

# Image Editing with Optimization (part I)

Jun-Yan Zhu

16-726, Spring 2022

# Style Transfer vs. Image-to-Image Translation

- Data (how to define Style)
  - A single image? A collection of images?
- Applications
  - Photo -> Painting (Neural Style Transfer, Image-to-Image Translation)
  - Photo -> Photo (Image-to-Image Translation, Photo Style Transfer (Color) )
  - Painting -> Photo (Image-to-Image Translation, Deep Image Analogy)
- Algorithms:
  - Patch-based method (i.e., correspondence between output and input)
  - Optimization-based method
  - Feed-forward network
- Loss functions
  - Style Loss: GAN loss, Gram matrix loss
  - Content Loss: Perceptual Loss (L2 reconstruction loss), identity loss, conditional GAN Loss, Cycle-consistency loss, Contrastive Loss (InfoNCE)

# Image Editing with Optimization

$$|\text{Gram}(\hat{y}) - \text{Gram}(y)|$$

optimized output
style image

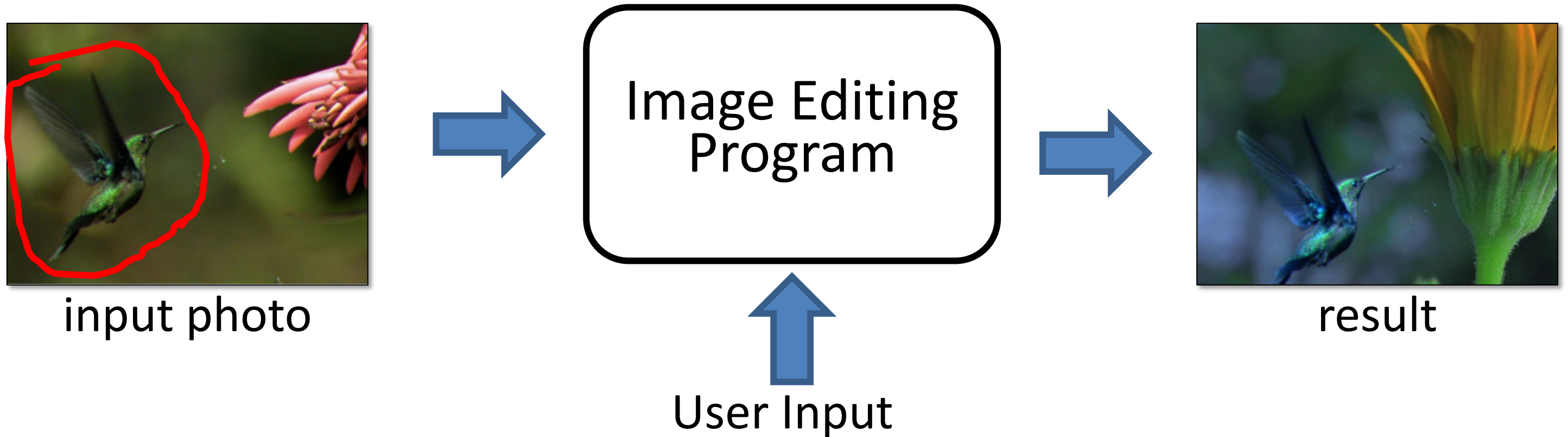
$$+ |\text{F}(\hat{y}) - \text{F}(x)|$$

optimized output
content image

$$\arg \min_{\hat{y}} \mathcal{L}_{\text{style}}(\hat{y}, y) + \lambda \mathcal{L}_{\text{content}}(\hat{y}, x)$$

User input
Input image

# Image Editing with Optimization

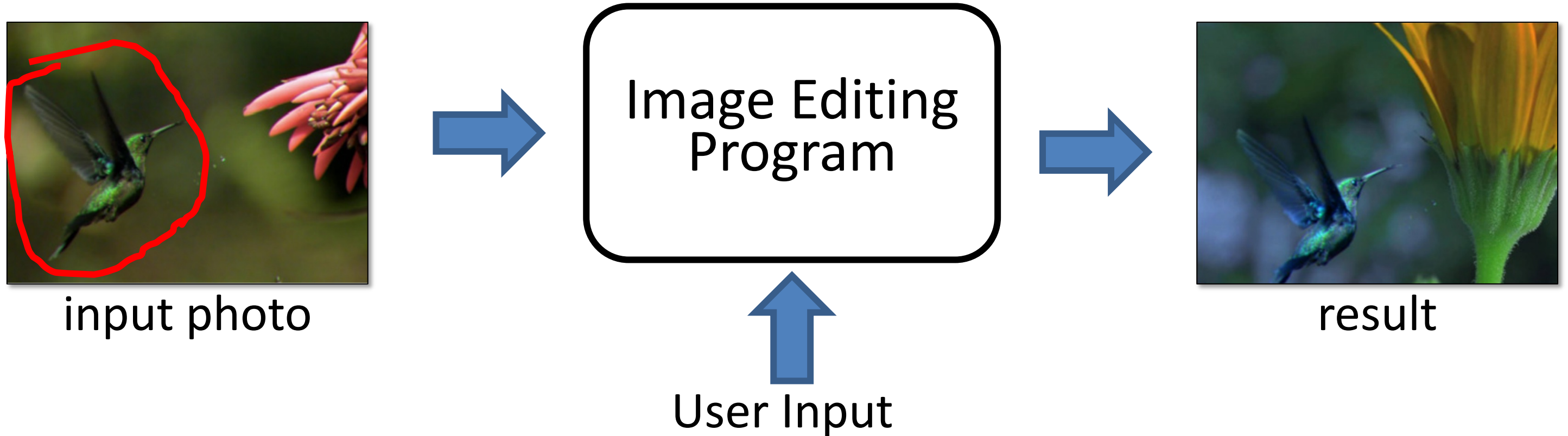


$$\arg \min_{\hat{y}} \mathcal{L}_{\text{background\_boundary}}(\hat{y}, y) + \lambda \mathcal{L}_{\text{source\_gradient}}(\hat{y}, x)$$

result      background      result      object



# Image Editing with Optimization



- Desired output:
- stay close to the input.
  - satisfy user's constraint.

# Image Editing with Optimization

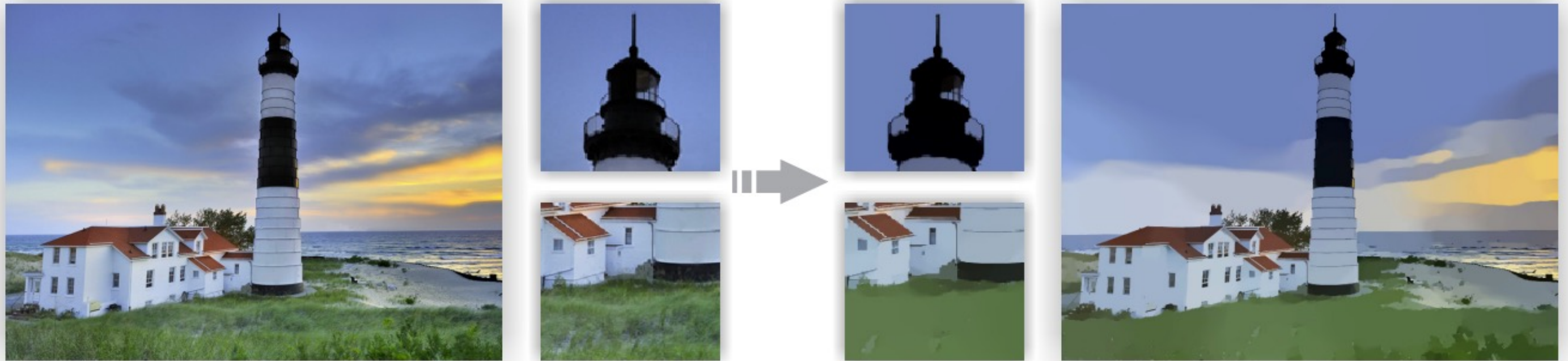


Image Smoothing via L0 Gradient Minimization [Xu et al., SIGGRAPH Asia 2011]

$$\arg \min_{\hat{y}} \{ ||\underset{\substack{\uparrow \\ \text{output}}}{\hat{y}}} - \underset{\substack{\uparrow \\ \text{input}}}{x} || + \lambda C(\hat{y}) \}$$

L0 norm on image gradients  
(the total number of nonzero elements)

# Image Editing with Optimization



Colorization using Optimization [Levin et al., SIGGRAPH 2004]

YUV color space (Y is fixed)

constant: scribbles

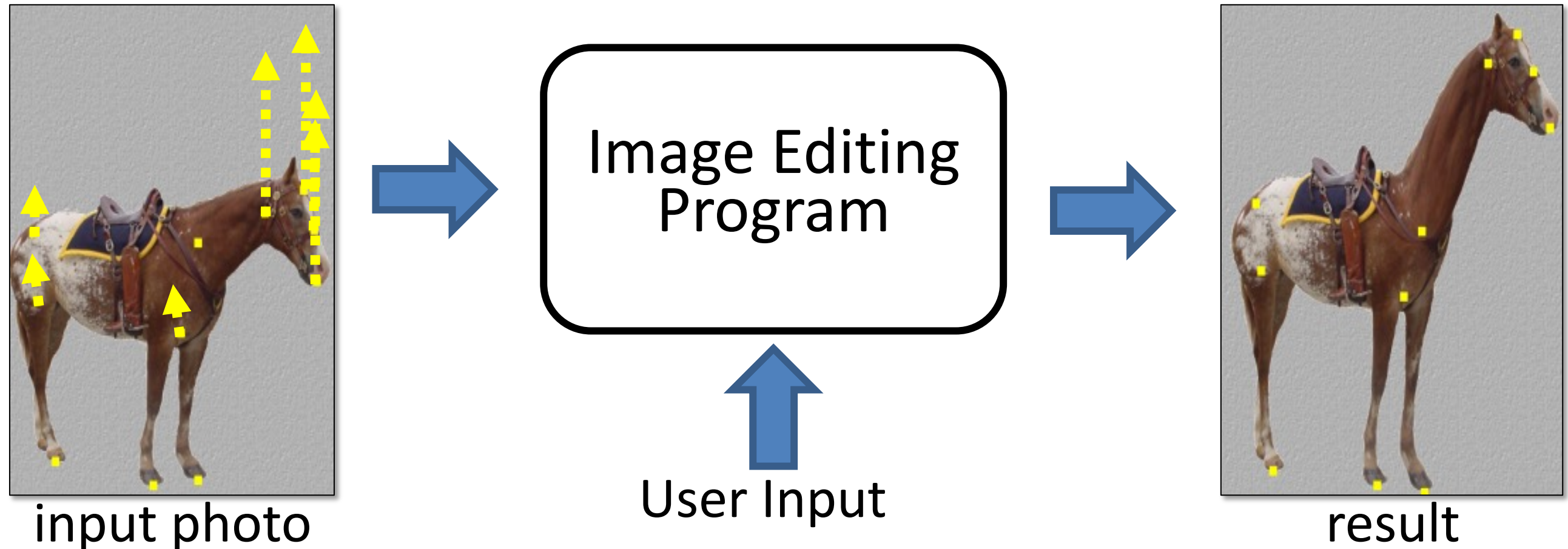
variables: rest of the pixels

$$\sum_{\mathbf{r}} \left( \underset{\substack{\uparrow \\ \text{the color of pixel } r}}{U(\mathbf{r})} - \sum_{\mathbf{s} \in N(\mathbf{r})} \underset{\substack{\nwarrow \\ \text{the color of pixel } s \text{ (} s \text{ is } r\text{'s neighbor)}}}{w_{\mathbf{r}\mathbf{s}} U(\mathbf{s})} \right)^2$$

visual similarity between  $\mathbf{r}$  and  $\mathbf{s}$   
Intensity, location, edge, motion, etc.



# Image Editing with Optimization



Moving least squares + transformation parameters.

- Desired output:
- stay close to the input.
  - satisfy user's constraint.



# So far so good

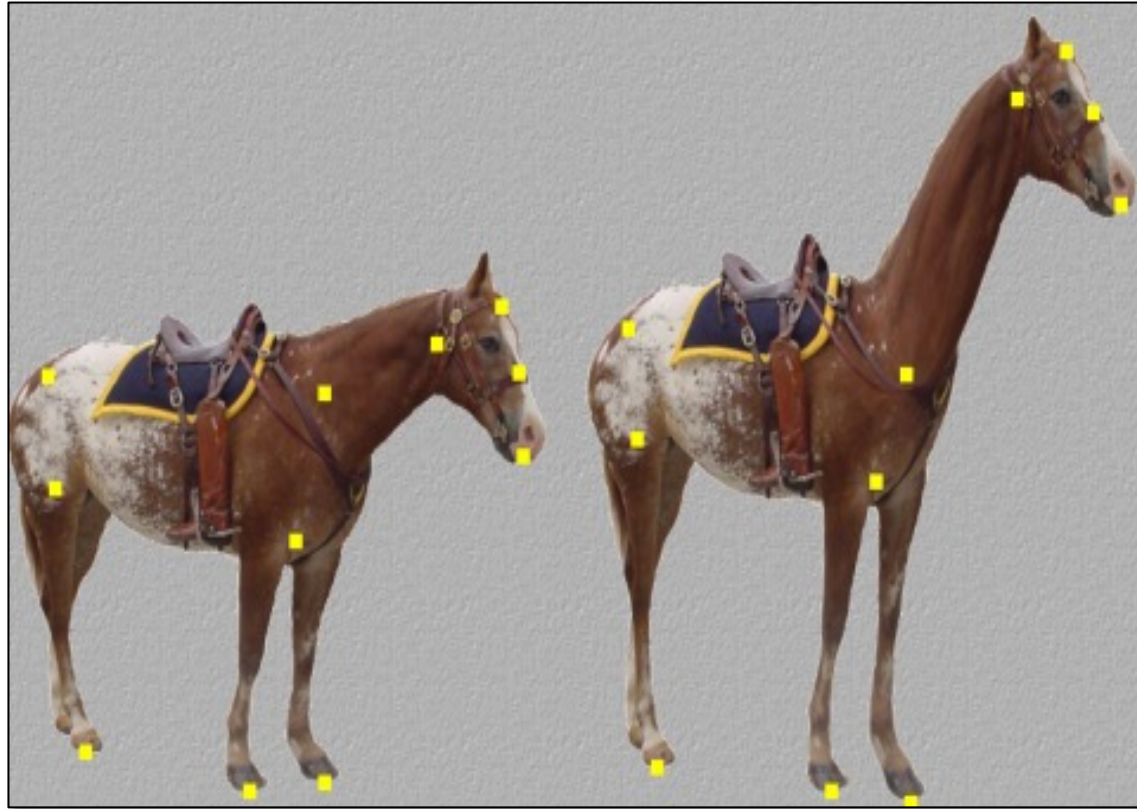


Image Warping



Image Composition

# Things can get really bad



Image Warping



Image Composition

## The lack of “safety wheels”

# Adding the “safety wheels”



Input Photo

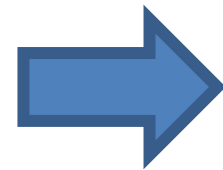
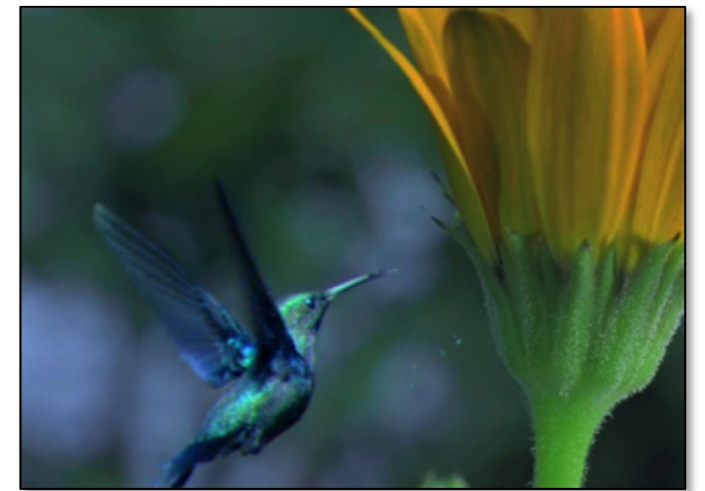
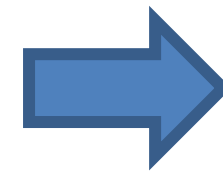


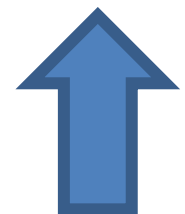
Image Editing  
Program



Output Result



User Input



Natural Image  
Manifold

A desired output:

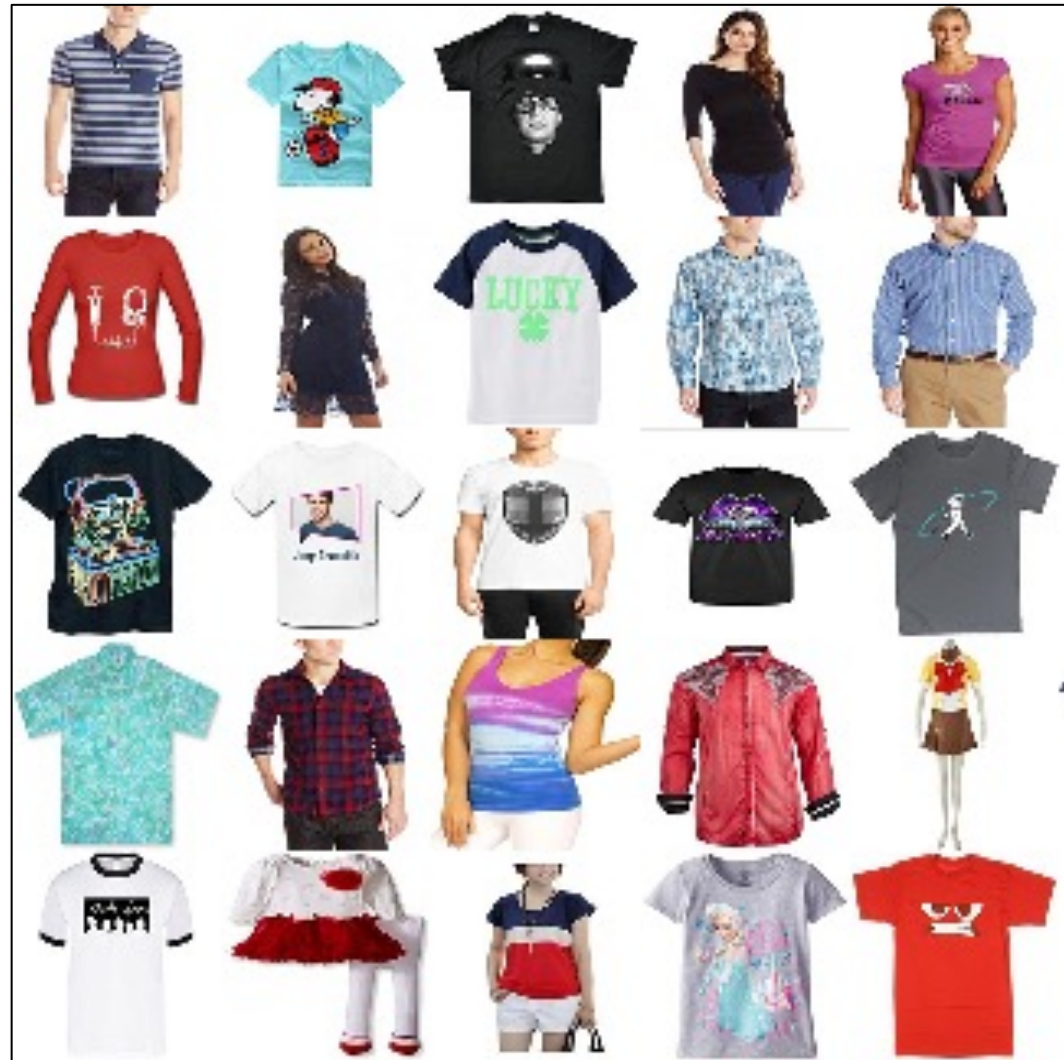
- stay close to the input.
- satisfy user's constraint.
- Lie on the natural image manifold

# Learning Natural Image Manifold

- Deep generative models:  $G(z) : z \rightarrow x$ 
  - Generative Adversarial Network (**GAN**)  
(e.g., DCGAN, StyleGAN2, BigGAN)
  - Variational Auto-Encoder (**VAE**)  
(e.g., VQ-VAE2)
  - Flow-based models (e.g., RealNVP, Glow)...
  - Diffusion models (e.g., DDPM, DDIM)
  - ...



# GAN as Manifold Approximation



Sample training images  
from “Amazon Shirts”



Random image samples  
from Generator  $G(z)$

# Traverse on the GAN Manifold

$G(z_0)$

Linear Interpolation in  $z$  space:  $G(z_0 + t \cdot (z_1 - z_0))$

$G(z_1)$



## Limitations of DCGAN:

- not photo-realistic enough, low resolution
- produce images randomly, no user control



# Changing Variables

- Traditional method: Optimizing the image

$$\hat{y}^* = \arg \min_{\hat{y}} \mathcal{L}(\underset{\substack{\uparrow \\ \text{input}}}{x}, \underset{\substack{\downarrow \\ \text{user constraint}}}{y}, \underset{\substack{\uparrow \\ \text{output}}}{\hat{y}})$$

- New method: Optimizing the latent code

$$z^* = \arg \min_z \mathcal{L}(\underset{\substack{\uparrow \\ \text{input}}}{x}, \underset{\substack{\downarrow \\ \text{user constraint}}}{y}, \underset{\substack{\uparrow \\ \text{Latent code}}}{G(z)})$$

Generator

# Projecting and Editing an Image



original photo

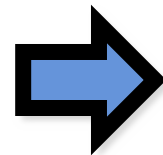


different degree of image manipulation

Project 



projection on manifold



 Edit Transfer



transition between the original and edited projection



# Projecting and Editing an Image



original photo

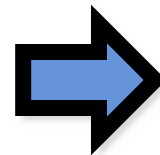


different degree of image manipulation

Project 



projection on manifold



 Edit Transfer



transition between the original and edited projection

# Projecting an Image into GAN Manifold

Input: real image  $x$   
Output: latent vector  $z$

**Optimization**

$$z^* = \arg \min_z \mathcal{L}(G(z), x)$$

Reconstruction loss

Generative model



0.196



0.238



0.332

# Projecting an Image into GAN Manifold

Input: real image  $x$   
Output: latent vector  $z$

**Optimization**

$$z^* = \arg \min_z \mathcal{L}(G(z), x)$$

**Inverting Network**  $z = E(x)$

$$E = \arg \min_E \mathbb{E}_x \underbrace{\mathcal{L}(G(E(x)), x)}_{\text{Auto-encoder with a fixed decoder}}$$

Auto-encoder  
with a fixed decoder



0.196



0.238



0.332



0.218



0.242



0.336

# Projecting an Image into GAN Manifold

Input: real image  $x$   
Output: latent vector  $z$

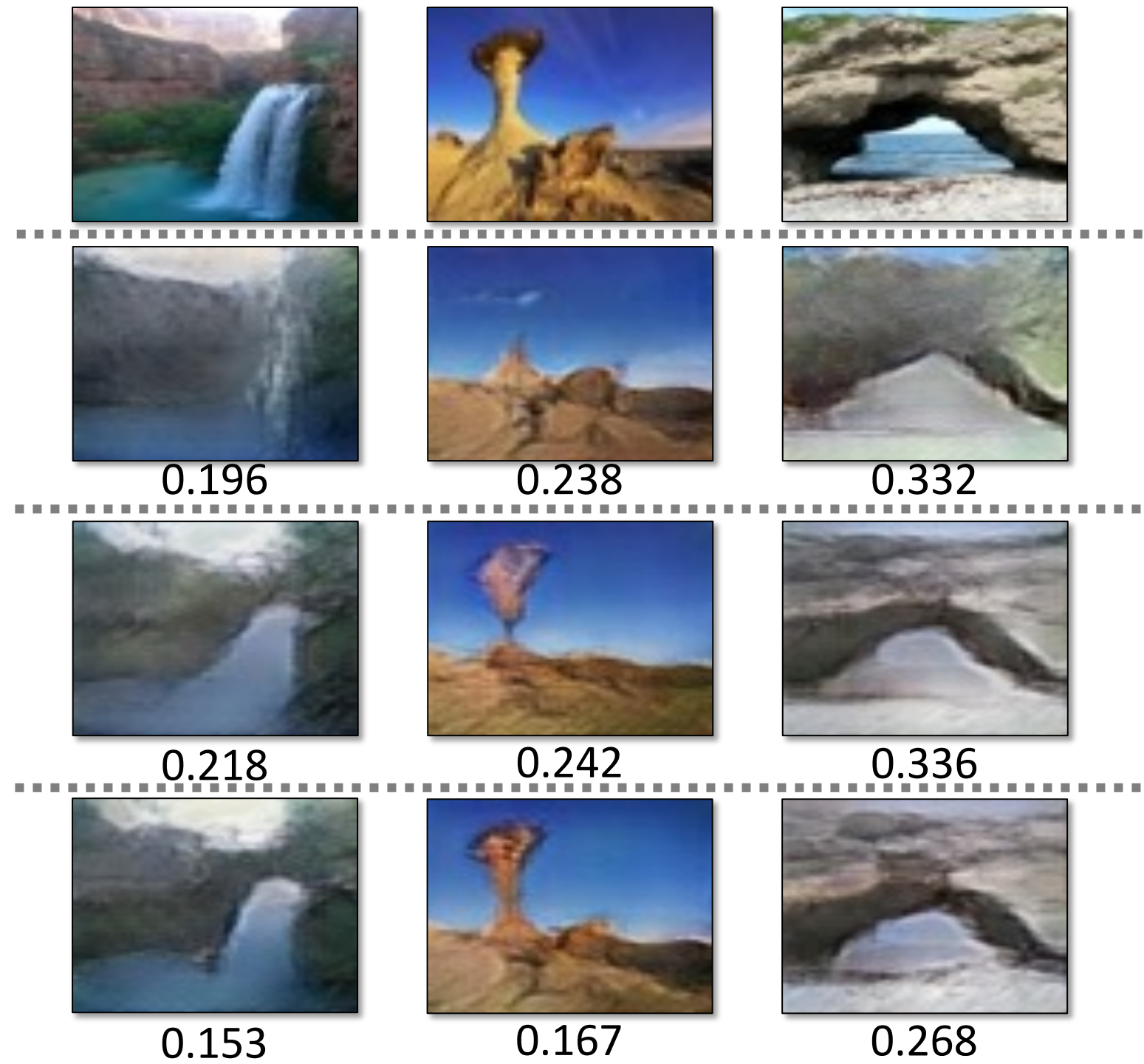
**Optimization**

$$z^* = \arg \min_z \mathcal{L}(G(z), x)$$

**Inverting Network**  $z = E(x)$

$$E = \arg \min_E \mathbb{E}_x \mathcal{L}(G(E(x)), x)$$

**Hybrid Method**  
Use the **network** as initialization  
for the **optimization** problem





# Manipulating the Latent Code



original photo



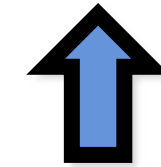
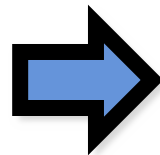
different degree of image manipulation

Project



projection on manifold

Editing UI



Edit Transfer



transition between the original and edited projection

# Manipulating the Latent Code

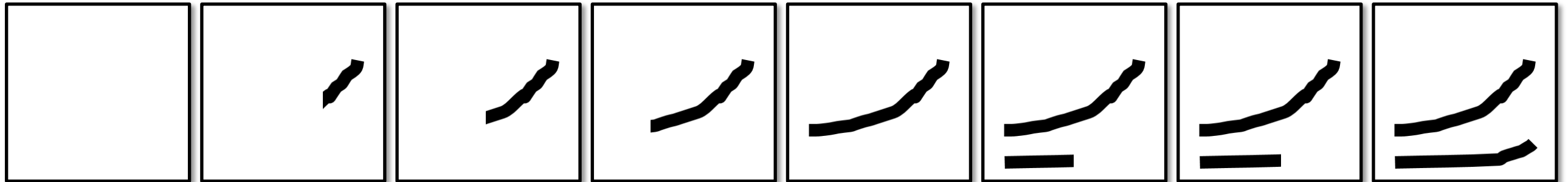
constraint violation loss  $L_g$

user guidance image

Objective: 
$$z^* = \arg \min_{z \in \mathbb{Z}} \left\{ \underbrace{\sum_g (\mathcal{L}_g(G(z)) \underbrace{v_g}_{\text{data term}})}_{\text{data term}} + \underbrace{\lambda_s \cdot \|z - z_0\|_2^2}_{\text{manifold smoothness}} \right\}.$$

Guidance

$v_g$



$G(z)$



$z_0$

# Post-Processing



original photo

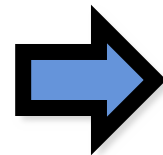


different degree of image manipulation

Project 



projection on manifold



 Edit Transfer

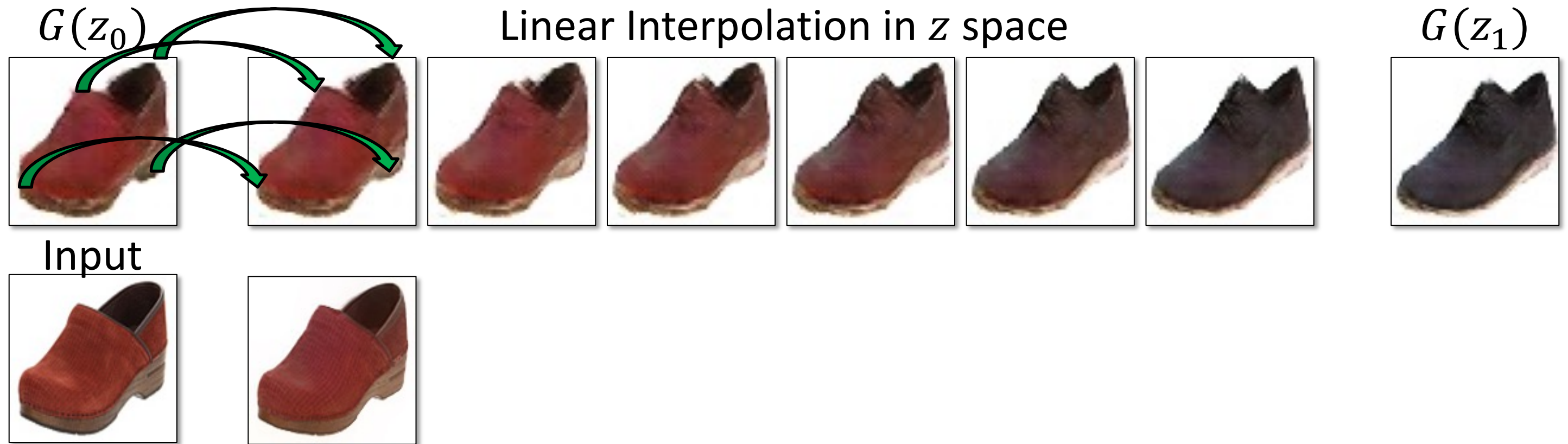


transition between the original and edited projection

# Edit Transfer

**Motion** ( $\mathbf{u}, \mathbf{v}$ ) + **Color** ( $A_{3 \times 4}$ ): estimate per-pixel geometric and color variation

$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x + u, y + v, t + 1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dx dy$$

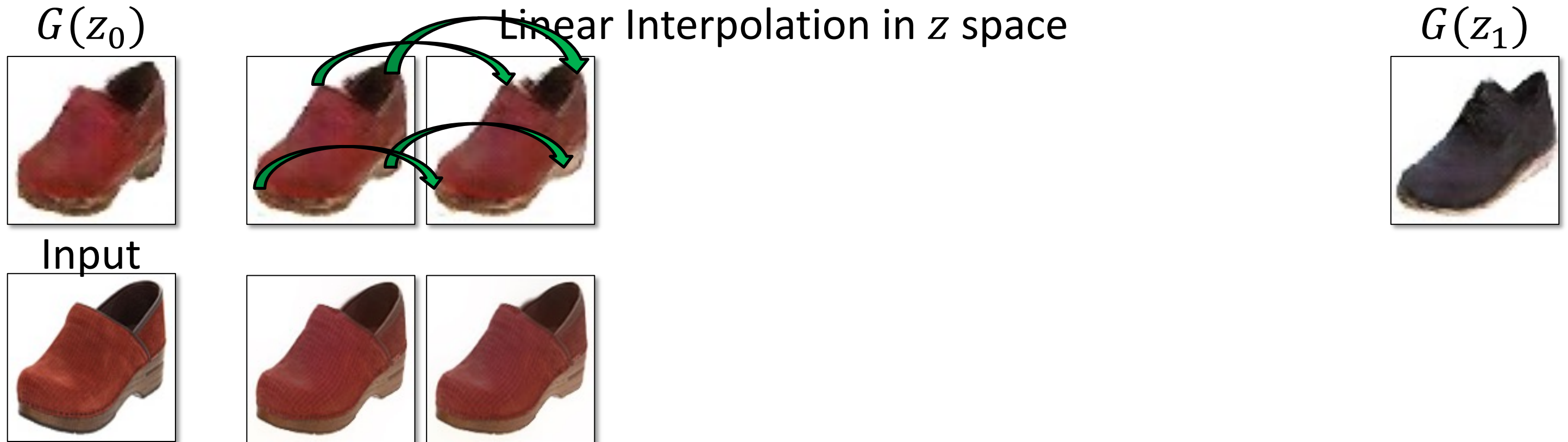




# Edit Transfer

**Motion** ( $\mathbf{u}, \mathbf{v}$ ) + **Color** ( $A_{3 \times 4}$ ): estimate per-pixel geometric and color variation

$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x + u, y + v, t + 1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dx dy$$



# Edit Transfer

**Motion** ( $\mathbf{u}, \mathbf{v}$ ) + **Color** ( $A_{3 \times 4}$ ): estimate per-pixel geometric and color variation

$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x + u, y + v, t + 1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dx dy$$

$G(z_0)$

Linear Interpolation in z space

$G(z_1)$



Input

Result



# Image Manipulation Demo

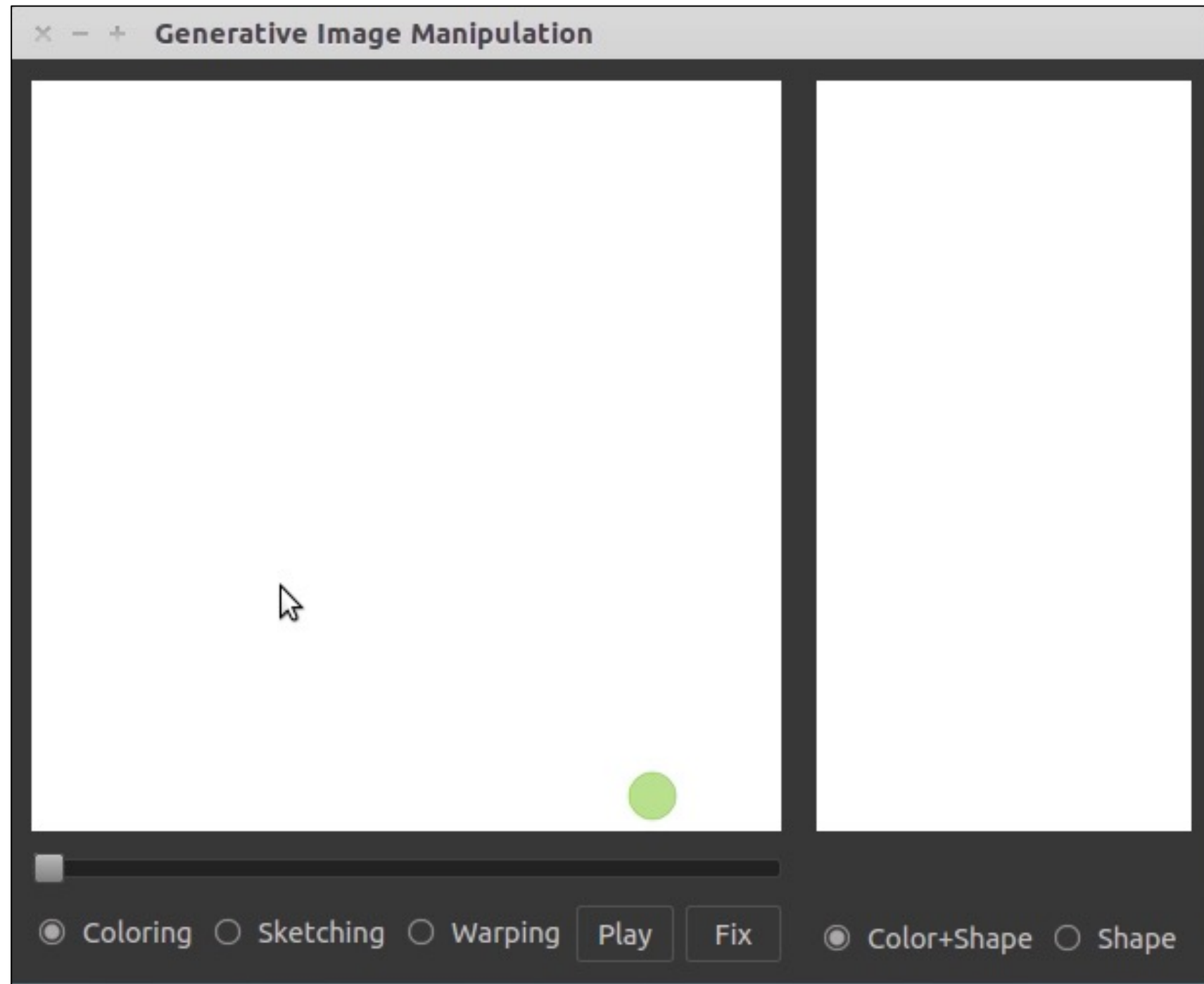


# Image Manipulation Demo

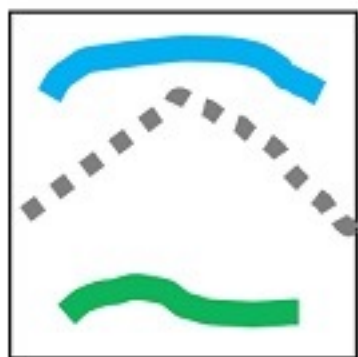
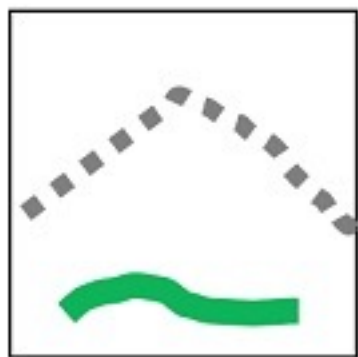
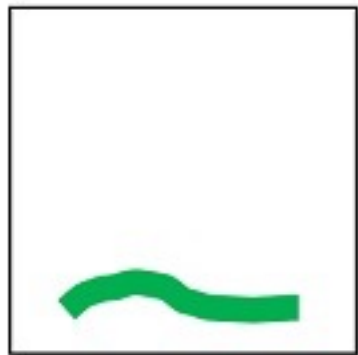




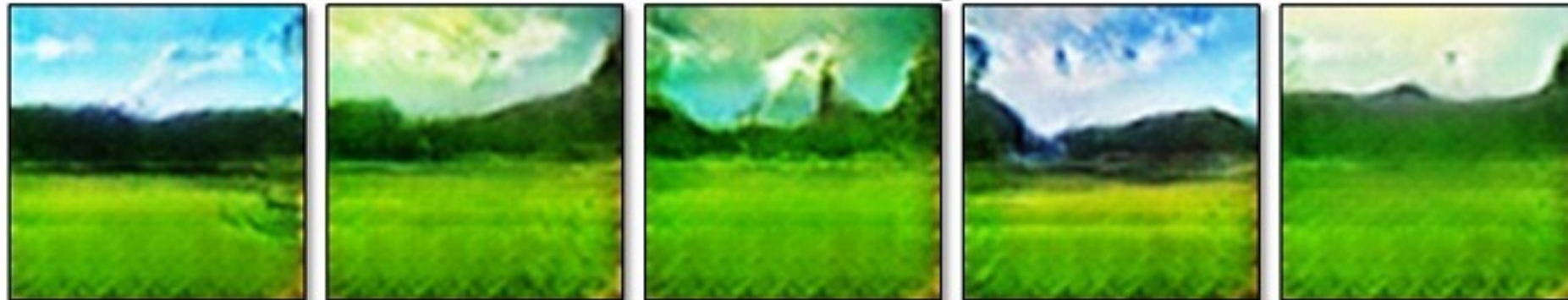
# Interactive Image Generation



User edits



Generated images



Color



Sketch

# Changing Variables

- Traditional method: Optimizing the image

$$\hat{y}^* = \arg \min_{\hat{y}} \mathcal{L}(\underset{\substack{\uparrow \\ \text{input}}}{x}, \underset{\substack{\downarrow \\ \text{user constraint}}}{y}, \underset{\substack{\uparrow \\ \text{result}}}{\hat{y}})$$

- New method: Optimizing the latent code

$$z^* = \arg \min_z \mathcal{L}(\underset{\substack{\uparrow \\ \text{input}}}{x}, \underset{\substack{\downarrow \\ \text{user constraint}}}{y}, \underset{\substack{\uparrow \\ \text{Latent code}}}{G(z)})$$

Generator

# Projecting and Editing an Image



original photo

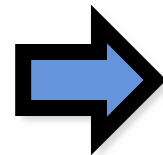


different degree of image manipulation

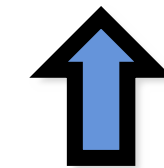
Project



projection on manifold



Editing UI



Post-processing



transition between the original and edited projection



# Image Editing with GANs

- Step 1: Image Projection/Reconstruction

$$z_0 = \arg \min_z \mathcal{L}(G(z), x)$$

- Step 2: Manipulating the latent code

$$z_1 = z_0 + \Delta z$$

- Step 3: Generate the edited result

$$G(z_1)$$

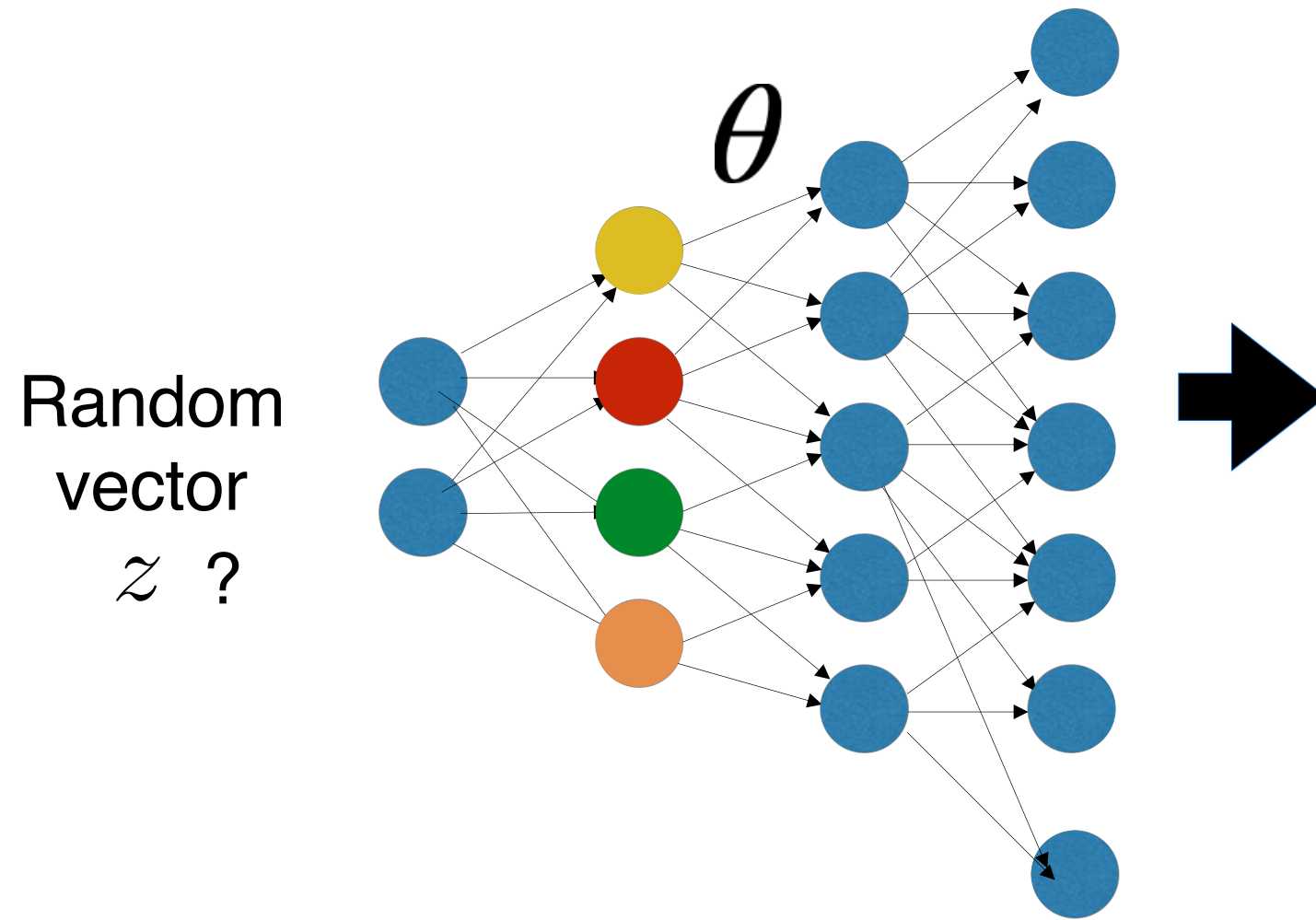
# Image Projection with GANs

# Image Reconstruction (high-res images, Big Models)



Original image  $x$

# Image Reconstruction (high-res images, Big Models)

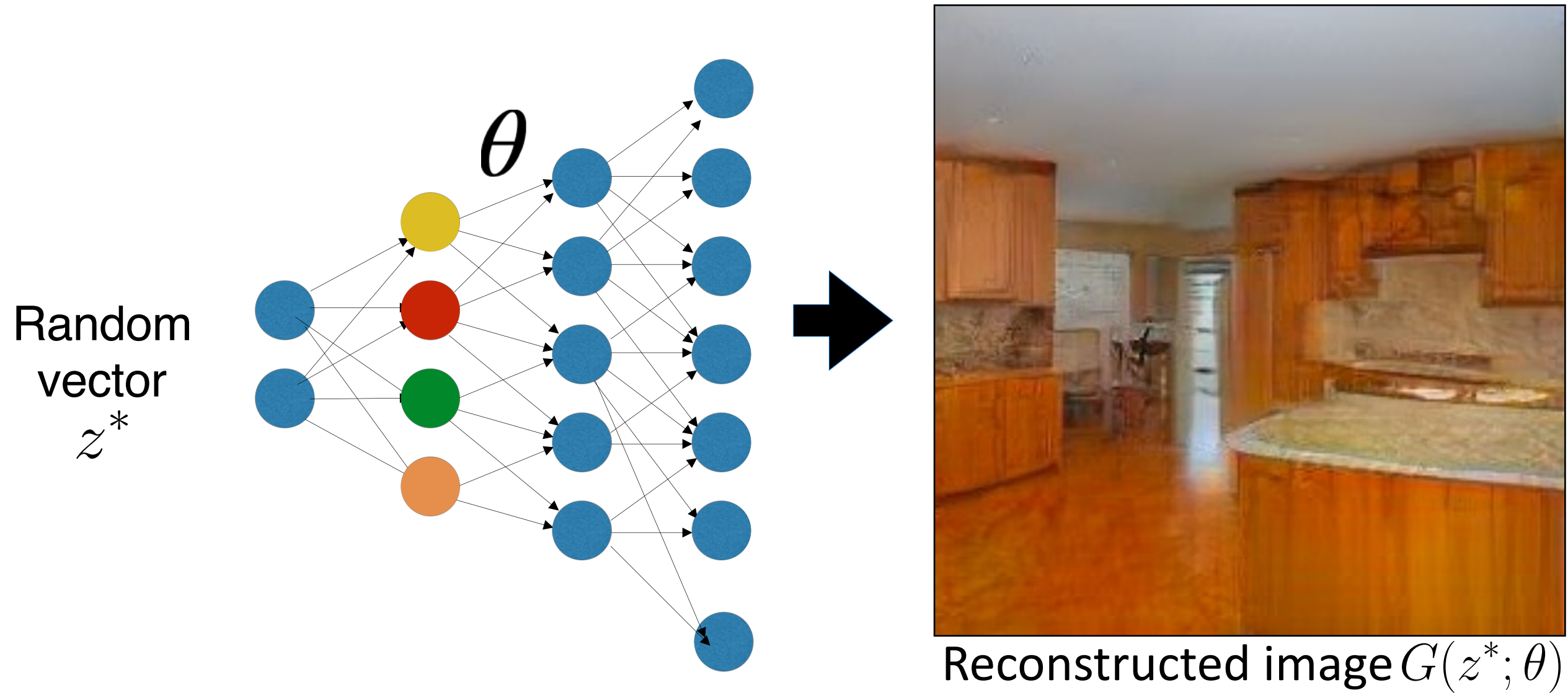


Original image  $x$

$$z^* = \arg \min_z \mathcal{L}(G(z; \theta), x)$$



# Image Reconstruction (high-res images, Big Models)



$$z^* = \arg \min_z \mathcal{L}(G(z; \theta), x)$$

# Find the Differences...



Original image



GAN reconstructed image



# Find the Differences...



Original image

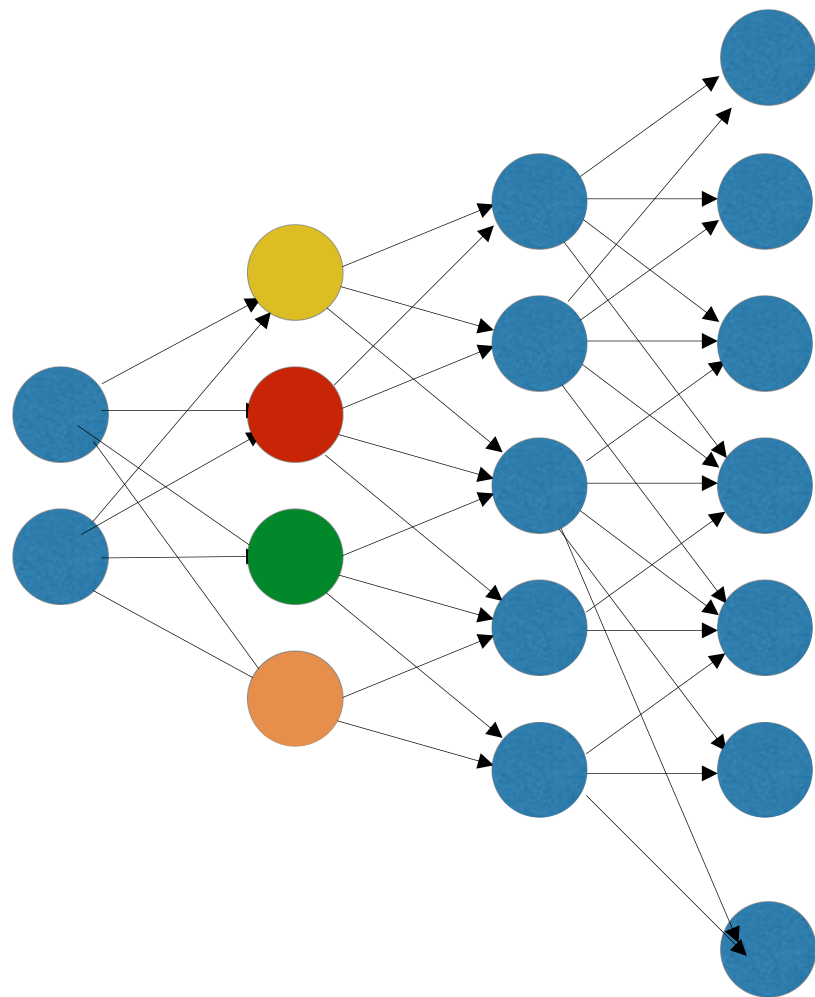


GAN reconstructed image



Original image

Random  
vector  
 $z^*$



Reconstructed image  $G(z^*; \theta)$

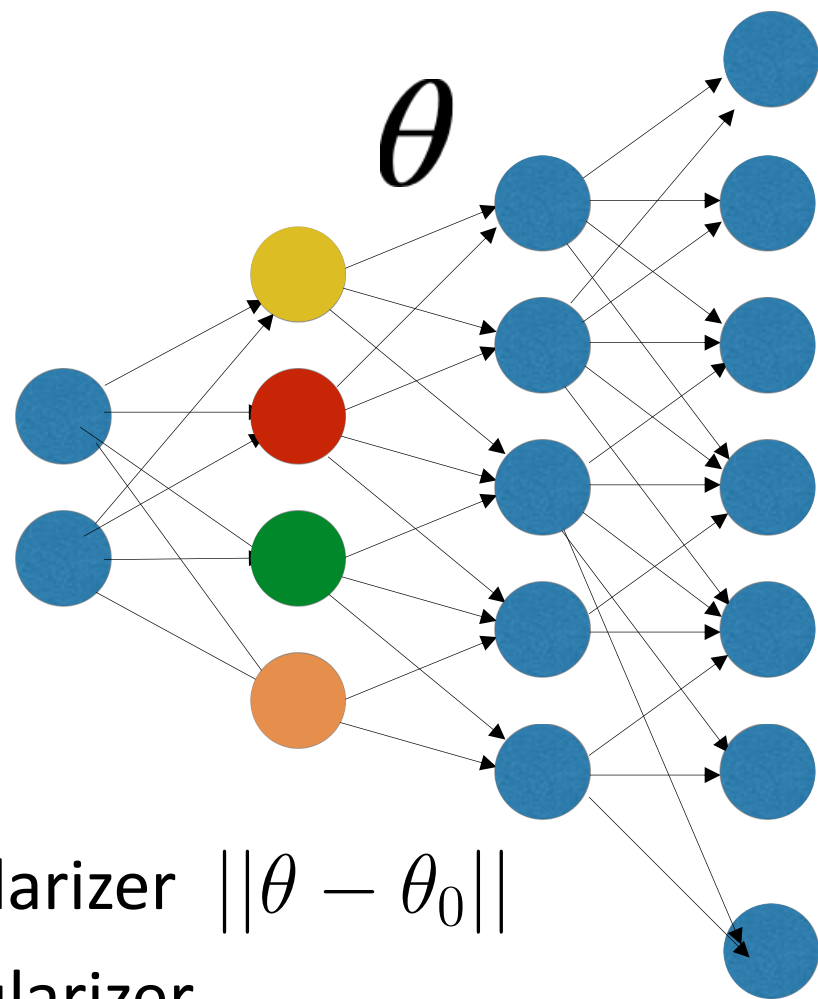
$$z^* = \arg \min_z \mathcal{L}(G(z; \theta), x)$$





Original image

Random  
vector  
 $z^*$



Weight space regularizer  $||\theta - \theta_0||$

Feature space regularizer



Reconstructed image  $G(z^*; \theta)$

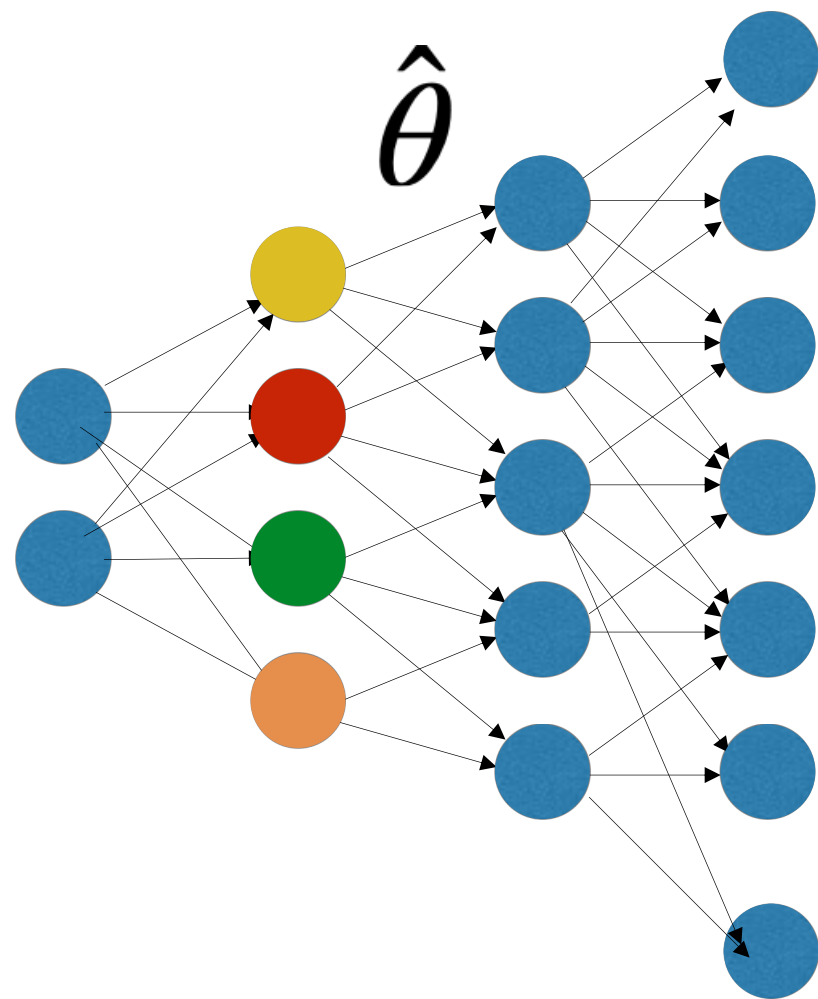
$$z^*, \theta^* = \arg \min_{z, \theta} \mathcal{L}(G(z; \theta), x)$$

← Regularizer



Original image

Random  
vector  
 $z^*$



Reconstructed image  $G(z^*; \theta^*)$

$$z^*, \theta^* = \arg \min_{z, \theta} \mathcal{L}(G(z; \theta), x) + R(\theta) \leftarrow \text{Regularizer}$$



# Reconstructing a Real Photo



Original image



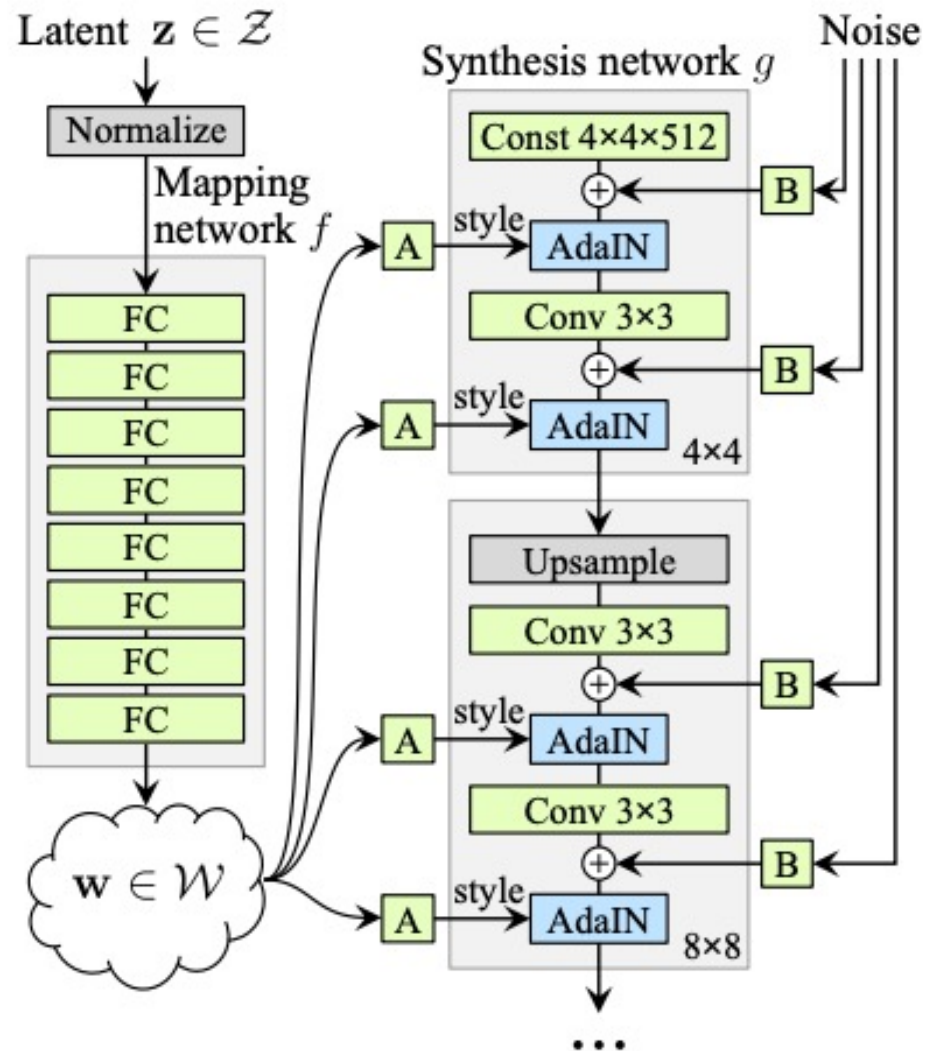
With  $z^*$



With  $z^*$  and  $\theta^*$

Semantic Photo Manipulation [Bau, Strobel, Peebles, Wulff, Zhou, Zhu, Torralba, SIGGRAPH 2019]  
Inspired by Deep Image Prior [Ulyanov et al.] and Deep Internal learning [Shocher et al.]

# Using Different Layers



Optimizing the latent code

$$z^* = \arg \min_z \mathcal{L}(G(z), x)$$

Optimizing the style code

$$w^* = \arg \min_w \mathcal{L}(g(w), x)$$

Optimizing the extended style code

$$w_+^* = \arg \min_{w_+} \mathcal{L}(g(w_+), x)$$



# Using Different Layers: w space



StyleGAN — generated images



StyleGAN2 — generated images

# Using Different Layers: w space



StyleGAN2 — real images



# Using Different Layers: w+ space



All the results are reconstructed using Face Model



# Reconstruction $\neq$ Editing



Interpolations between two images



# Reconstruction $\neq$ Editing



Interpolations between two images