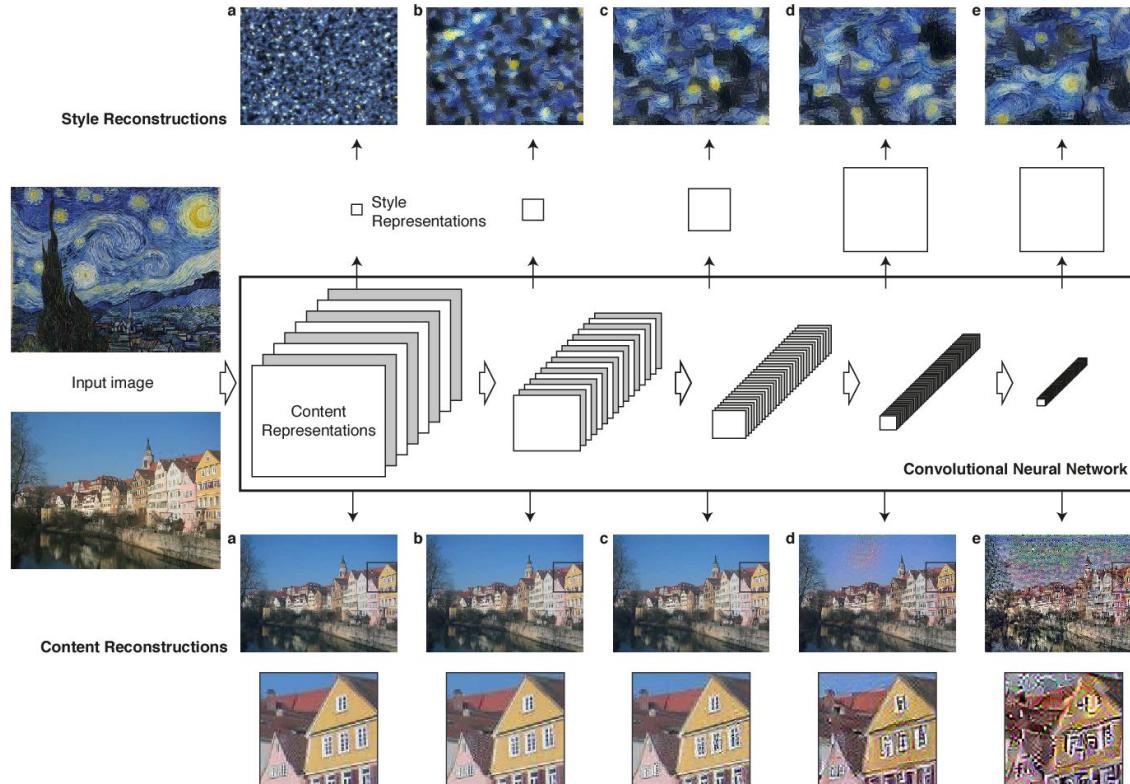
[Image Source](#)

Style Transfer



[Image Source](#)

Deep Feature Spaces



Gatys, L.A., Ecker, A.S., & Bethge, M. (2016). Image Style Transfer Using Convolutional Neural Networks. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2414-2423.

Deep Features as Perceptual Metric



Content Loss

$$\mathcal{L}_{content}(\vec{p}, \vec{x}, l) = \frac{1}{2} \sum_{i,j} (F_{ij}^l - P_{ij}^l)^2$$

\vec{p} The original image

\vec{x} The generated image

l Layer

F_{ij}^l Activation of the i^{th} filter at position j in the feature representation of \vec{x} in l

P_{ij}^l Activation of the i^{th} filter at position j in the feature representation of \vec{p} in l

SRGAN

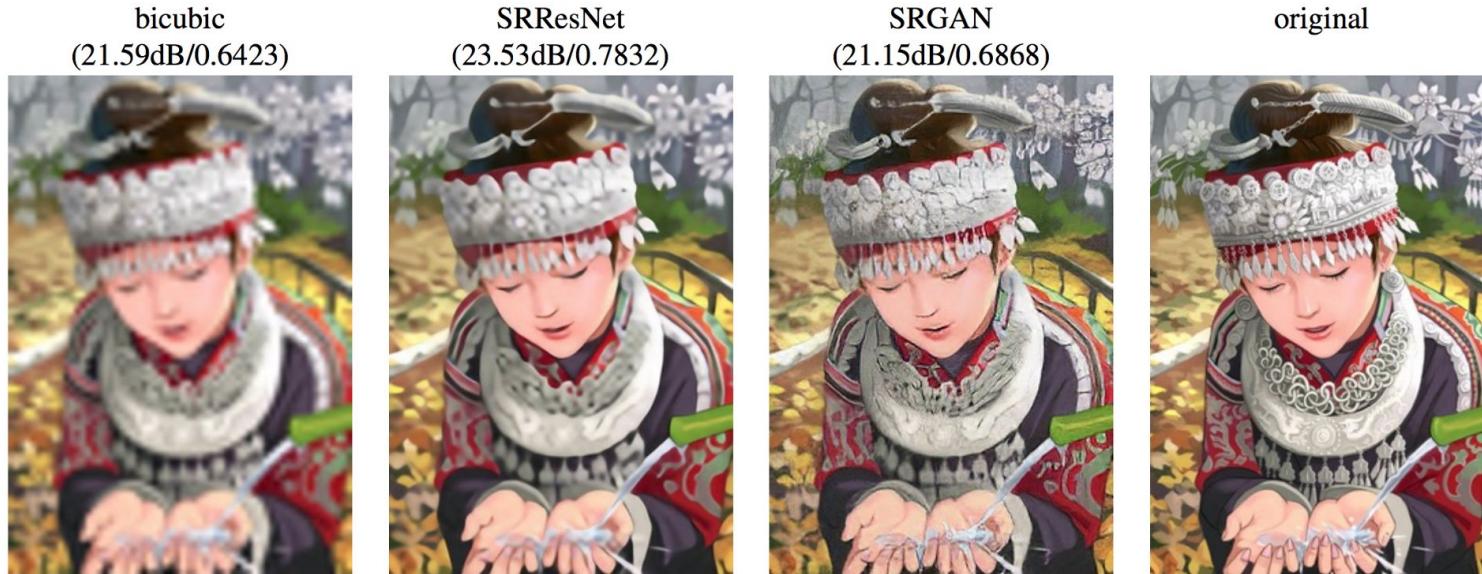
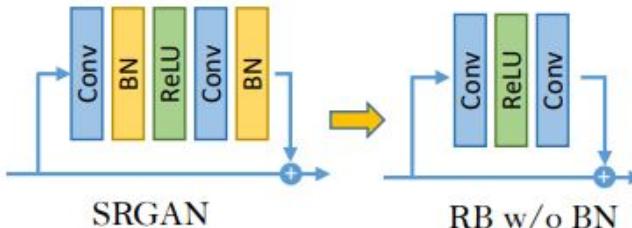


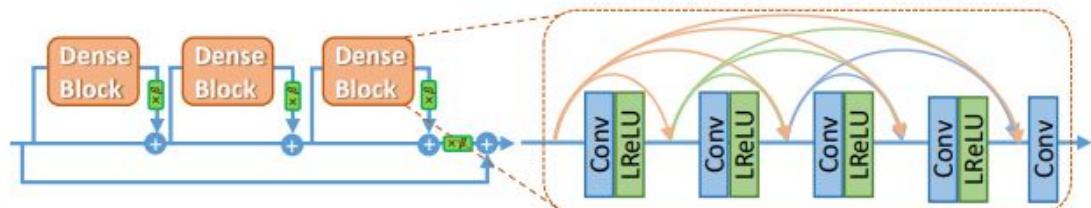
Figure 2: From left to right: bicubic interpolation, deep residual network optimized for MSE, deep residual generative adversarial network optimized for a loss more sensitive to human perception, original HR image. Corresponding PSNR and SSIM are shown in brackets. [4× upscaling]

ESRGAN

Residual Block (RB)



Residual in Residual Dense Block (RRDB)



$$D(x_r) = \sigma(C(\text{Real})) \rightarrow 1 \quad \text{Real?}$$

$$D(x_f) = \sigma(C(\text{Fake})) \rightarrow 0 \quad \text{Fake?}$$

a) Standard GAN

$$D_{Ra}(x_r, x_f) = \sigma(C(\text{Real}) - \mathbb{E}[C(\text{Fake})]) \rightarrow 1$$

More realistic than fake data?

$$D_{Ra}(x_f, x_r) = \sigma(C(\text{Fake}) - \mathbb{E}[C(\text{Real})]) \rightarrow 0$$

Less realistic than real data?

b) Relativistic GAN

Remastering Old Video Games



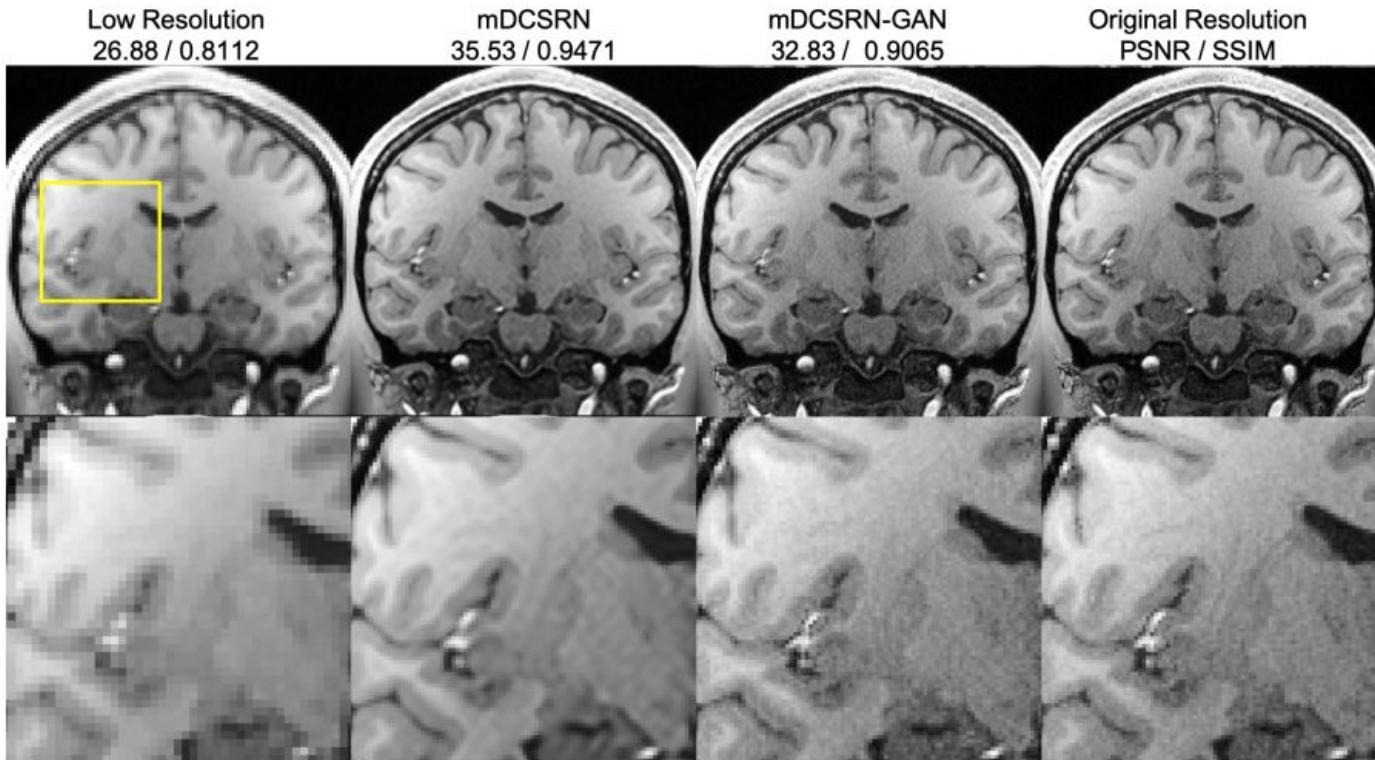
[Source](#)

Remastering Old Video Games



[Source](#)

Improving MRI Quality



Chen, Y., Christodoulou, A.G., Zhou, Z., Shi, F., Xie, Y., & Li, D. (2020). MRI Super-Resolution with GAN and 3D Multi-Level DenseNet: Smaller, Faster, and Better. *ArXiv, abs/2003.01217*.

Making Satellite Images Clearer



[Source](#)

Detour: Content Loss Application

Search Visually Similar Products

Home / Dresses / KAMARIA SHORT SLEEVE BUTTON UP DRESS



ID: #002ed317f65bc3760e9381f0

Kamaria Short Sleeve Button

\$280.00

Culture Phit Size Chart: Bring some fun to your warm weather! Flattering stripes on stretch jersey knit fabrication. Round neck sleeves. Straight hemline. 95% rayon; 5% spandex. Machine washable.

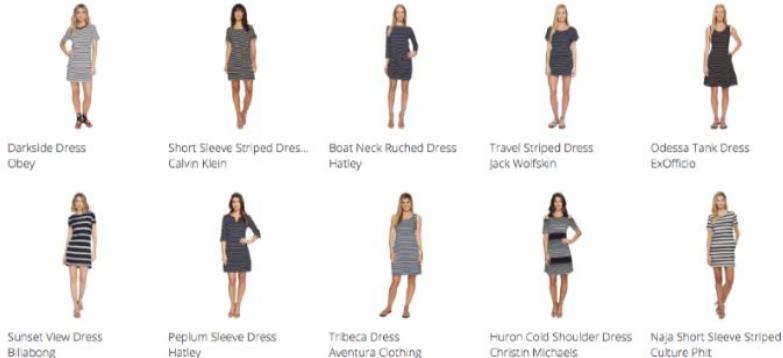
Measurements: Length: 34 in Product measurements were taken when garment was laid flat. Measurements may vary by size. Length: 34 in

Similarity visual

[ORDER NOW](#)

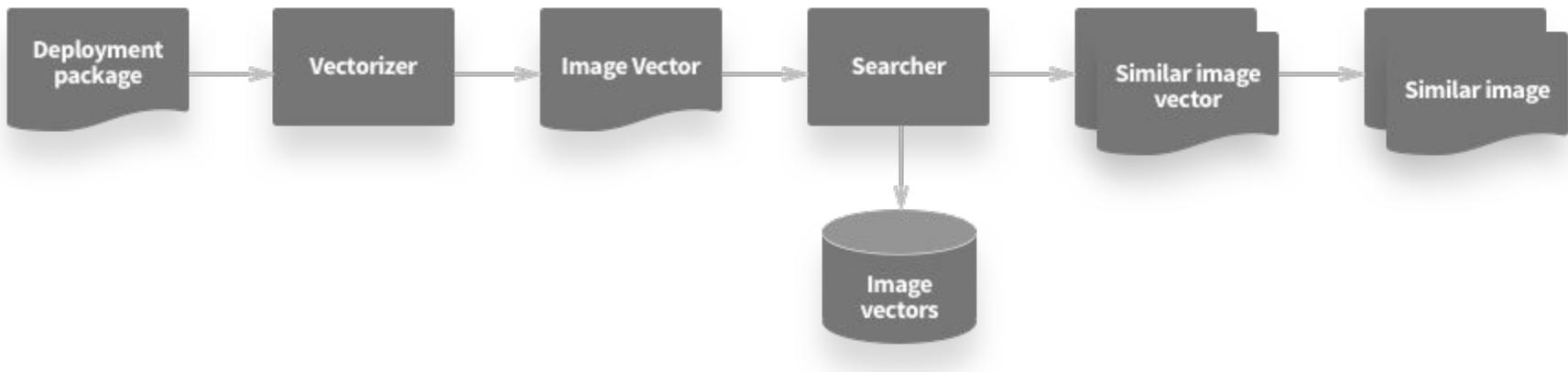
Brand	Culture Phit
Gender	Women
Color category	Blue
Dress length	short
Dress type	shift
Material	spandex
Neckline style	v-neck
Occasion	casual

WE FOUND FOR YOU 10 SIMILAR PRODUCTS

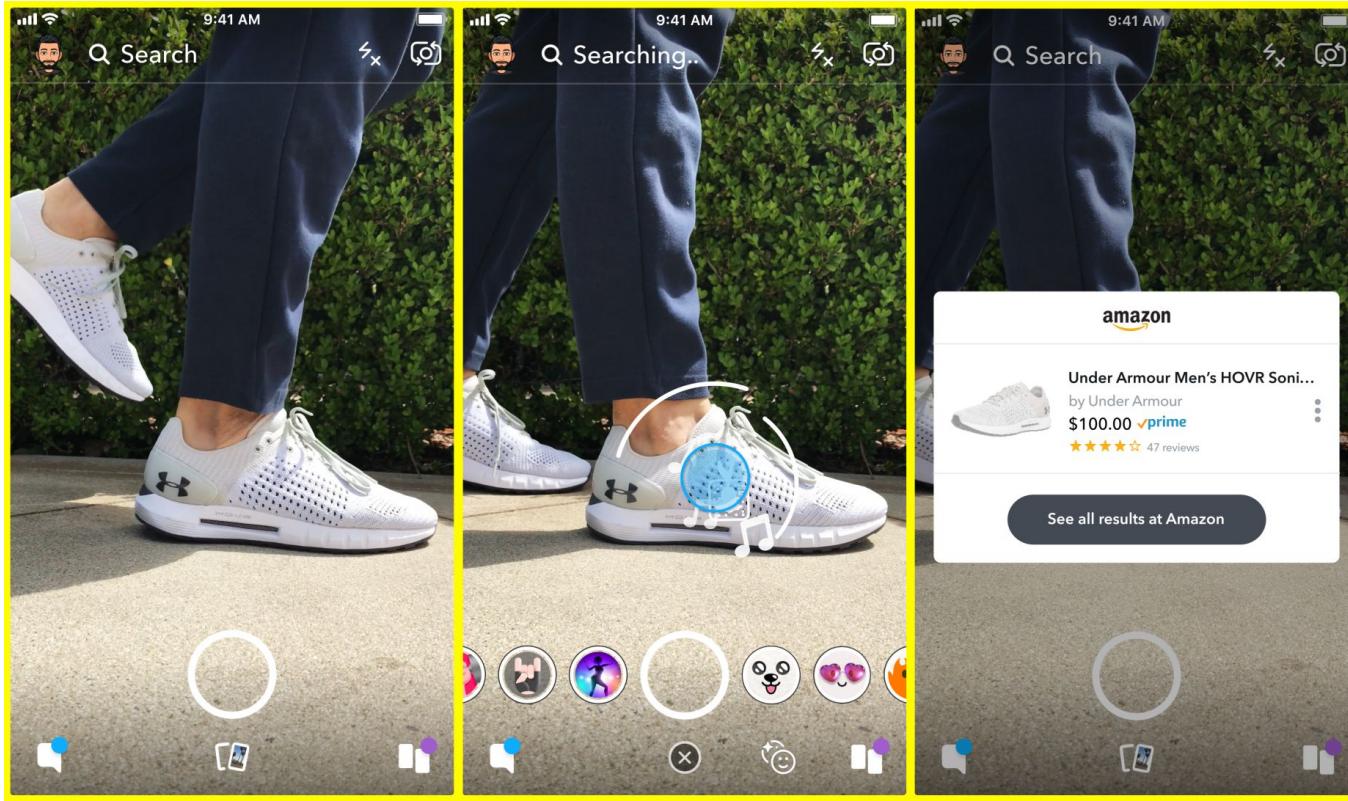


[Source](#)

Under the Hood



Product Search via Visual Query



[Source](#)

Image Restoration

Bringing Old Photos Back to Life



Standard Level



\$25.00 per photo

[View Examples](#)

Services applied:

- Light Scratches Removal
- Sepia Issue Fixing
- Dull Colors Editing
- Color Correction
- Small Damaged Areas Restoration
- Blemishes Removal
- Overall Enhancement
- Up to Three Persons on the Image

[Order Now](#)

Premium Level



\$35.00 per photo

[View Examples](#)

Services applied:

- Standard Restoration
- Medium Scratches Retouching
- Missing Parts Fixing
- Complete Restoration of an Every Detail
- Up to Three Persons on the Image

[Order Now](#)

Extreme Level



\$50.00 per photo

[View Examples](#)

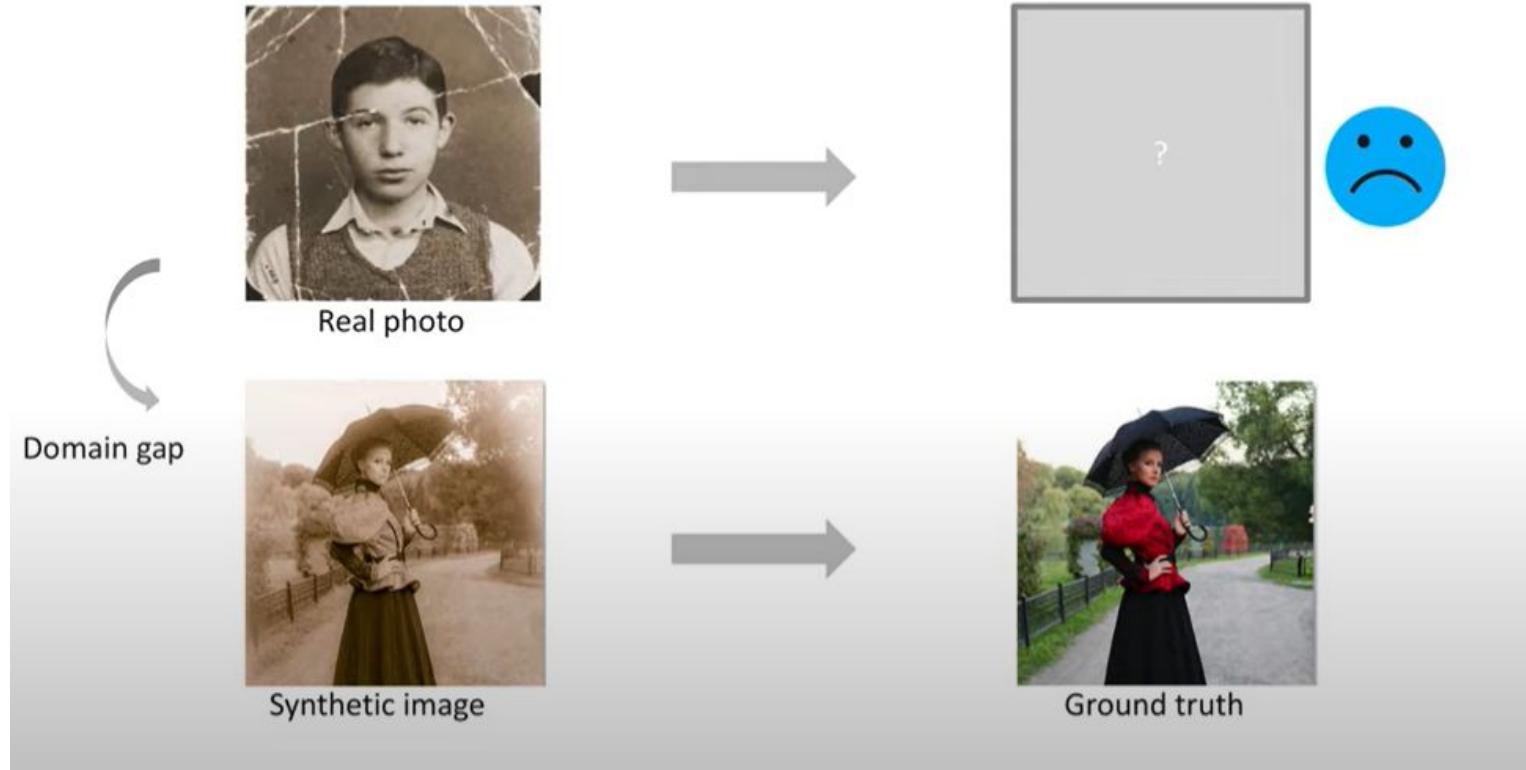
Services applied:

- Premium Restoration
- Deeply Damaged Photos Restoration
- Ruined Edges Fixing
- Photos Merging
- Image Coloration
- Up to Three Persons on the Image

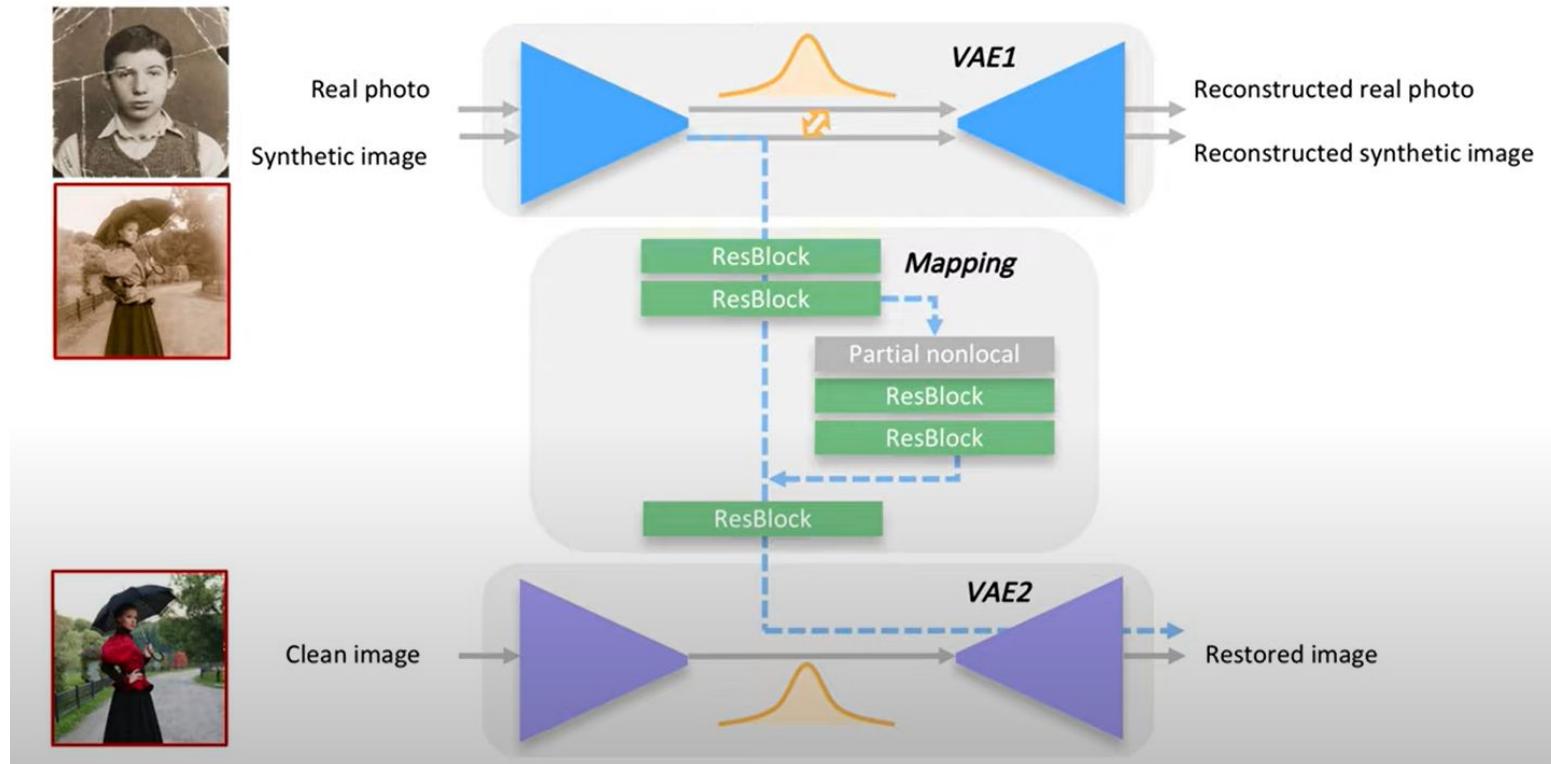
[Order Now](#)

[Source](#)

Using Synthetic Data



Solution Via Latent Space Mapping



Results



Wan, Z., Zhang, B., Chen, D., Zhang, P., Chen, D., Liao, J., & Wen, F. (2020). Bringing Old Photos Back to Life. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2744-2754.

Frame Interpolation

How are 2D Animations Made



[Source](#)

Costs

3.8K
Per Minute

- Our entry level animation
- Color elements



5.5K
Per Minute

- Full color
- Some character & accent animation



8K
Per Minute

MOST POPULAR

- Motion design
- Art director
- Script writer



10K
Per Minute

- Character animation
- Premium script writer
- Advanced camera transitions



Source

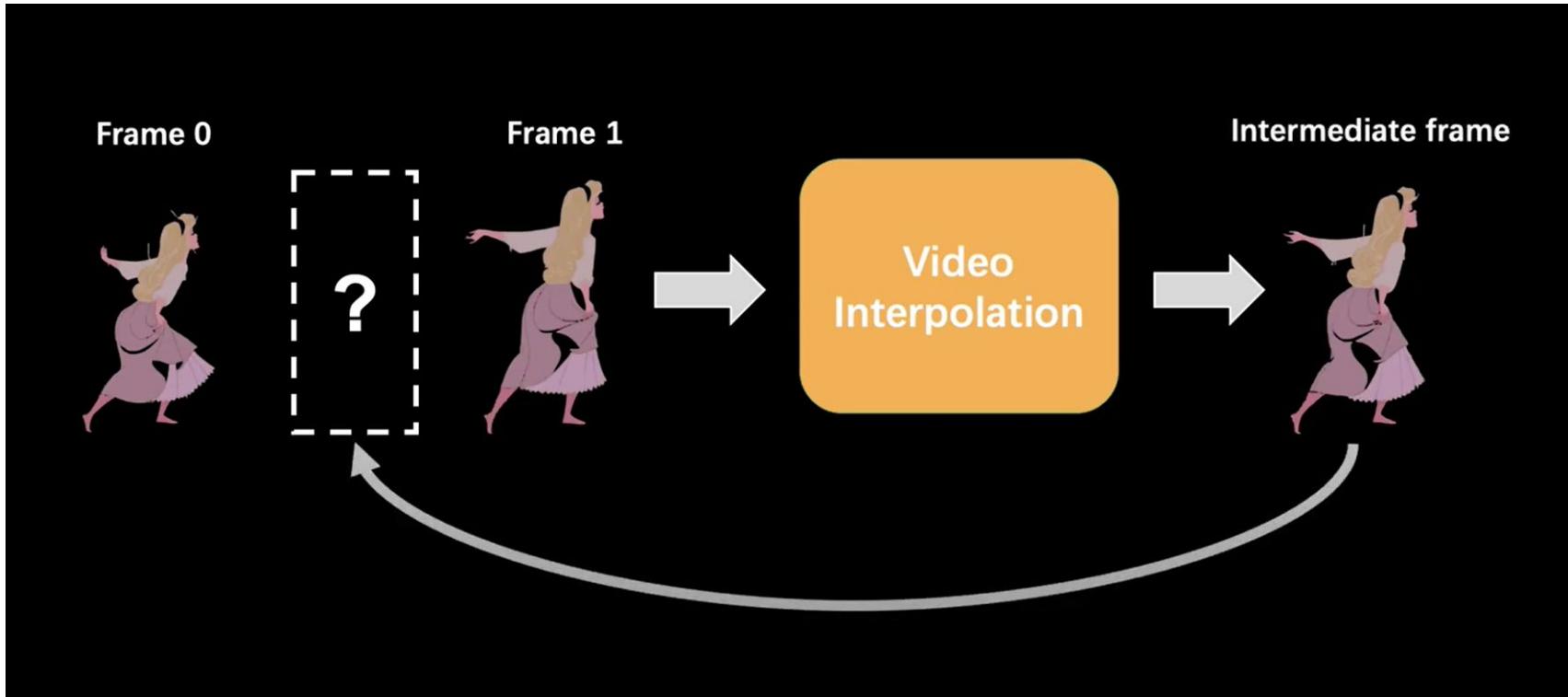
Saving Cost and Time



“on twos” 24 fps → 12 fps

“on threes” 24 fps → 8 fps

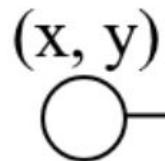
Frame Interpolation



[Source](#)

Optical Flow

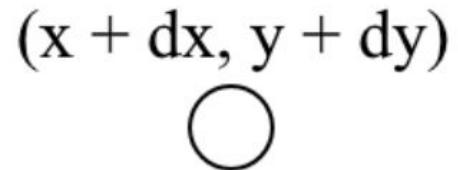
$I(x, y, t)$



displacement = (dx, dy)

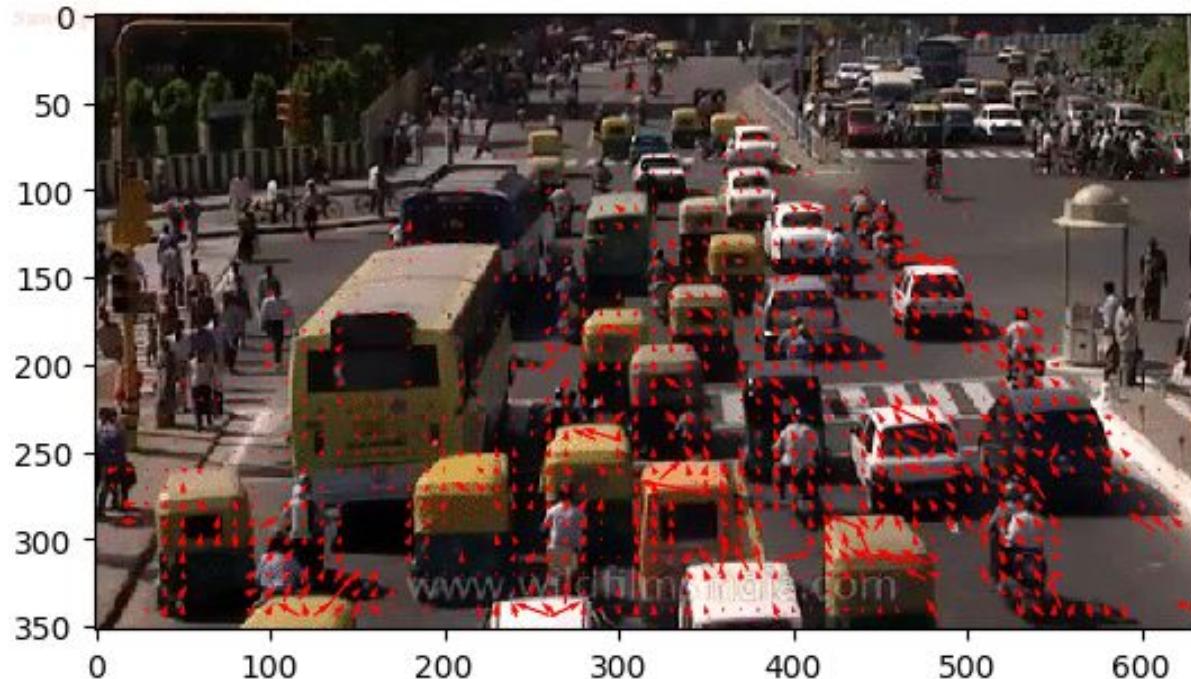
time = t

$I(x + dx, y + dy, t + dt)$



time = $t + dt$

Optical Flow



[Source](#)

Anime Interp

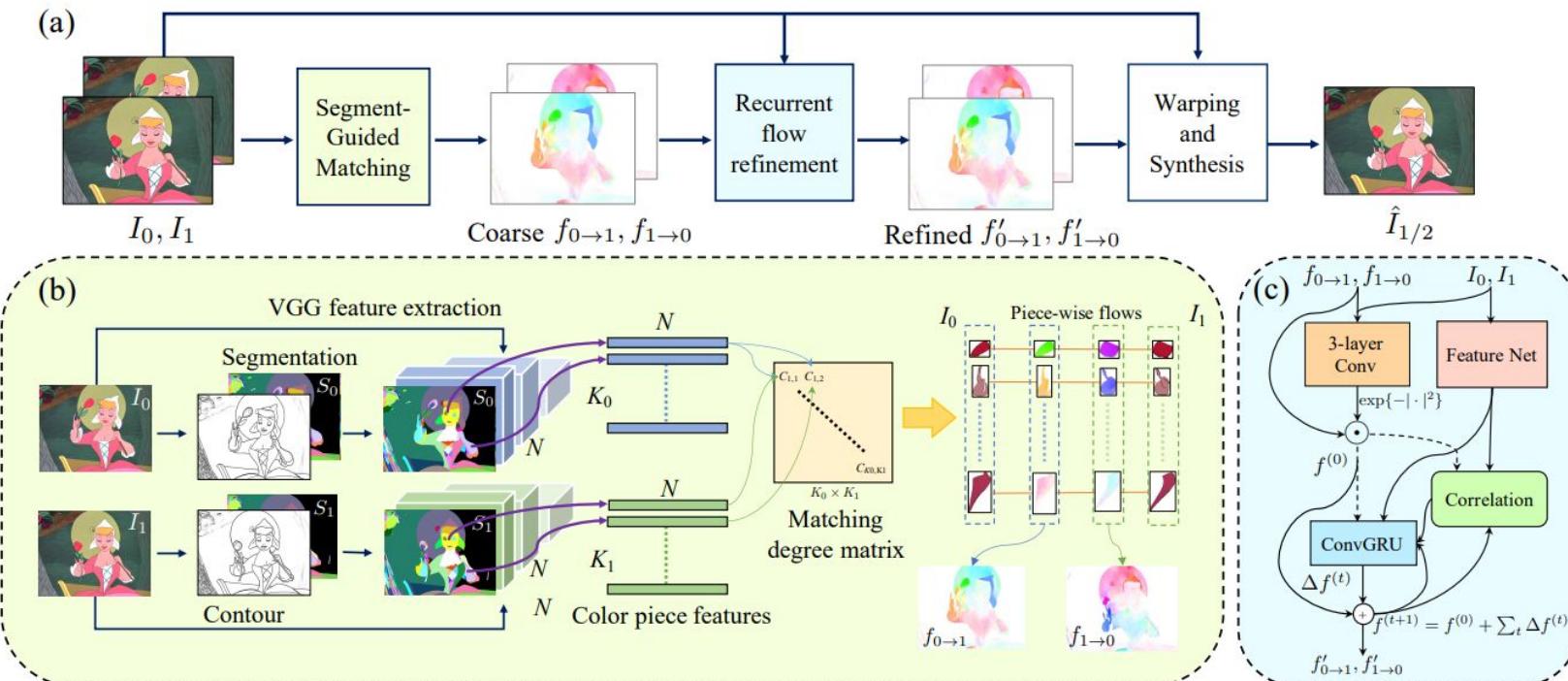


Image Matting

Image Matting

- Separating an image into foreground and background
- Matting equation

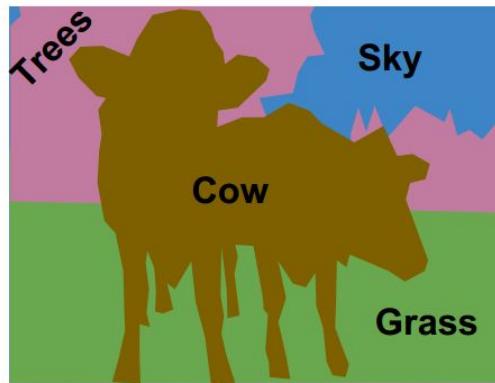
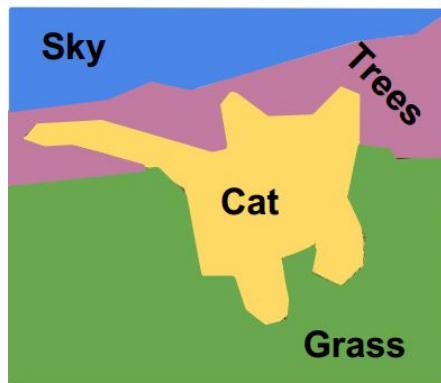
$$C = F * \alpha + B * (1 - \alpha)$$

- C is the scene, F are the foreground pixels, B are background pixels and alpha is transparency

Semantic Segmentation



[This image is CC0 public domain](#)



[Image Source](#)

Problem With Segmentation



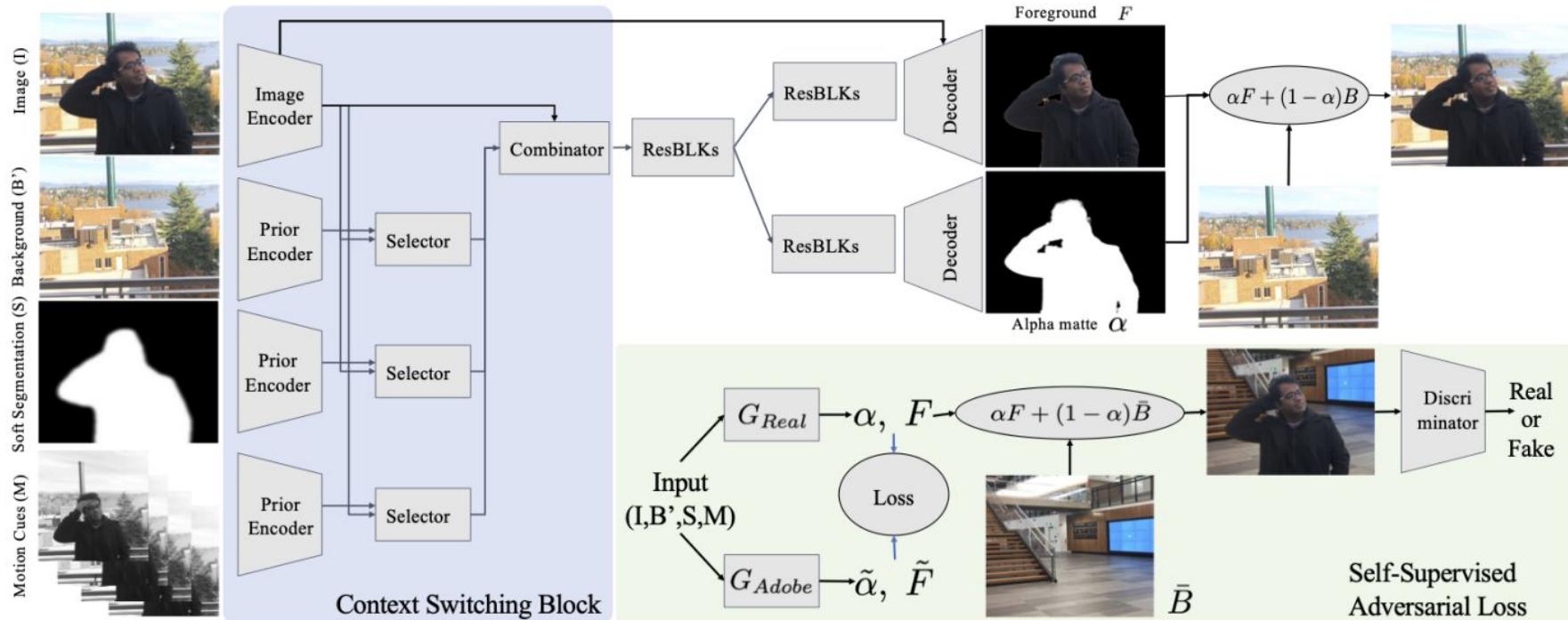
Segmentation



Matting (Ours)

[Image Source](#)

Alpha Matte Learning



Steps

- Supervised Learning Step
- Finetuning Using a GAN
- Using the Supervised Network as a Teacher to prevent trivial solutions

Results

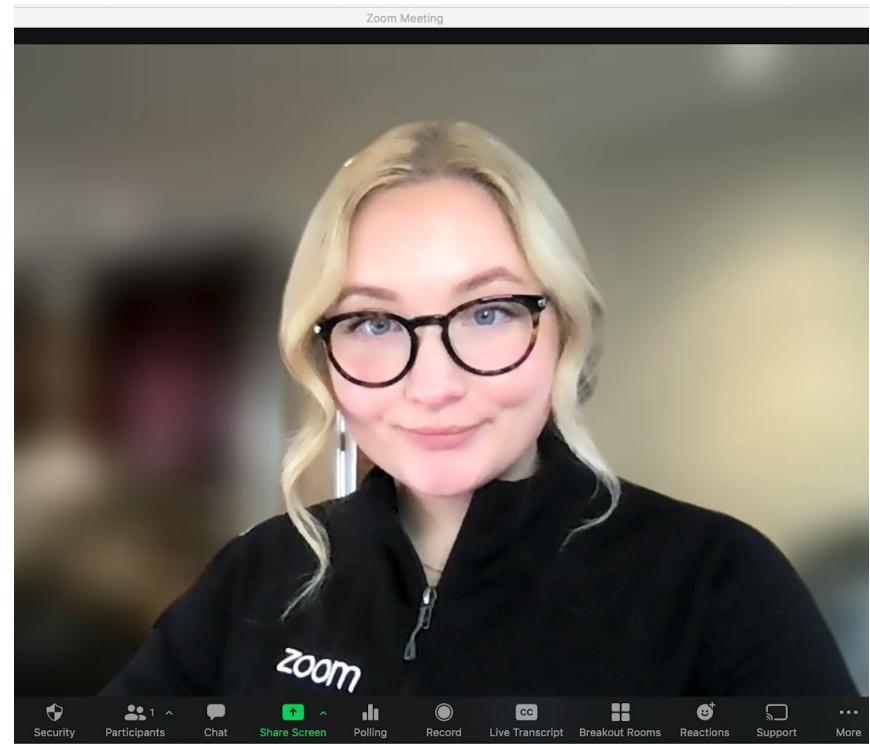


Sengupta, S., Jayaram, V., Curless, B., Seitz, S.M., & Kemelmacher-Shlizerman, I. (2020). Background Matting: The World Is Your Green Screen. 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2288-2297.

Applications



[Image Source](#)



[Image Source](#)

Applications

Kissicve
FOR THE CREATIVES



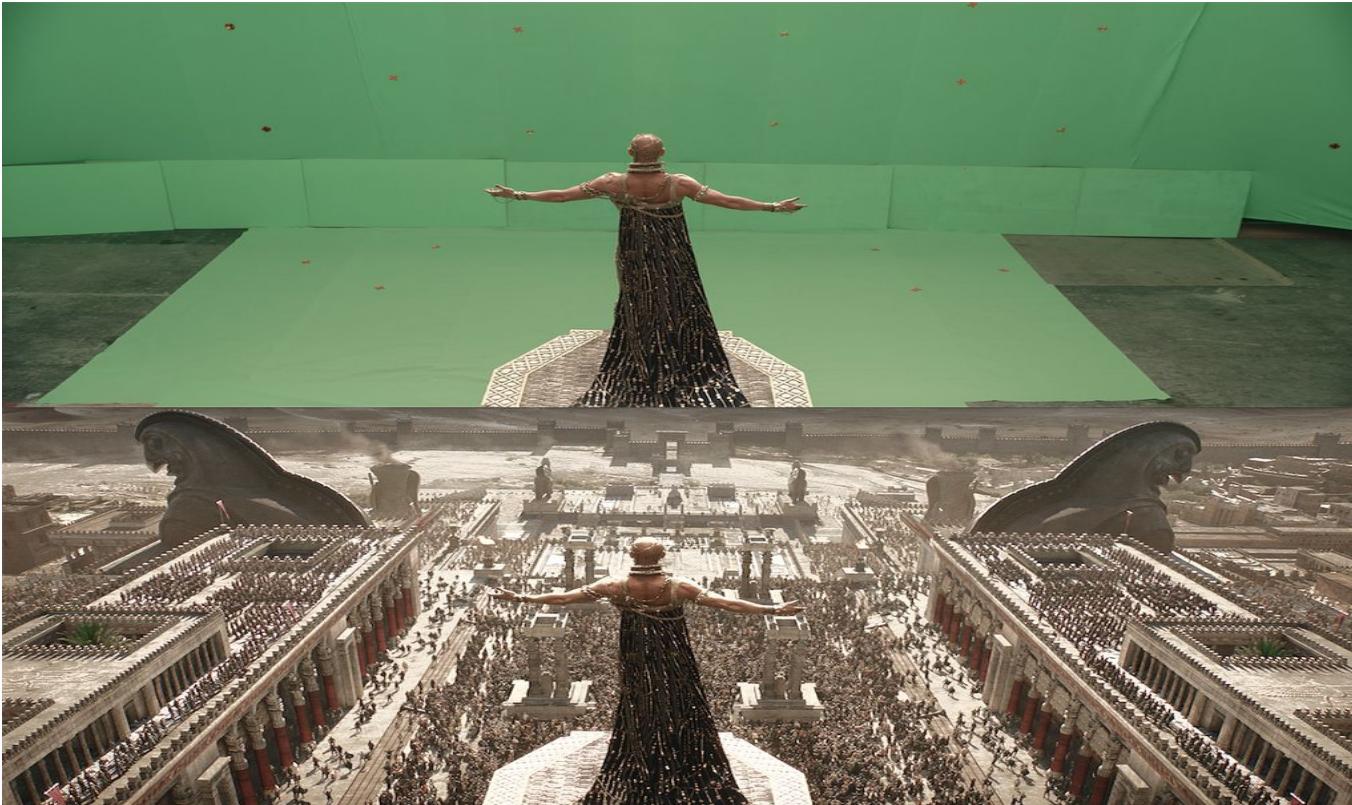
before



after

[Image Source](#)

Applications



[Image Source](#)

Data Compression

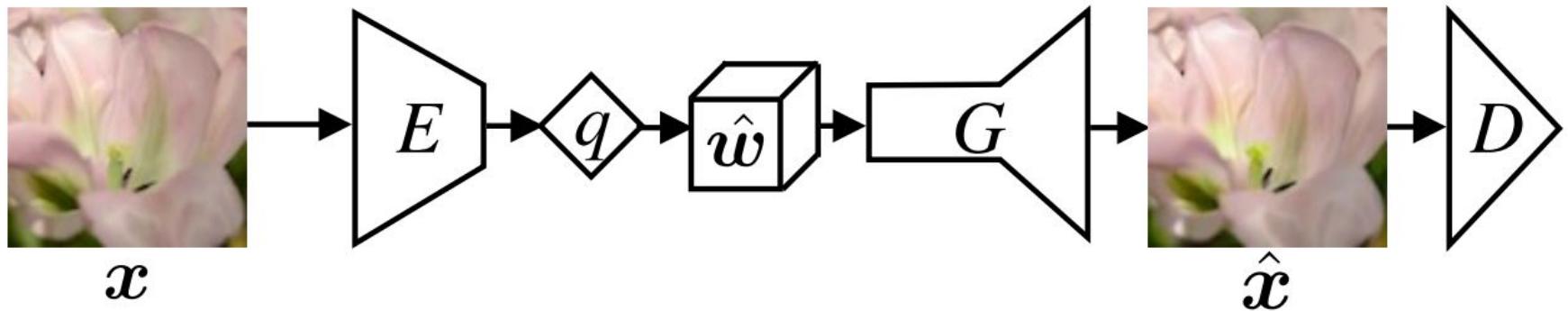
Motivation

- Handcrafted algorithms lack the “understanding” of the content they are compressing
- All pixels are treated equally thus causing distortions and artifacts

JPG Artifacts



Data Compression



Loss for Compression

$$\min_{E,G} \quad \mathcal{L}_{\text{GAN}} + \lambda \mathbb{E}[d(\mathbf{x}, G(\mathbf{z}))] + \beta H(\hat{\mathbf{w}})$$

d: L1/L2 distance

H: Entropy

HiFiC vs JPG

HiFiC



JPG



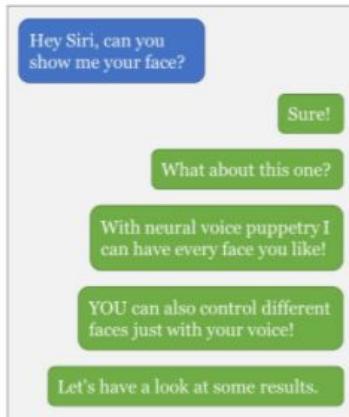
Original



Dynamic Ad Creation

Facial Reenactment

Facial Animation For Digital Assistants

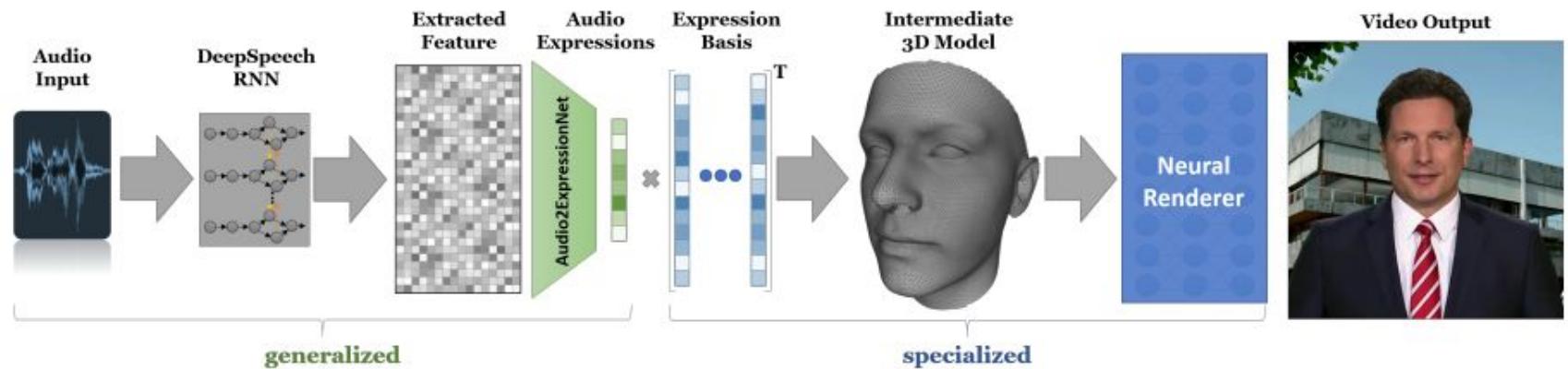


Audio-driven Facial Reenactment

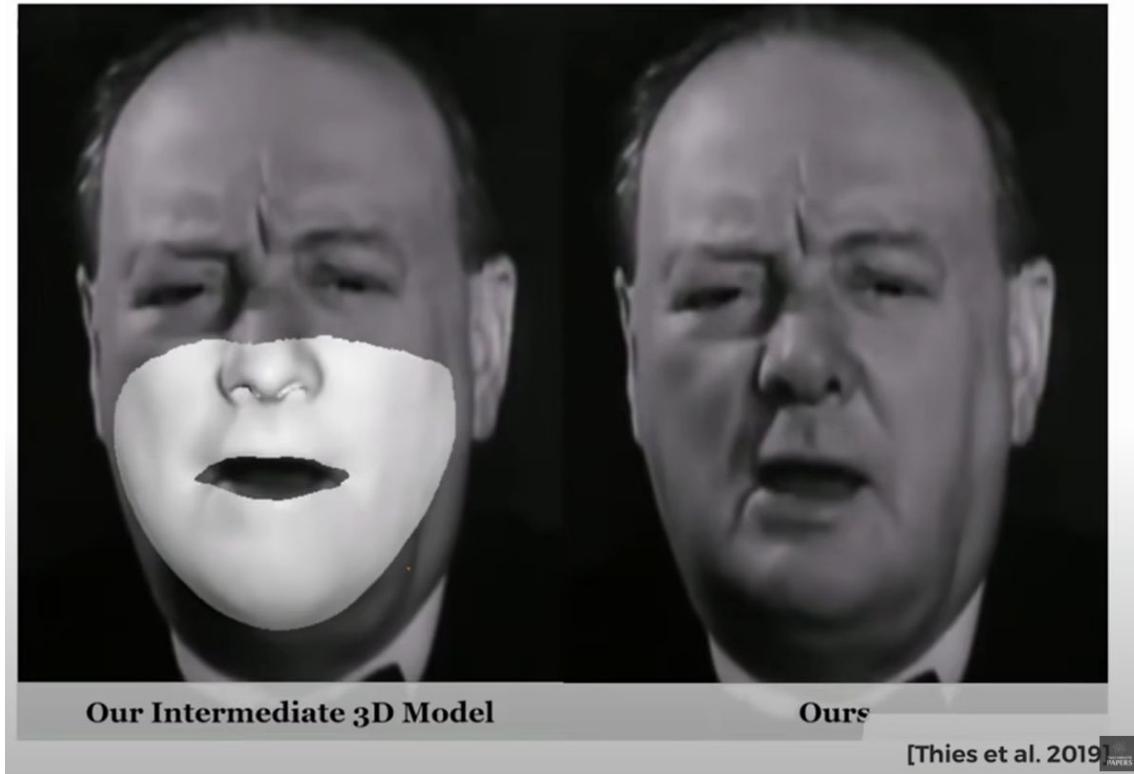


[Image Source](#)

Neural Voice Puppetry

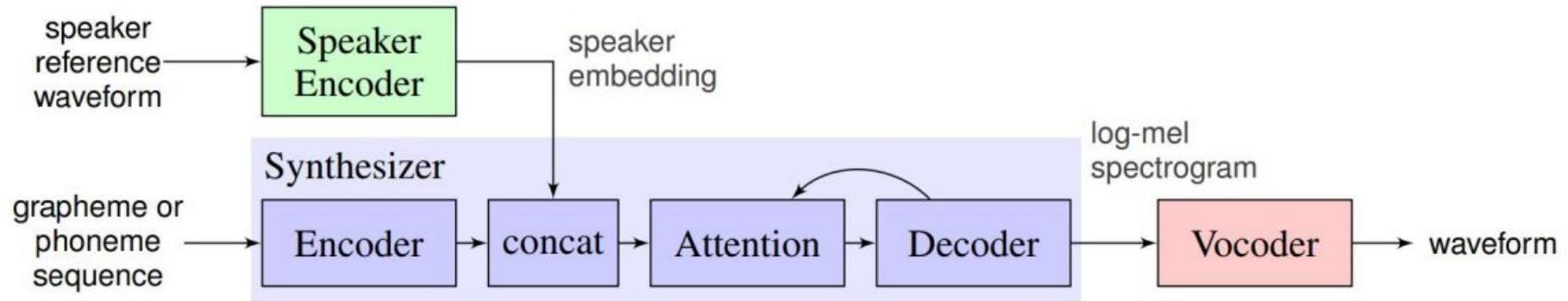


Intermediate 3D Model



[Image Source](#)

Voice Cloning



Charming
Celebrations



Empire Footwear, you've obviously heard of them? Make sure that you buy your shoes from there.

Thank you