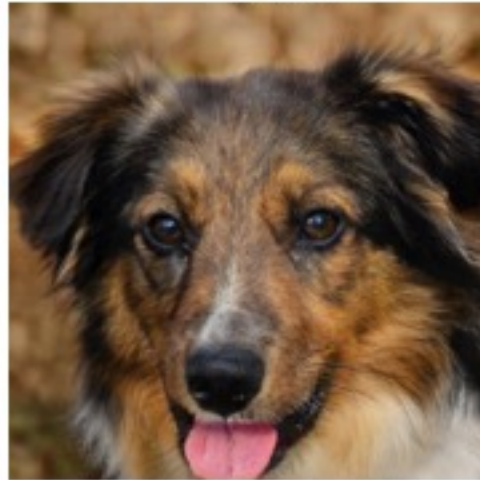


Input



Happy



Big Eyes



Golden Fur



Bulldog



Image Editing with Optimization (part II)

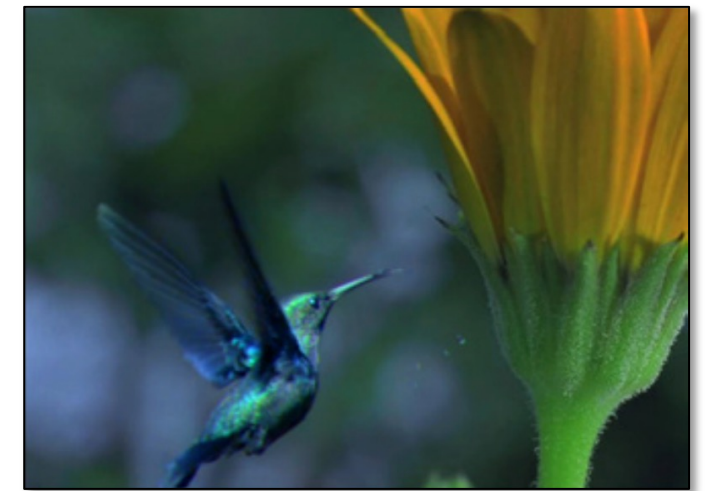
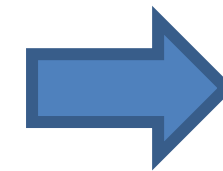
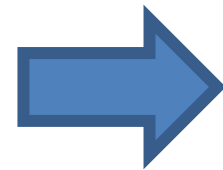
Jun-Yan Zhu

16-726, Spring 2023

Image Editing with Optimization



input photo



result

User Input

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

Image Editing with Optimization

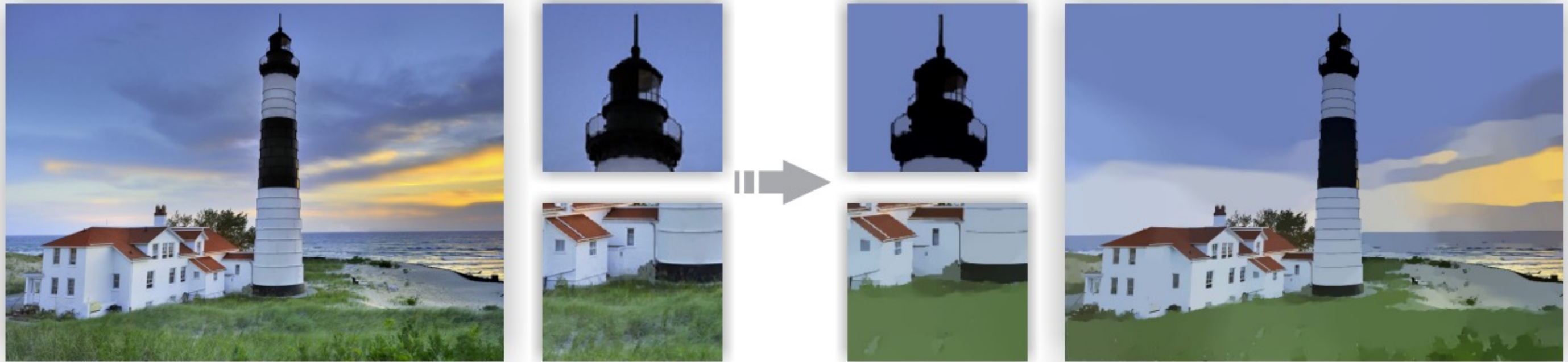


Image Smoothing via L0 Gradient Minimization [Xu et al., SIGGRAPH Asia 2011]

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

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

Things can get really bad



Image Warping



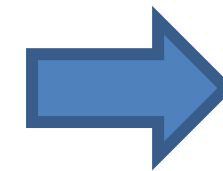
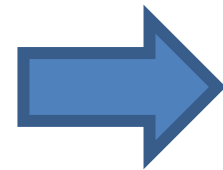
Image Composition

The lack of “safety wheels”

Adding the “safety wheels”



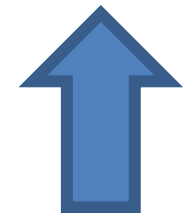
Input Photo



Output Result



User Input



A desired output:

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

Changing Variables

- Traditional method: Optimizing the image

$$\hat{y}^* = \arg \min_{\hat{y}} \mathcal{L}(x, y, \hat{y})$$

input user constraint output

- New method: Optimizing the latent code

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

input user constraint Latent code

Generator

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

Baseline

- Baseline: Optimizing the latent code

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

z^* and z_0 are used interchangeably

Find the Differences...



Original image



GAN reconstructed image

Baseline

- Baseline: Optimizing the latent code

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

How to Improve GANs Projection

- Baseline: Optimizing the latent code

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

- Generator fine-tuning:

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

Generator Fine-tuning (Progressive GANs)



Original image



With z^*

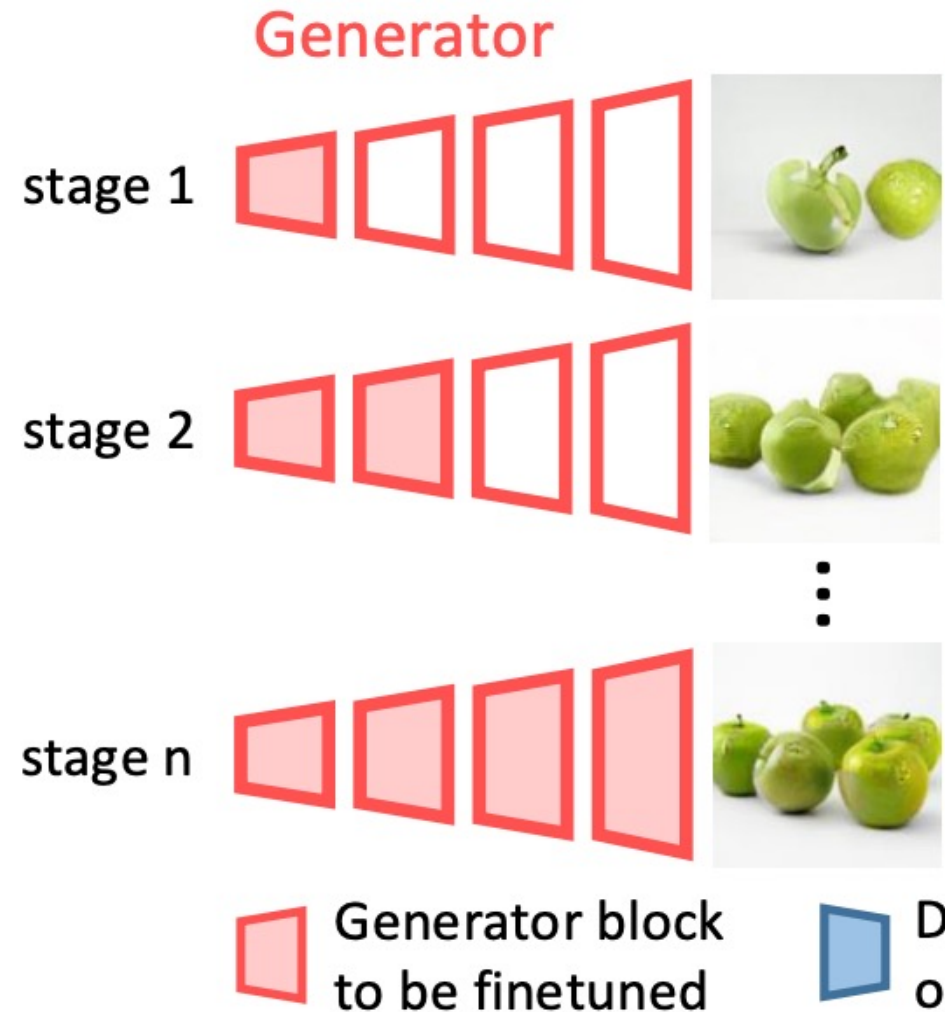


With z^* and θ^*

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.]

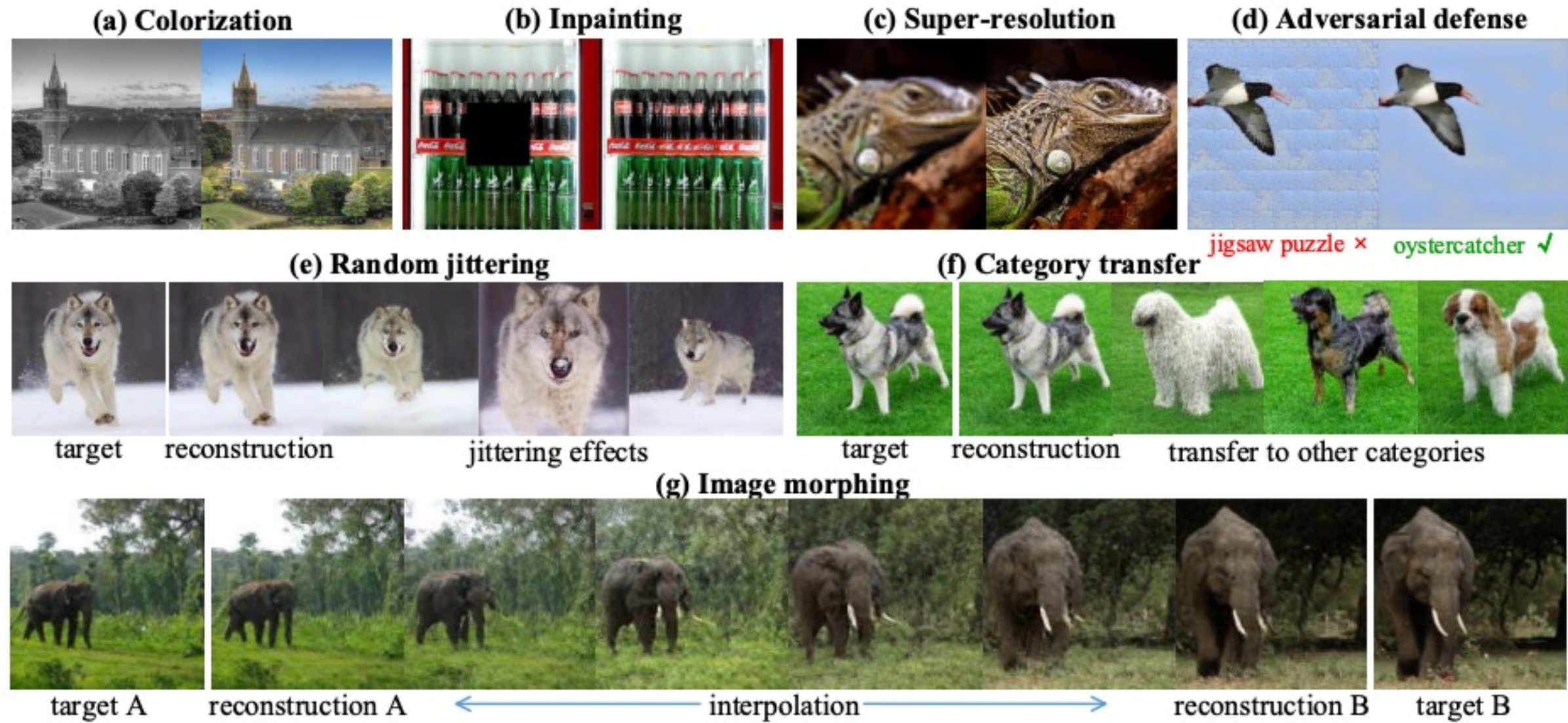
Generator Fine-tuning (BigGAN)



Progressive Reconstruction

- First match semantics
- Then match color and textures

Generator Fine-tuning (BigGAN)



Exploiting Deep Generative Prior for Versatile Image Restoration and Manipulation

[Pan et al., ECCV 2020]

How to Improve GANs Projection

- Baseline: Optimizing the latent code

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

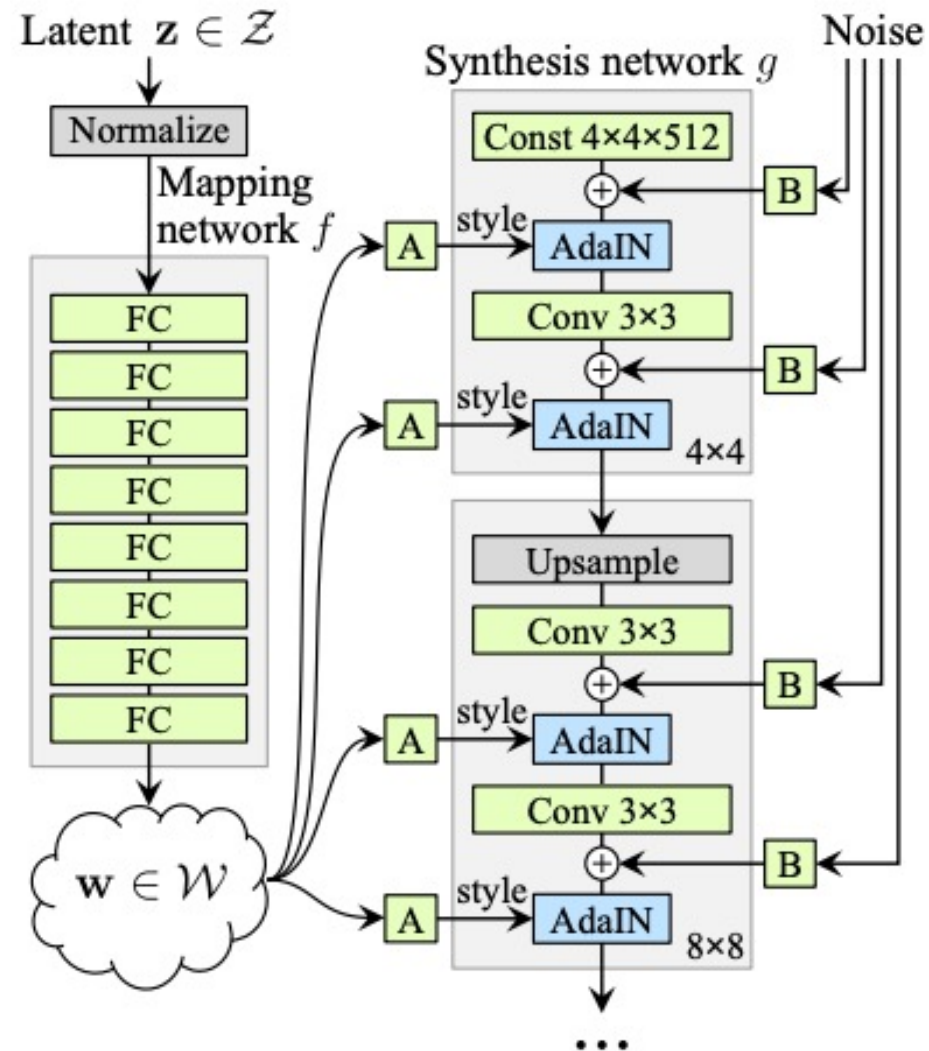
- Generator fine-tuning:

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

- Optimizing intermediate features

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

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

Input



Reconstruction

Using Different Layers: w+ Space



All the results are reconstructed via the StyleGAN Face model.

How to Improve GANs Projection

- Baseline: Optimizing the latent code

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

- Generator fine-tuning:

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

- Optimizing intermediate features

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



Used together

Generator Fine-tuning with $w+$ Space



Pivotal Tuning for Latent-based Editing of Real Images [Roich et al., 2021]

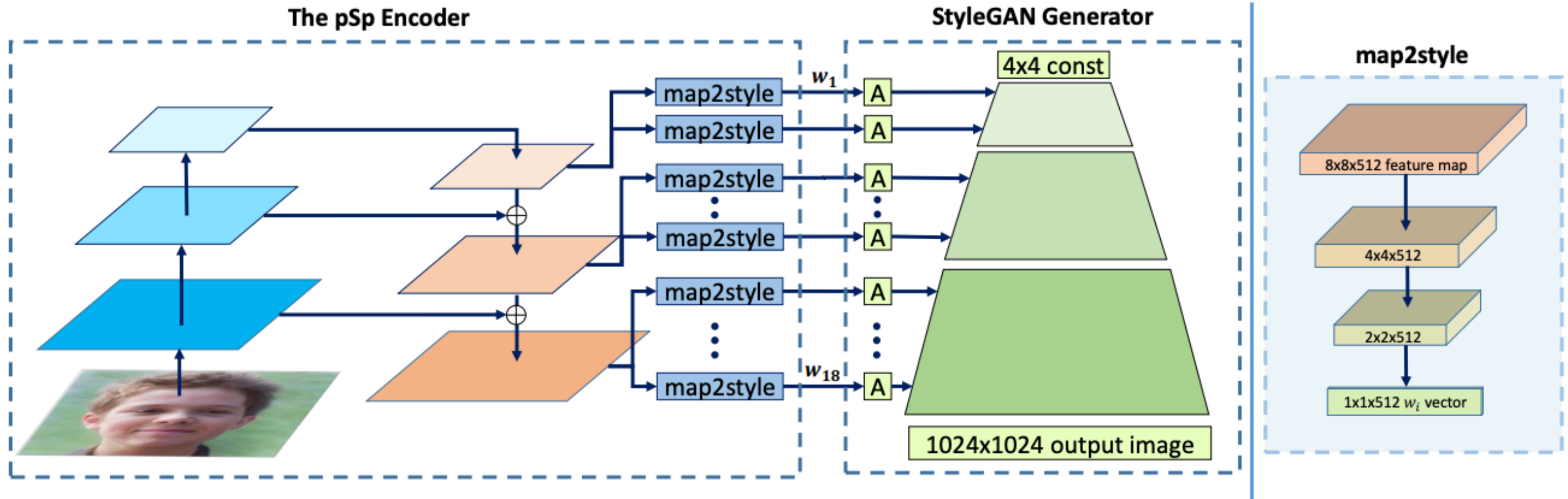
How to Improve GANs Projection

- Baseline: Optimizing the latent code

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

- Training an encoder $E(x)$. Advantages?
 - Faster inference
 - More reliable initialization
- Encoder design depends on
 - Generator architecture.
 - Which latent space: z , w , $w+$.
 - Pre-trained network weights.

Example: An StyleGAN Encoder



Encoding in Style: a StyleGAN Encoder for Image-to-Image Translation
[Richardson et al., CVPR 2021]

Example: An StyleGAN Encoder



Input

W

Naive $W+$

pSp

Debugging GANs Projection (HW5)

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

- What can go wrong?
 - Generator: G (cannot generate the image or too deep)
 - Reconstruction loss: L (not a good image distance)
 - Optimization method: SGD, ADAM (local minimum)
 - (1) use a more advanced solver: e.g., L-BFGS (Quasi-Newton)
 - (2) train an encoder to initialize the latent code. E(x)
- Debugging steps:
 - Reconstruct a generated image
 - Reconstruct a training set real image
 - Reconstruct a validation/test set real image
 - Reconstruct an in-the-wild image (e.g., Internet photo, camera roll)

Reconstruction \neq Editing



Interpolations between two images

Image Editing with GANs

- Step 1: Image Projection/Reconstruction

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

- Step 2: Manipulating the latent code

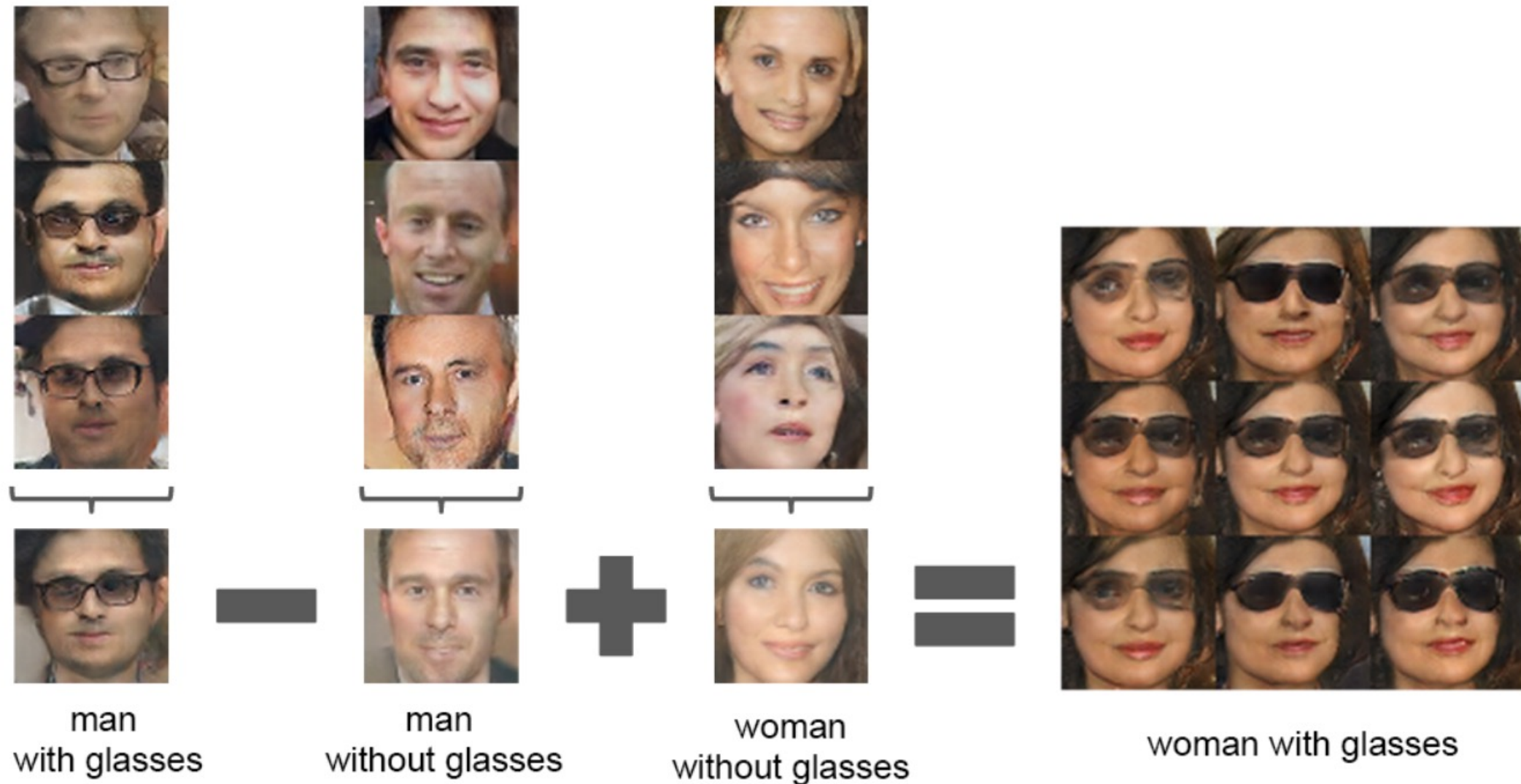
$$z_1 = z_0 + \Delta z$$

- Step 3: Generate the edited result

$$G(z_1)$$

Manipulating Latent code/layer
(computing directions offline)

Compute Δz



Step 1: annotate images (manually or via a pre-trained classifier)

Step 2: compute directions

Manipulating Latent code/layer (PCA directions)

GANSpace: Discovering PCA directions

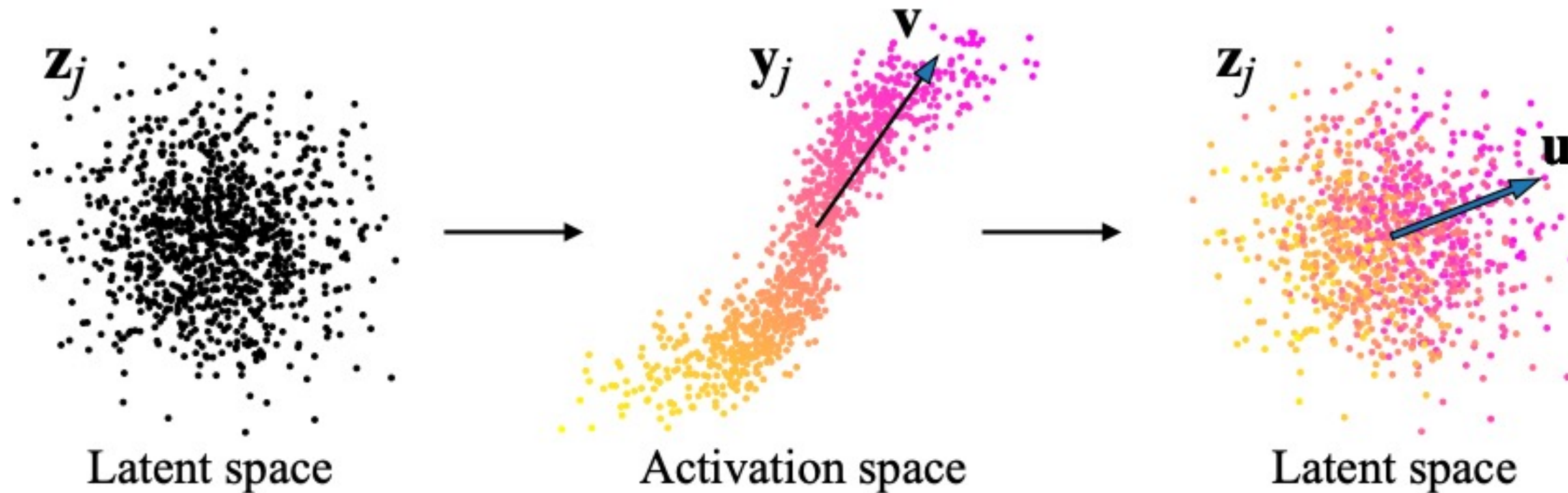


First find compute potential directions (PCA), then name them

GANSpace: Discovering PCA directions

z : latent codes. y : intermediate features.

v : PCA direction in feature space, u : PCA direction in latent space



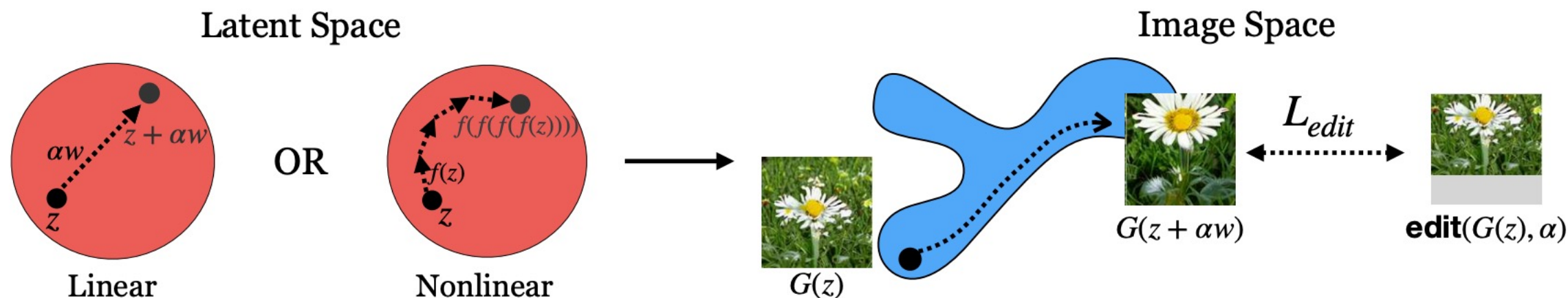
Also see “Editing in Style: Uncovering the Local Semantics of GANs”, Collins et al., CVPR 2020
“Closed-Form Factorization of Latent Semantics in GANs”, Shen and Zhou. CVPR 2021

GANSpace: Discovering PCA directions



Manipulating Latent code/layer (offline optimization)

Offline optimization



Given a pre-defined function **edit** and a pre-trained generator **G**

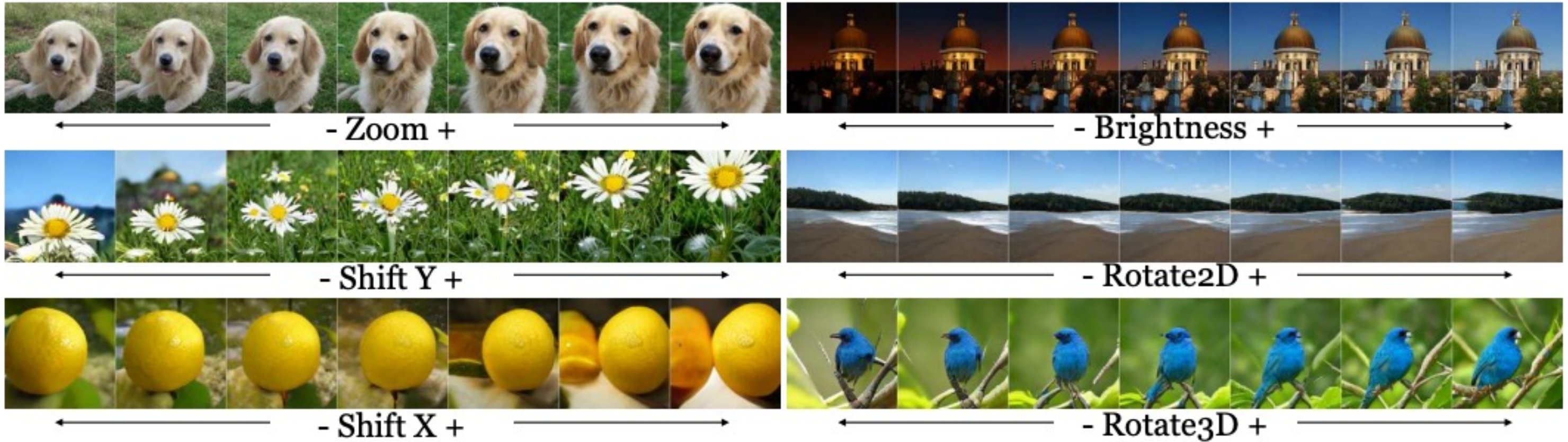
Linear case:
(w is a vector)

$$\arg \min_w \mathbb{E}_{z, \alpha} [\mathcal{L}(G(z + \alpha w), \text{edit}(G(z), \alpha))] \quad \nearrow \text{strength}$$

Non-linear case:
(f is a function)
apply it n times

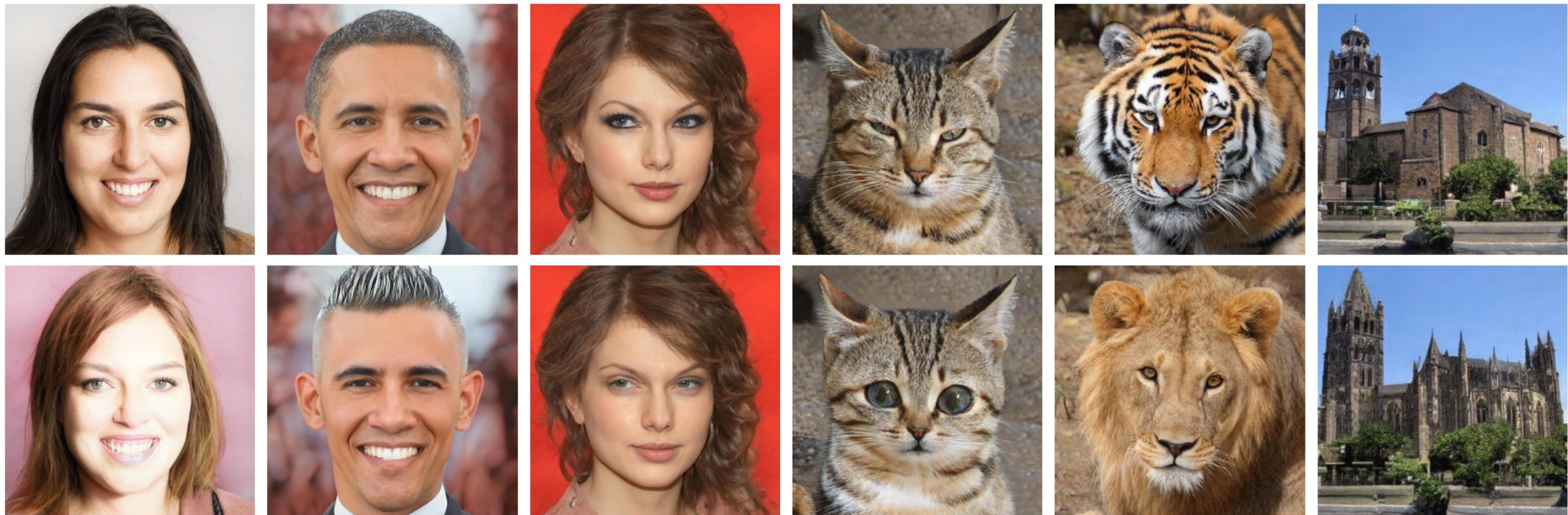
$$\arg \min_f \mathbb{E}_{z, n} [|||G(f^n(z)) - \text{edit}(G(z), n\epsilon)|||], \quad \nearrow \text{strength}$$

Offline optimization



Requirement: A known **edit** function
(e.g., shift, zoom, rotate)

CLIP-guided Directions



“Emma Stone”

“Mohawk hairstyle”

“Without makeup”

“Cute cat”

“Lion”

“Gothic church”

$$\arg \min_{w \in \mathcal{W}^+} D_{\text{CLIP}}(G(w), t) + \lambda_{\text{L2}} \|w - w_s\|_2 + \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(w)$$

$w \in \mathcal{W}^+$ Output is close to the text

Close to the original latent

Output is close to input

CLIP: Connecting Text and Images

FOOD101

guacamole (90.1%) Ranked 1 out of 101 labels



- ✓ a photo of **guacamole**, a type of food.
- ✗ a photo of **ceviche**, a type of food.
- ✗ a photo of **edamame**, a type of food.
- ✗ a photo of **tuna tartare**, a type of food.
- ✗ a photo of **hummus**, a type of food.

SUN397

television studio (90.2%) Ranked 1 out of 397



- ✓ a photo of a **television studio**.
- ✗ a photo of a **podium indoor**.
- ✗ a photo of a **conference room**.
- ✗ a photo of a **lecture room**.
- ✗ a photo of a **control room**.

YOUTUBE-BB

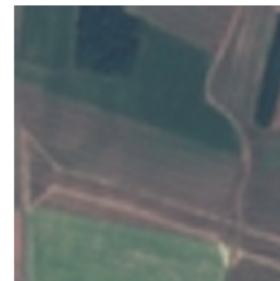
airplane, person (89.0%) Ranked 1 out of 23



- ✓ a photo of a **airplane**.
- ✗ a photo of a **bird**.
- ✗ a photo of a **bear**.
- ✗ a photo of a **giraffe**.
- ✗ a photo of a **car**.

EUROSAT

annual crop land (12.9%) Ranked 4 out of 10



- ✗ a centered satellite photo of **permanent crop land**.
- ✗ a centered satellite photo of **pasture land**.
- ✗ a centered satellite photo of **highway or road**.
- ✓ a centered satellite photo of **annual crop land**.
- ✗ a centered satellite photo of **brushland or shrubland**.

Input: an image and a caption.

Output: similarity between the text embedding and the image embedding

CLIP-guided Directions



$$\arg \min_{w \in \mathcal{W}} D_{\text{CLIP}}(G(w), t) + \lambda_{\text{L2}} \|w - w_s\|_2 + \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(w)$$

$w \in \mathcal{W} +$ Output is close to the text

Close to the original latent

Output is close to input

Manipulating network weights

Image Editing with GANs

- Step 1: Image Projection/Reconstruction

$$z_0, \theta_0 = \arg \min_{z, \theta} \mathcal{L}(G(z; \theta), x)$$

- Step 2: Manipulating the network weights

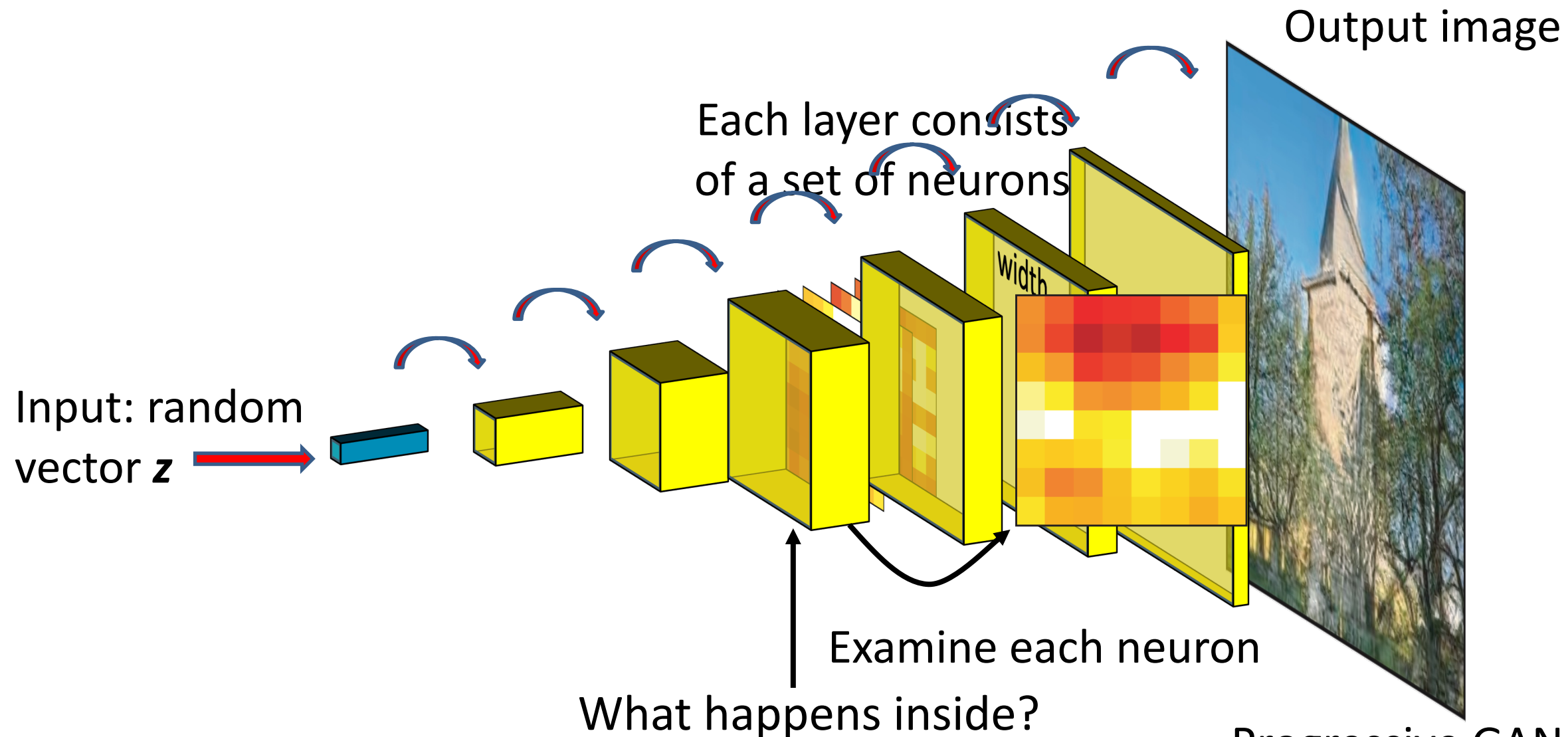
$$\theta_1 = \theta_0 + \Delta\theta$$

- Step 3: Generate the edited result

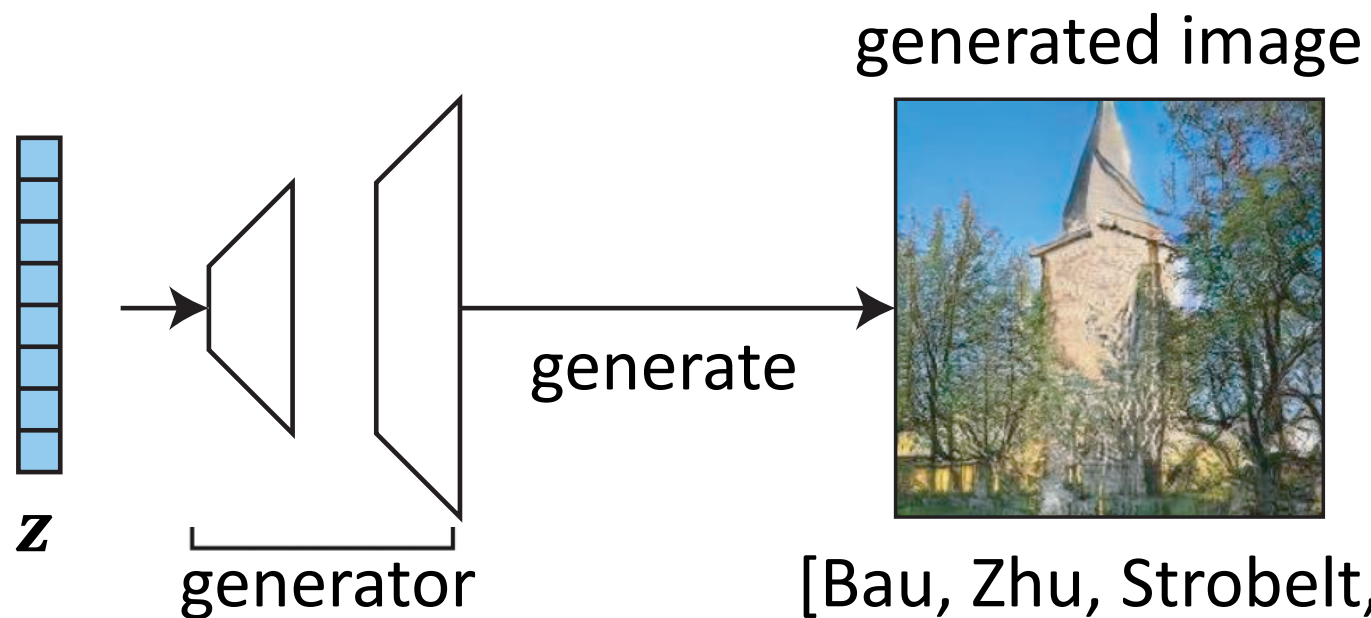
$$G(z_0; \theta_1)$$

Understanding a Generator

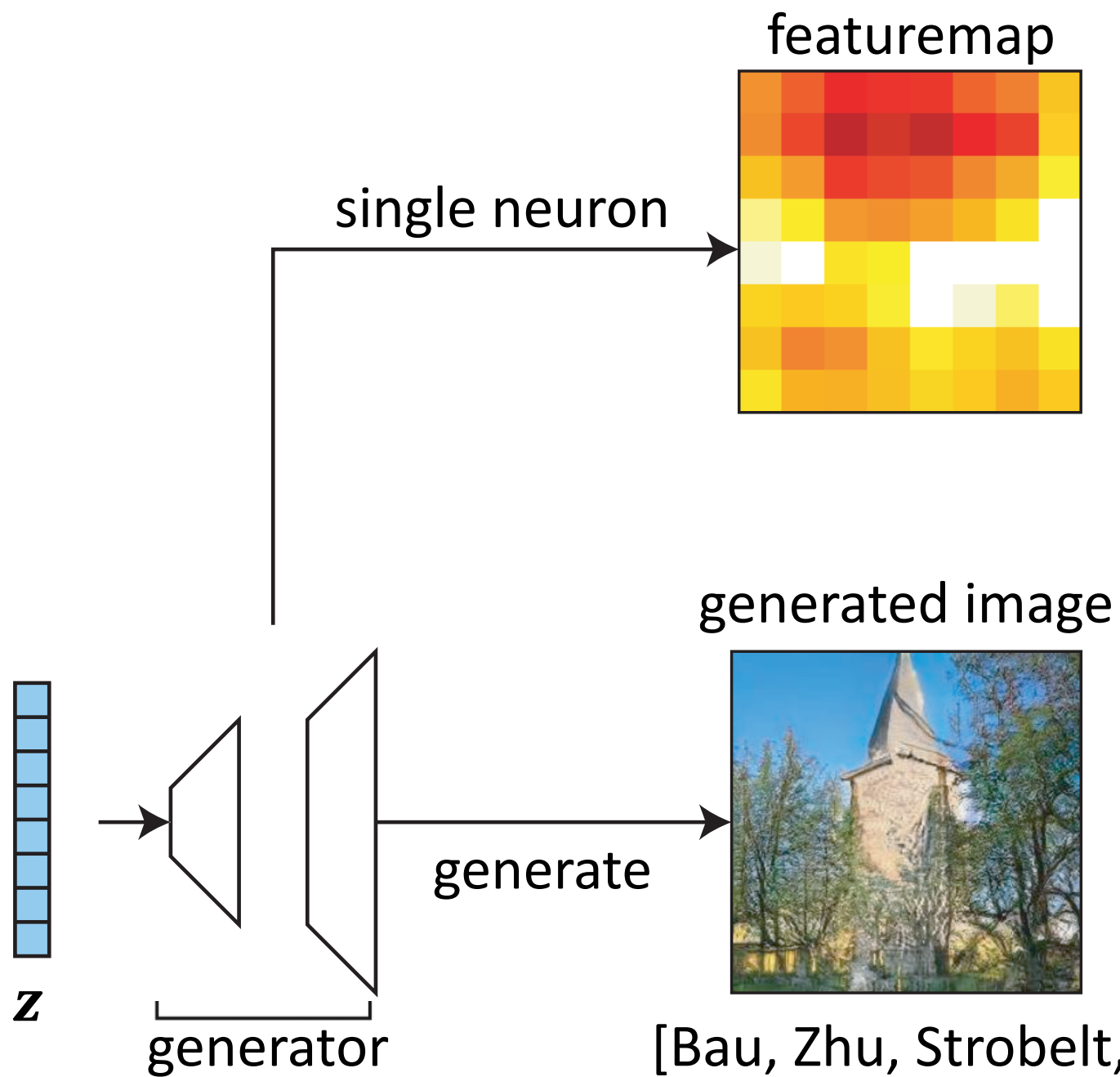
Each step:
Increases spatial resolution



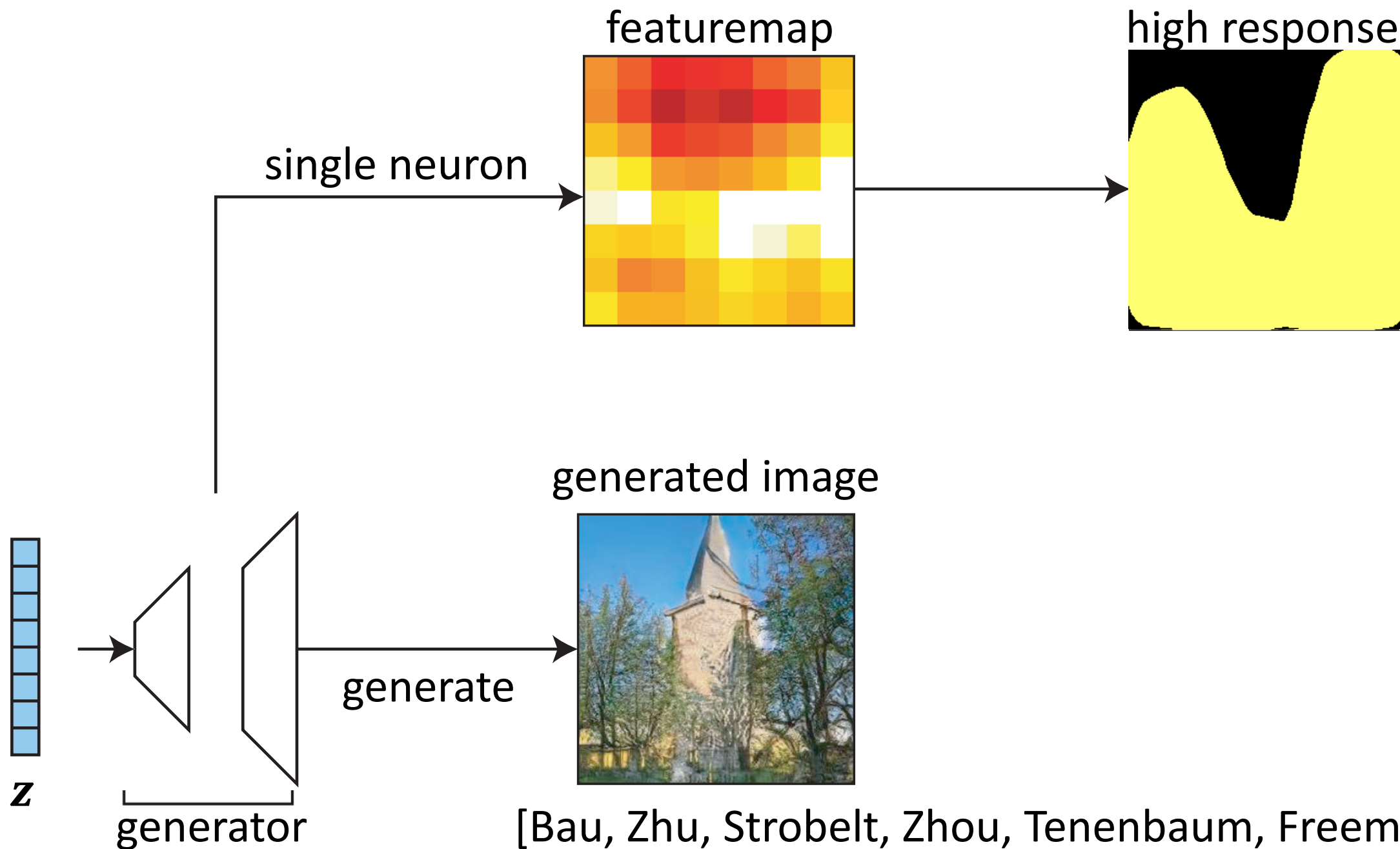
Which neurons **correlate** to an object class?



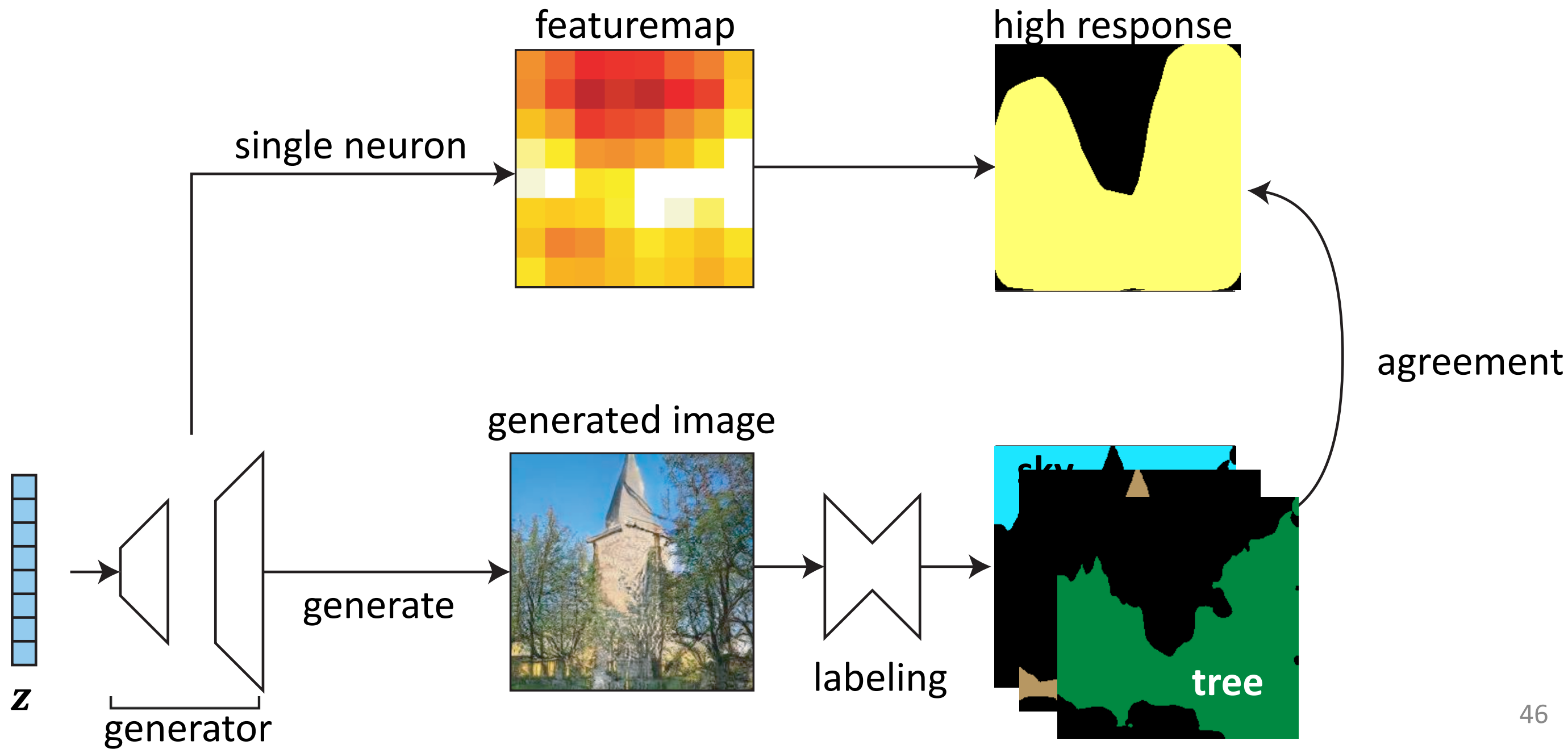
Which neurons **correlate** to an object class?



Which neurons correlate to an object class?



Which neurons **correlate** to an object class?



Which neurons correlate to an object class?

Church samples



Tree
Neuron

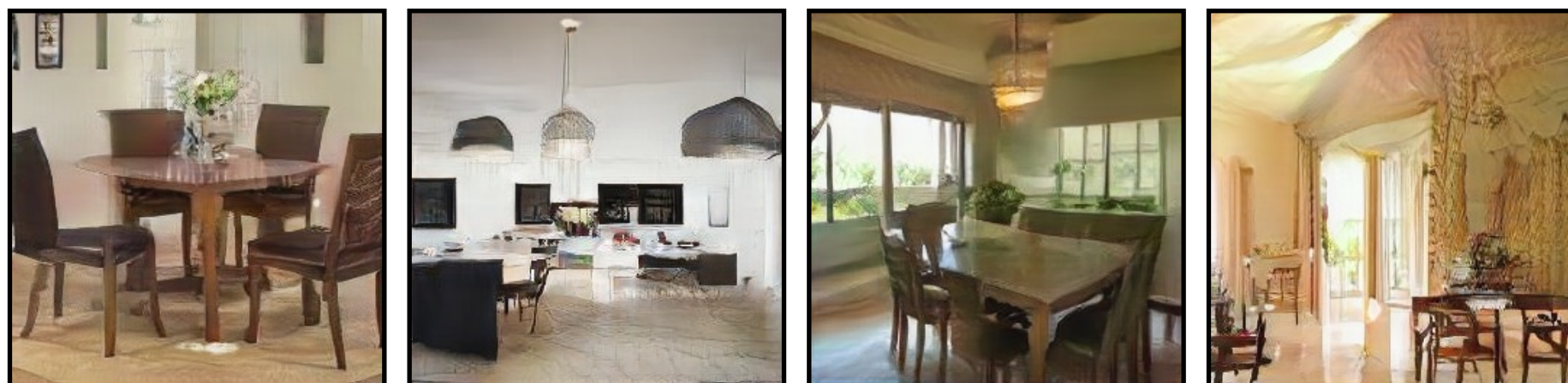


Dome
Neuron



Which neurons correlate to an object class?

Dining room samples



252 out of 512 neurons are correlated to objects, part, and materials

Window
Neuron

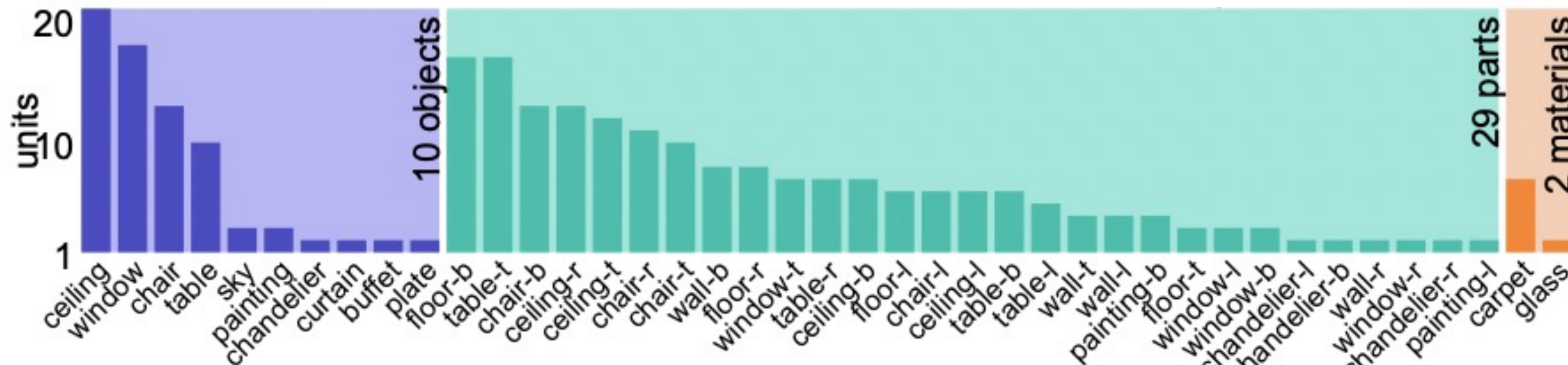
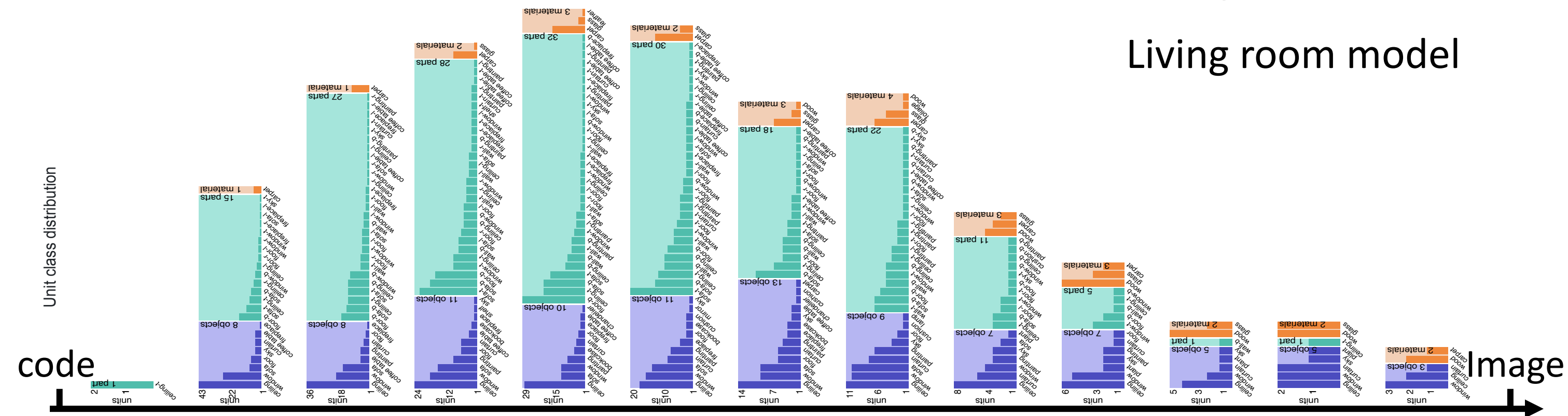


Table
Neuron

Which neurons correlate to an object class?

Living room model



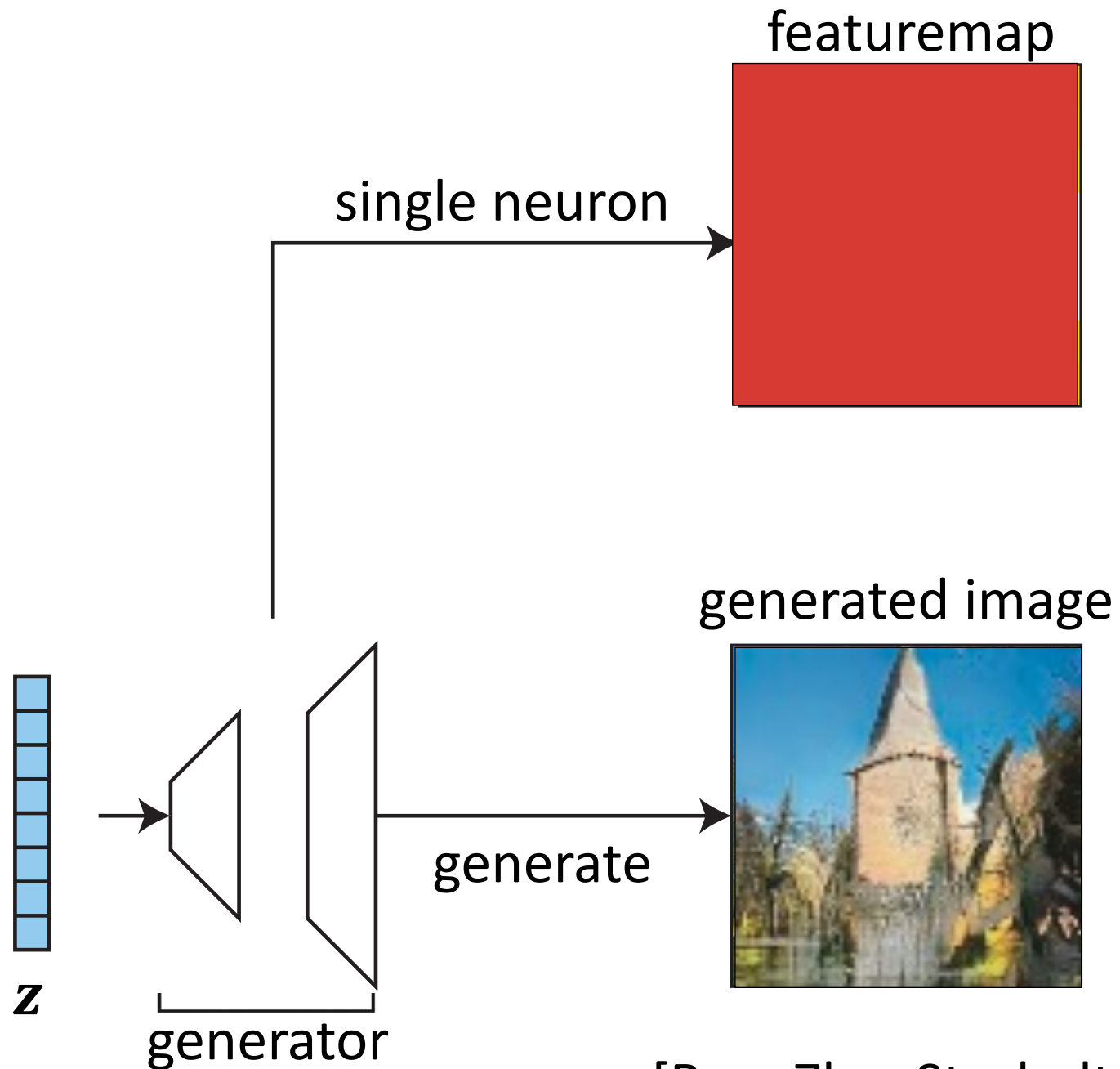
Layout

Object and parts

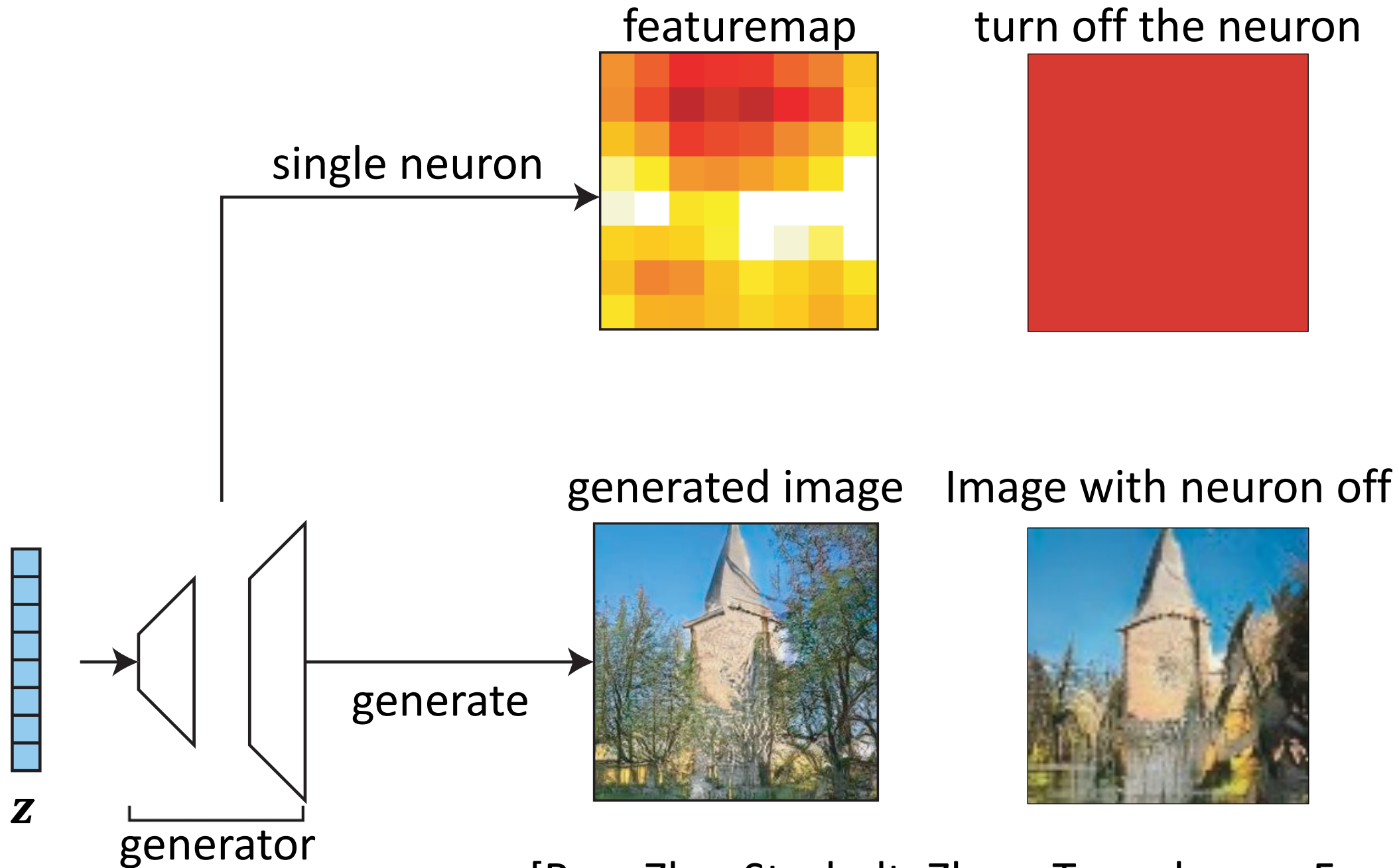
Edges, textures, local structure



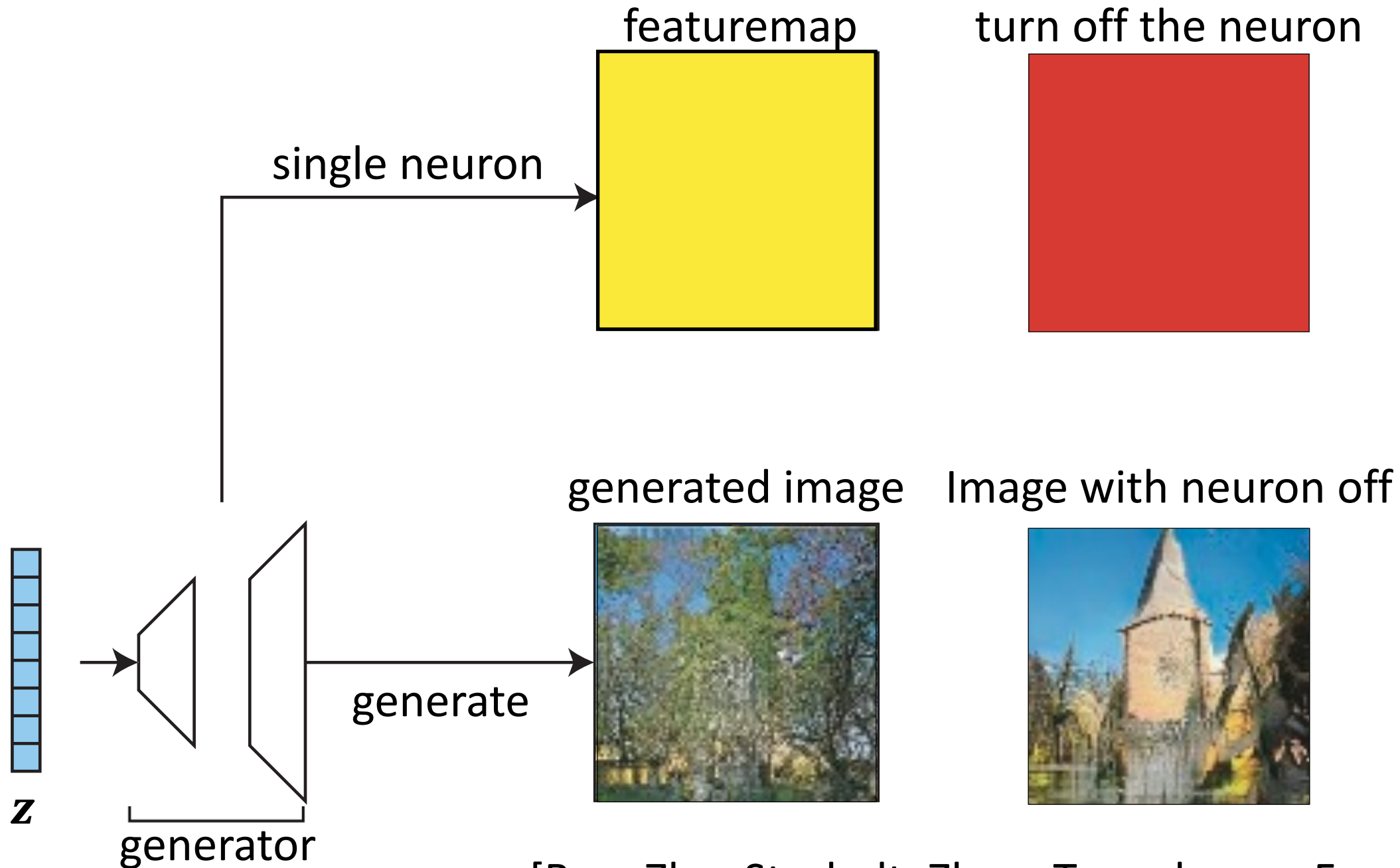
Which neurons cause an object class?



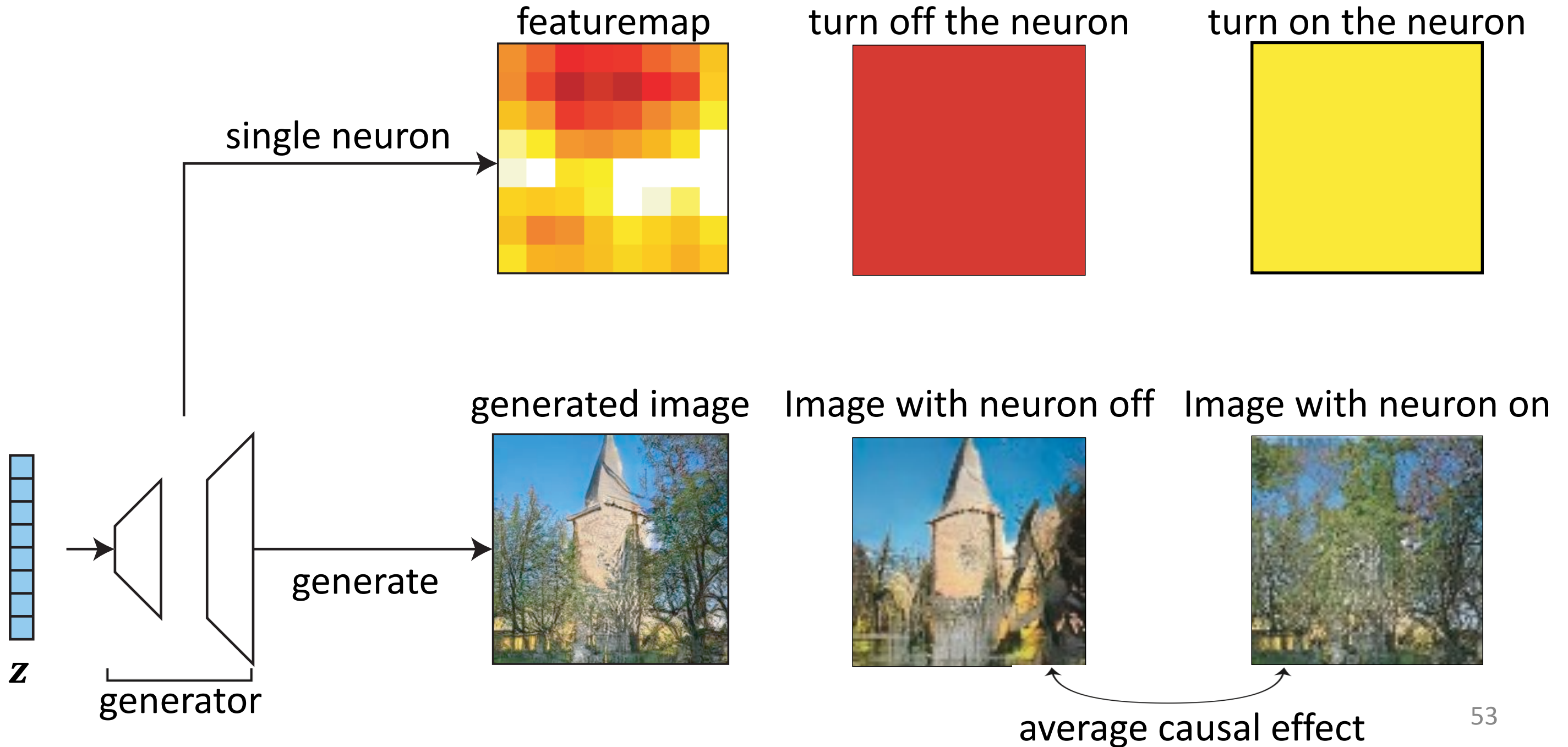
Which neurons cause an object class?



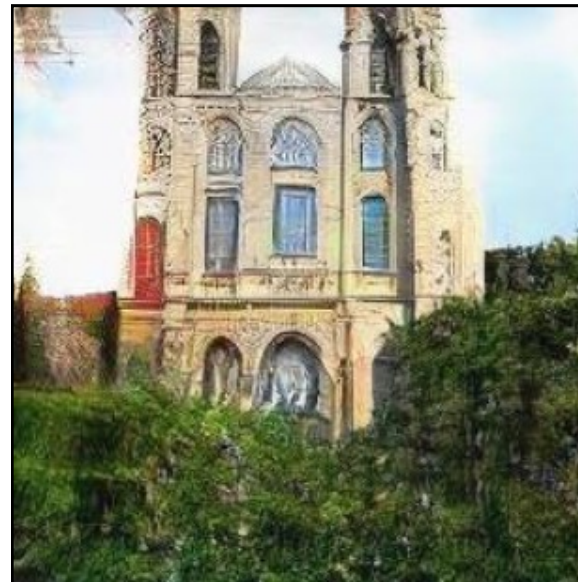
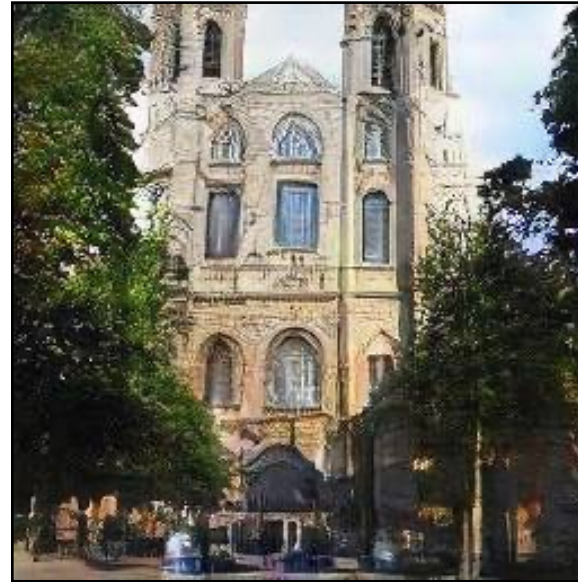
Which neurons cause an object class?



Which neurons cause an object class?



Which neurons cause an object class?



Interactive Painting

Select a feature brush & strength and enjoy painting:

tree

grass

door

sky

cloud

brick

dome

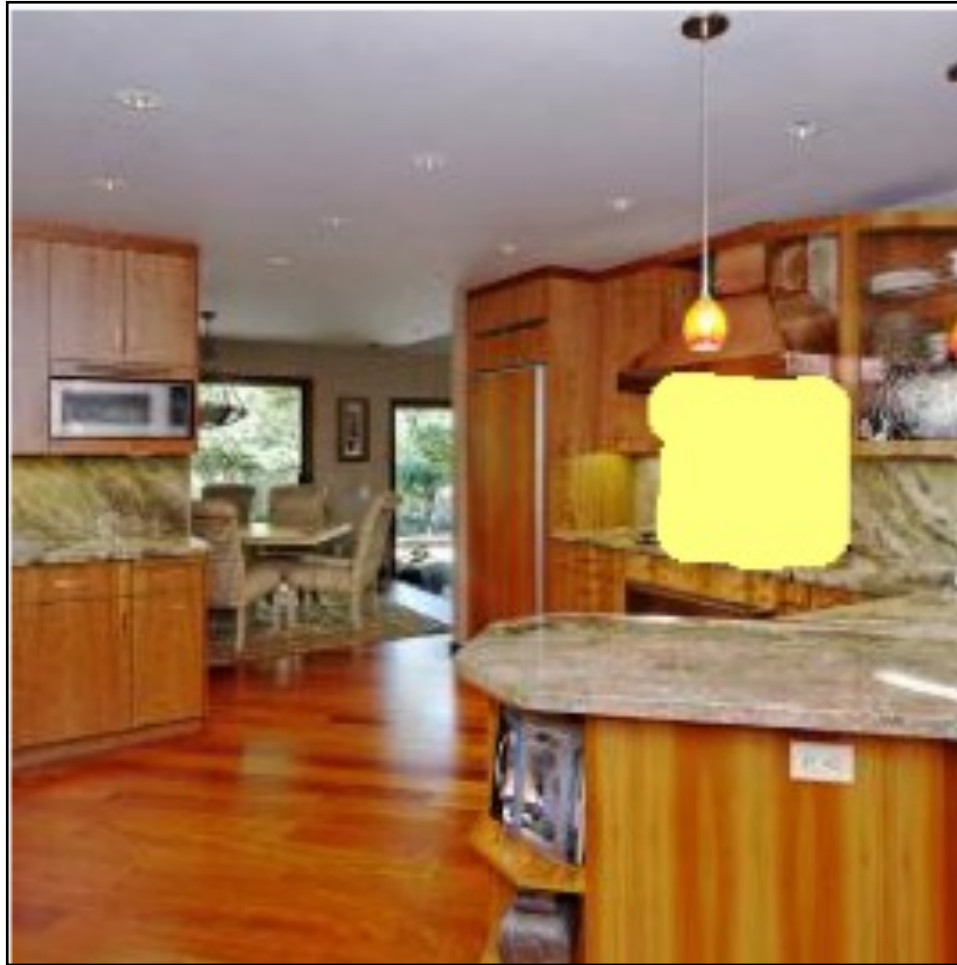


Online Demo

<http://bit.ly/ganpaint>



Manipulating a Real Photo



Original image + edits



Editing with \hat{z}

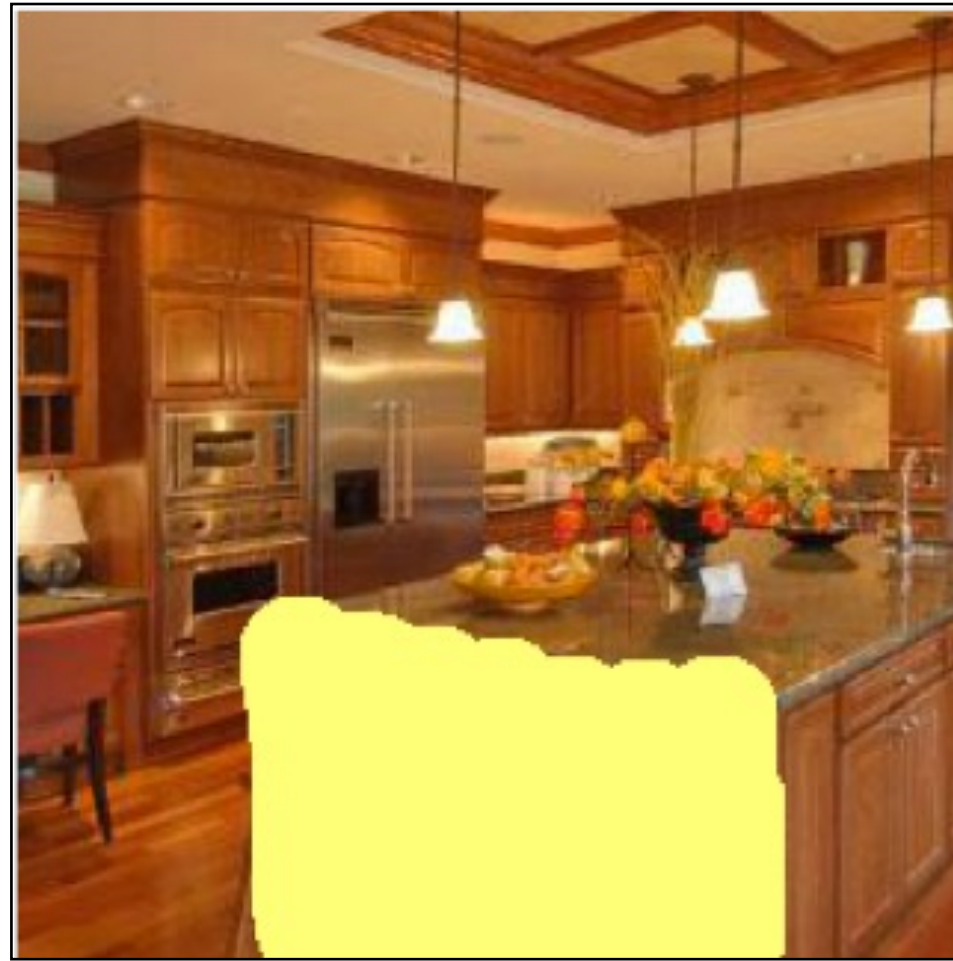


Editing with \hat{z} and $\hat{\theta}$

Manipulating a Real Photo



Input image



Remove chairs

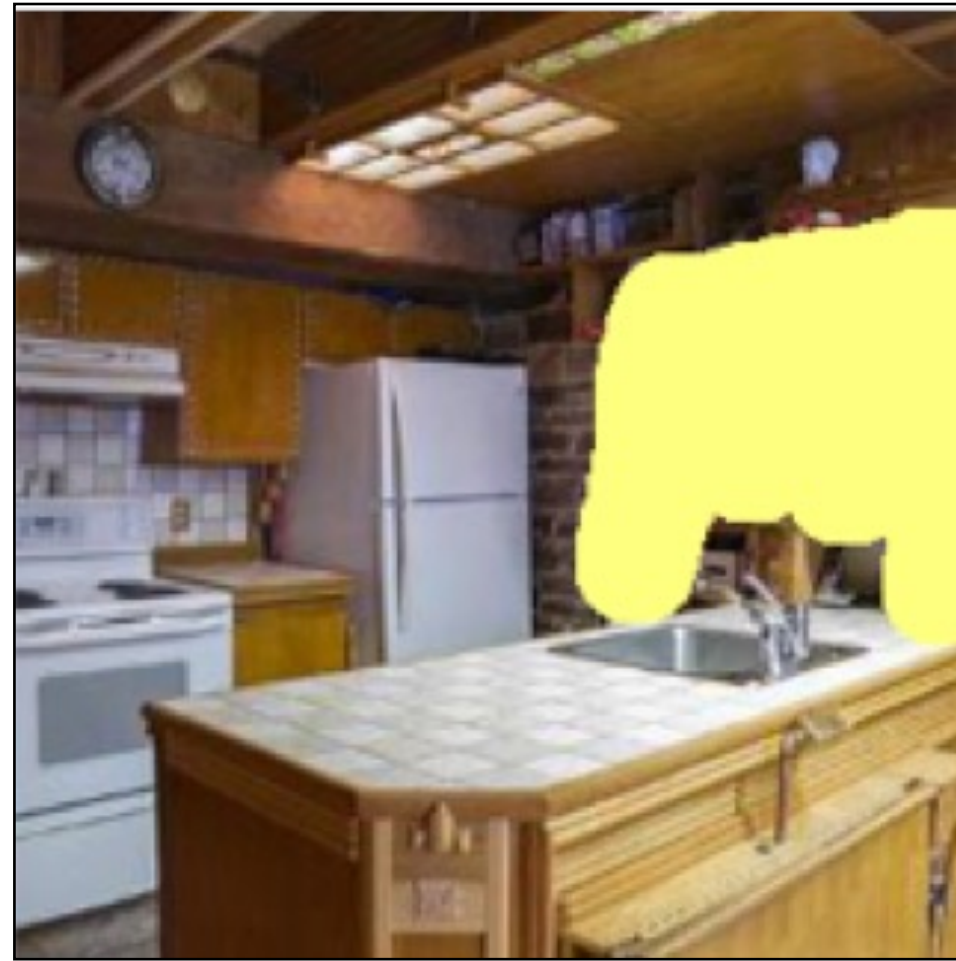


Output result

Manipulating a Real Photo



Input image



Add windows

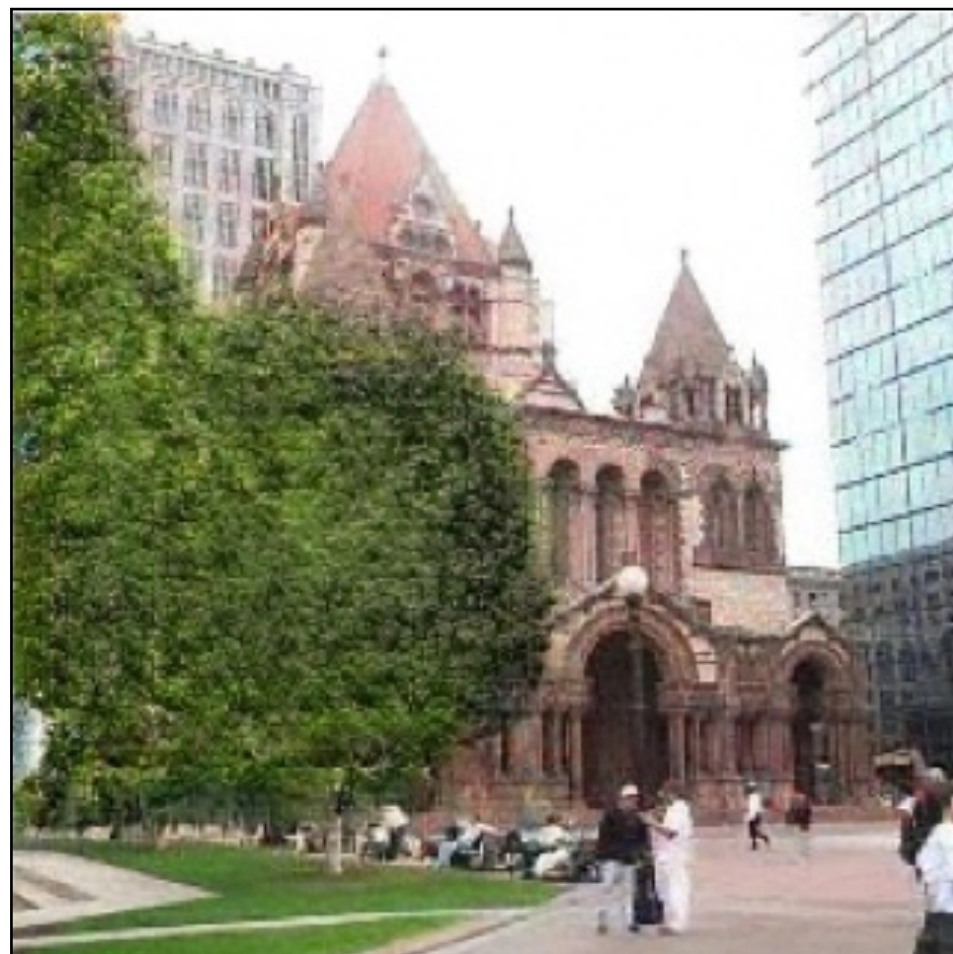


Output result

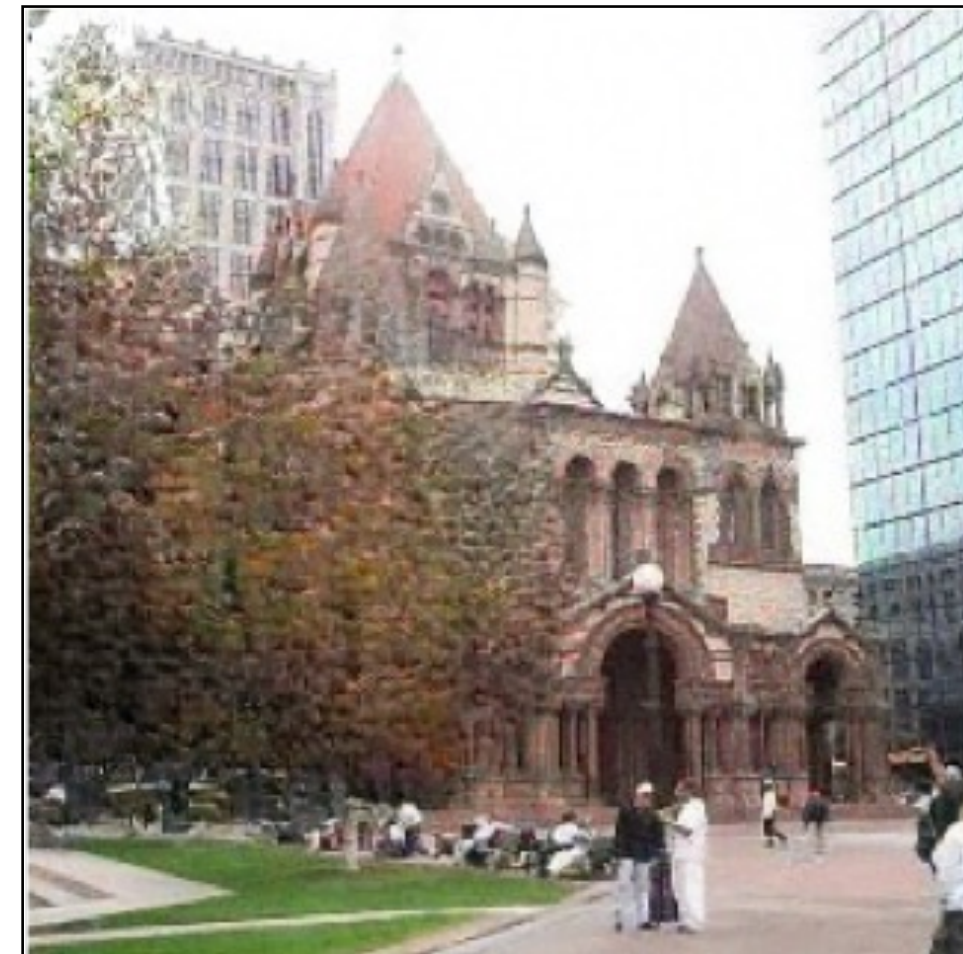
Manipulating a Real Photo via GAN Dissection



Input image



Restyle trees for spring



Restyle trees for autumn

Upload your image:

Choose File No file chosen

Draw:



tree

grass

door

dome

sky

cloud



low

med

high



undo reset

Optimization with Text-to-Image Diffusion Models

Text-to-image isn't perfect...

Stable
Diffusion



Photo of a [moongate](#)

Text-to-image isn't perfect...

Stable
Diffusion



Photo of a [moongate](#)

Customization

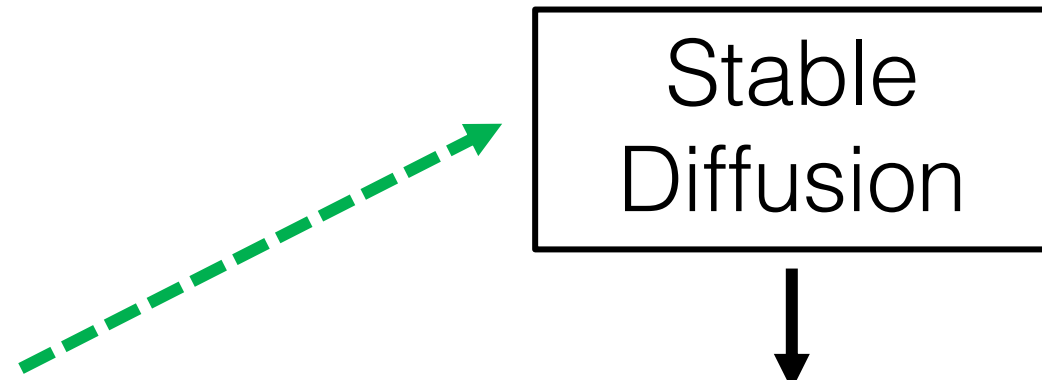
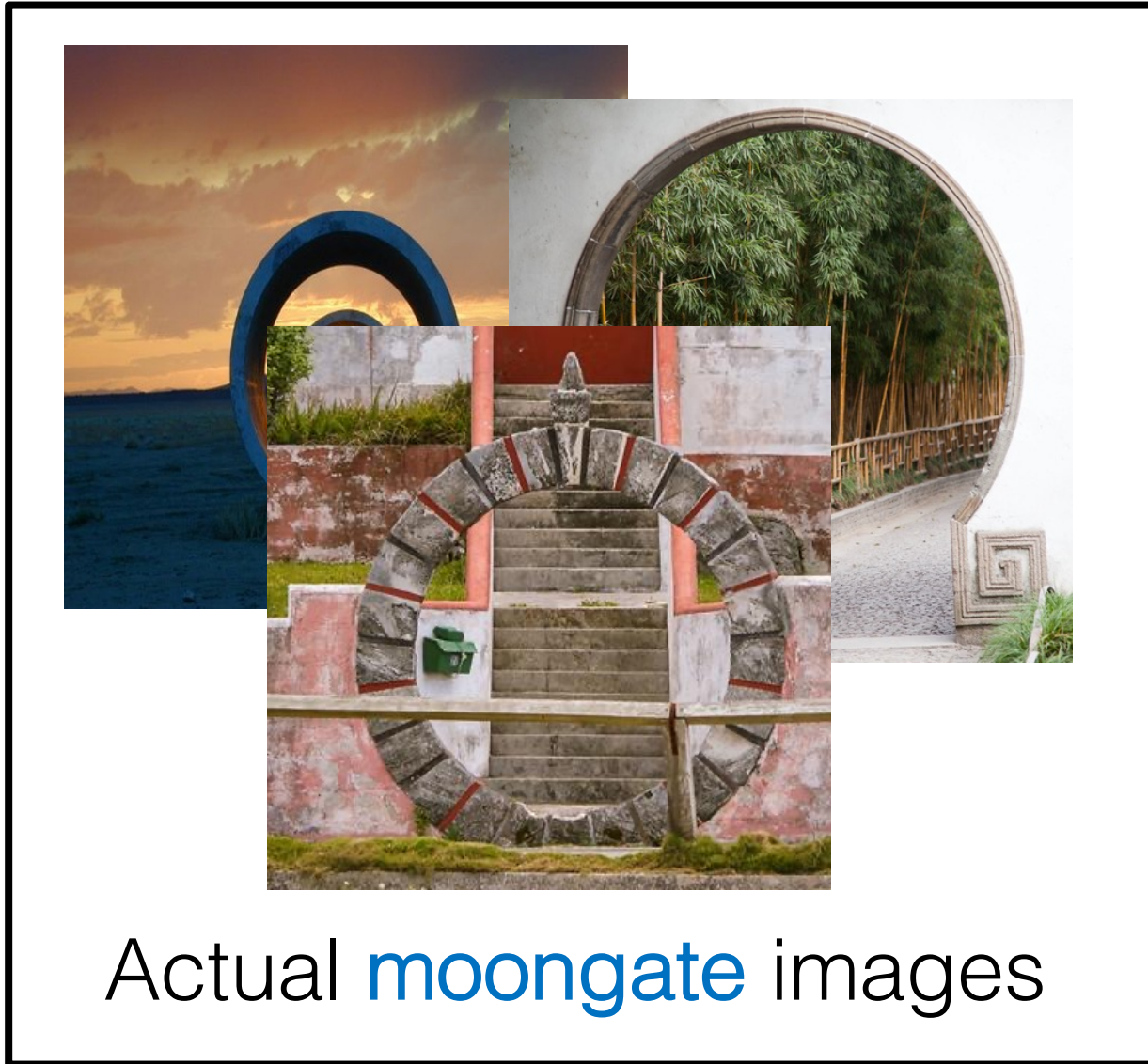


Photo of a **moongate**

Customization

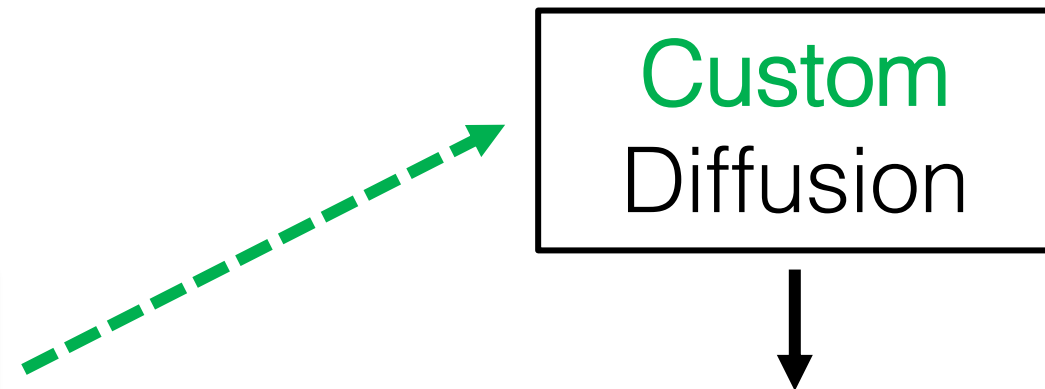
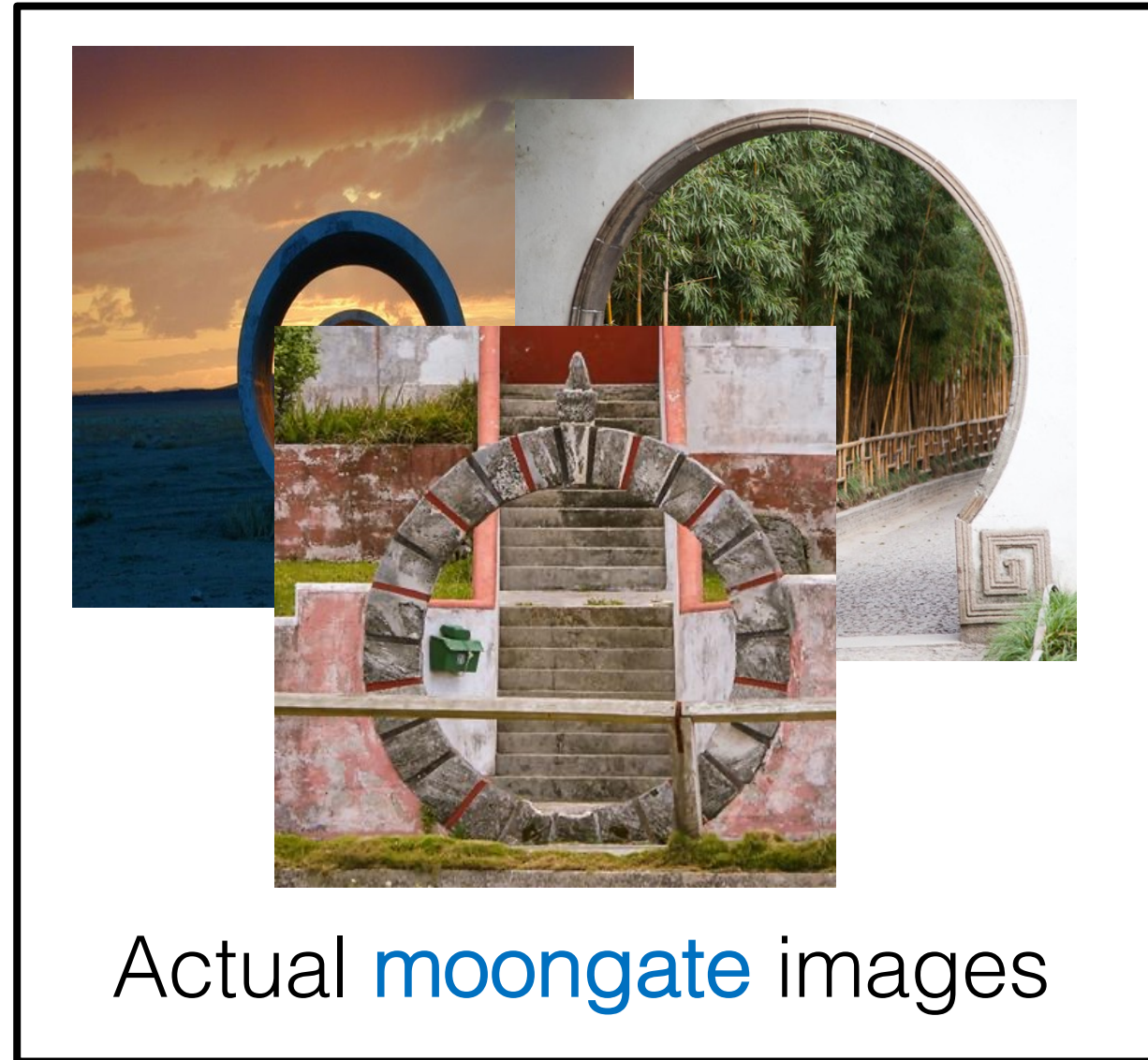


Photo of a **moongate**

Customization

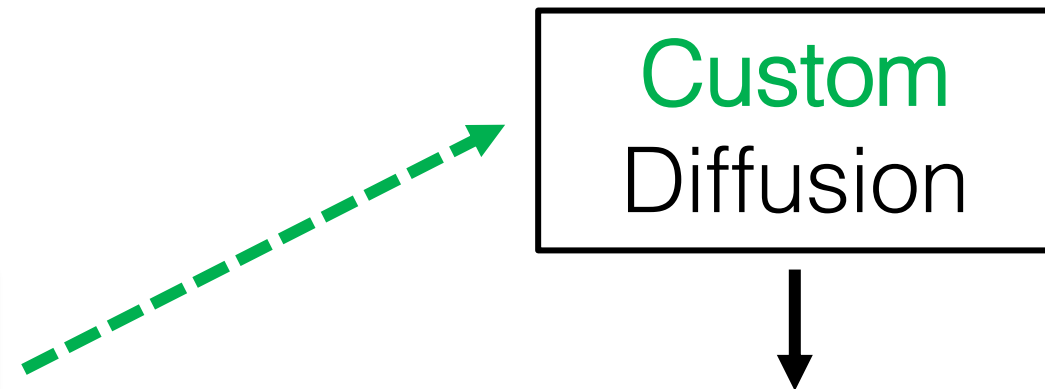
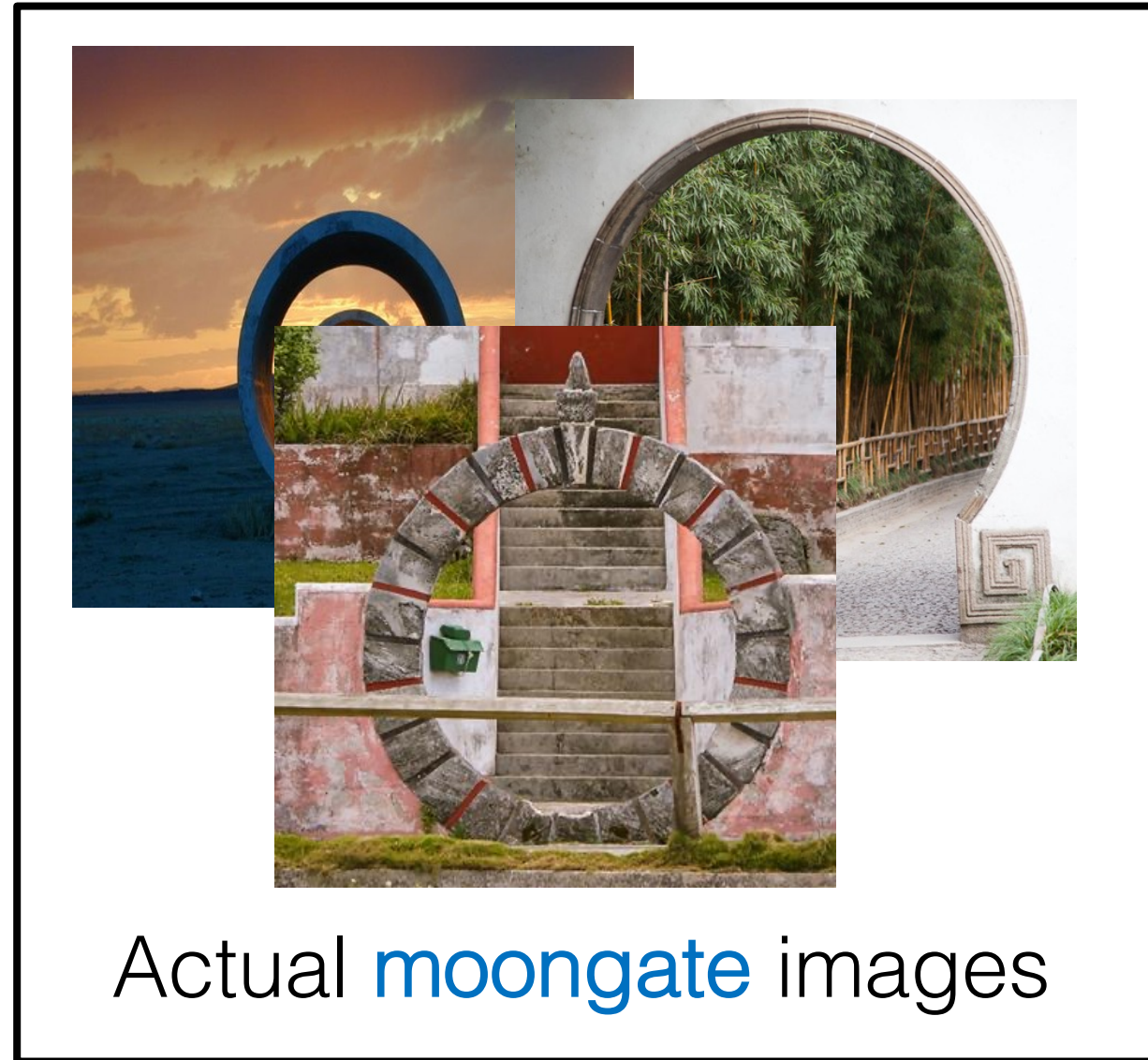
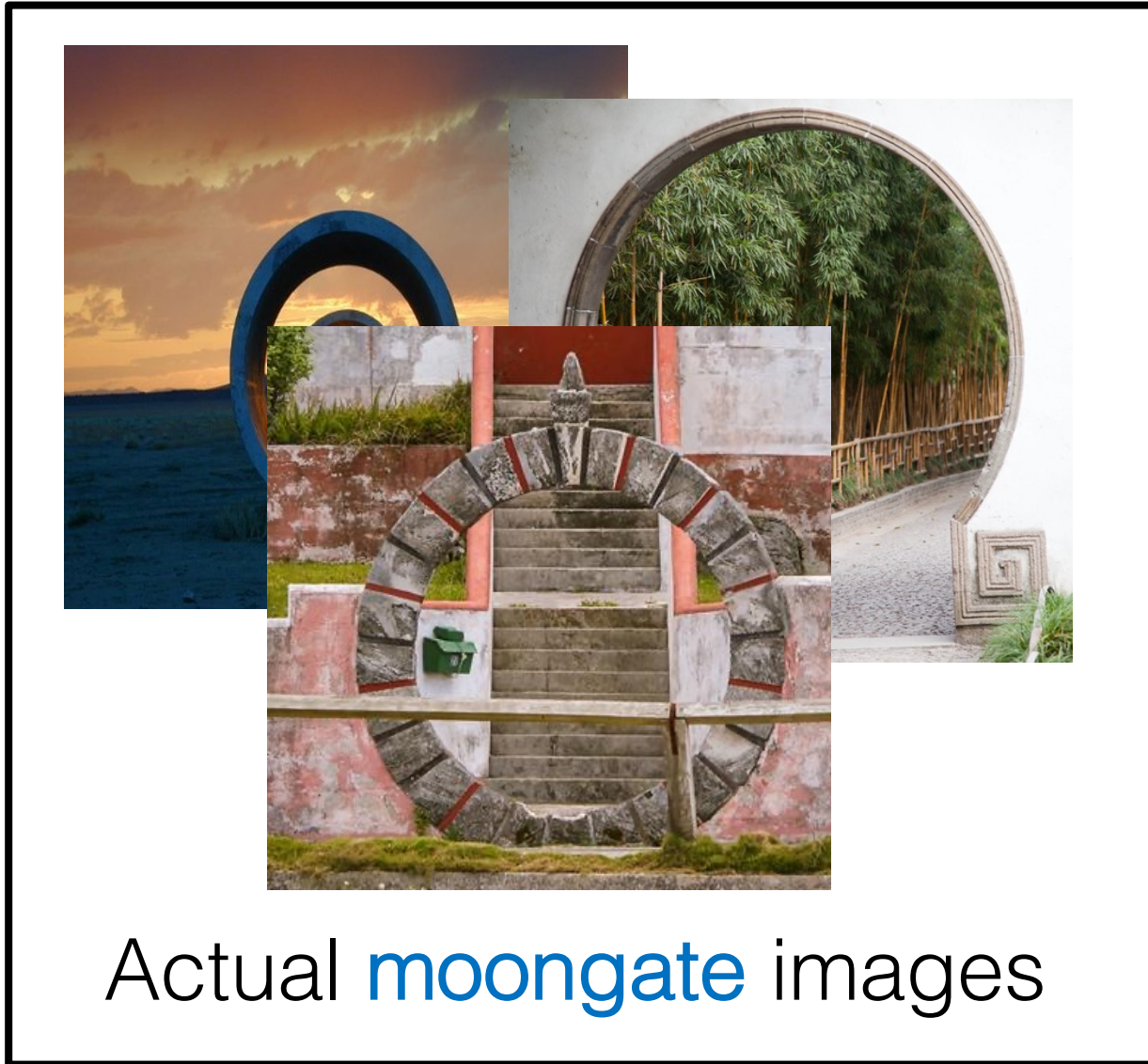


Photo of a **moongate**

Unseen contexts

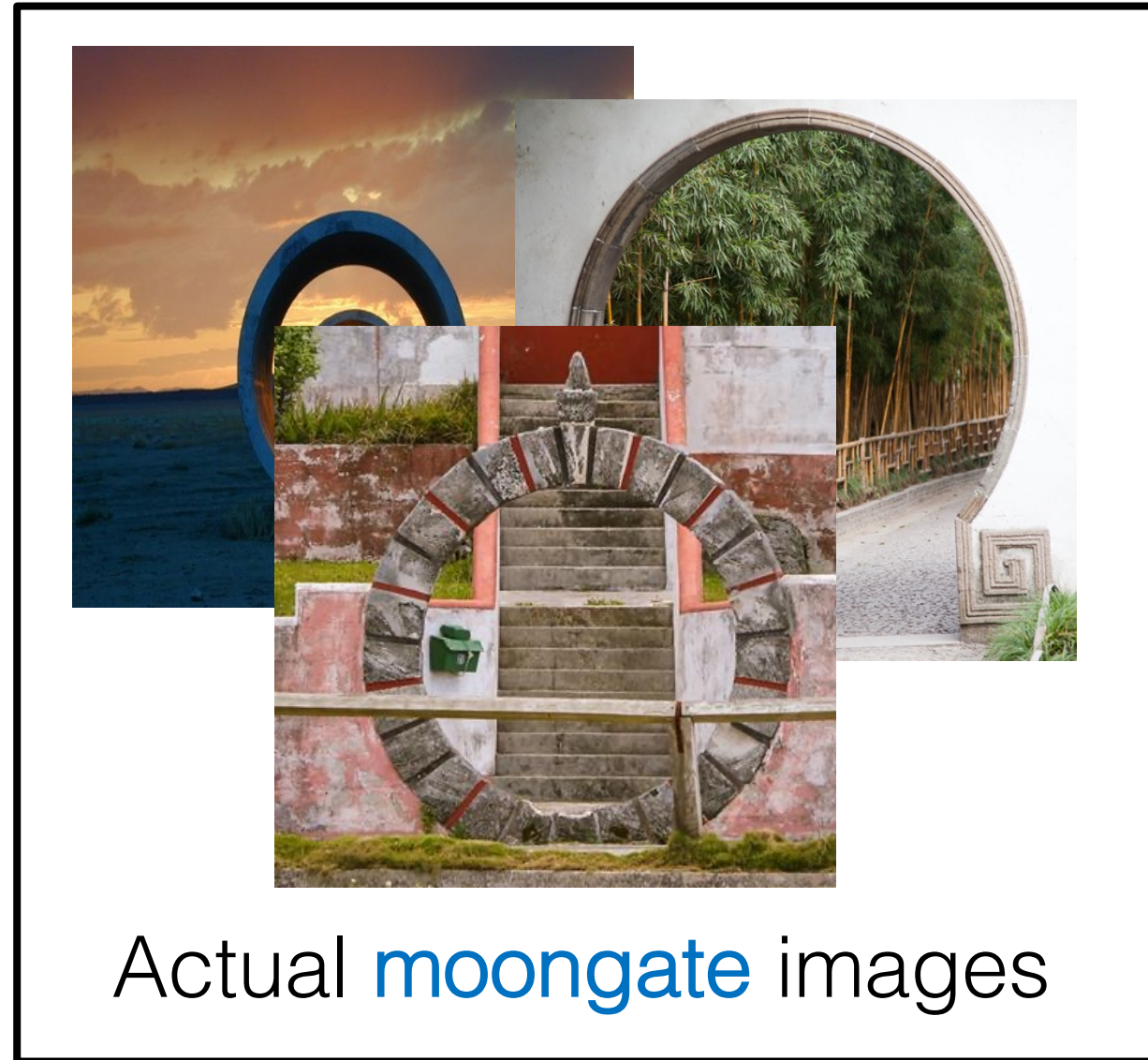


Custom
Diffusion



Moongate in the middle of highway

Unseen contexts

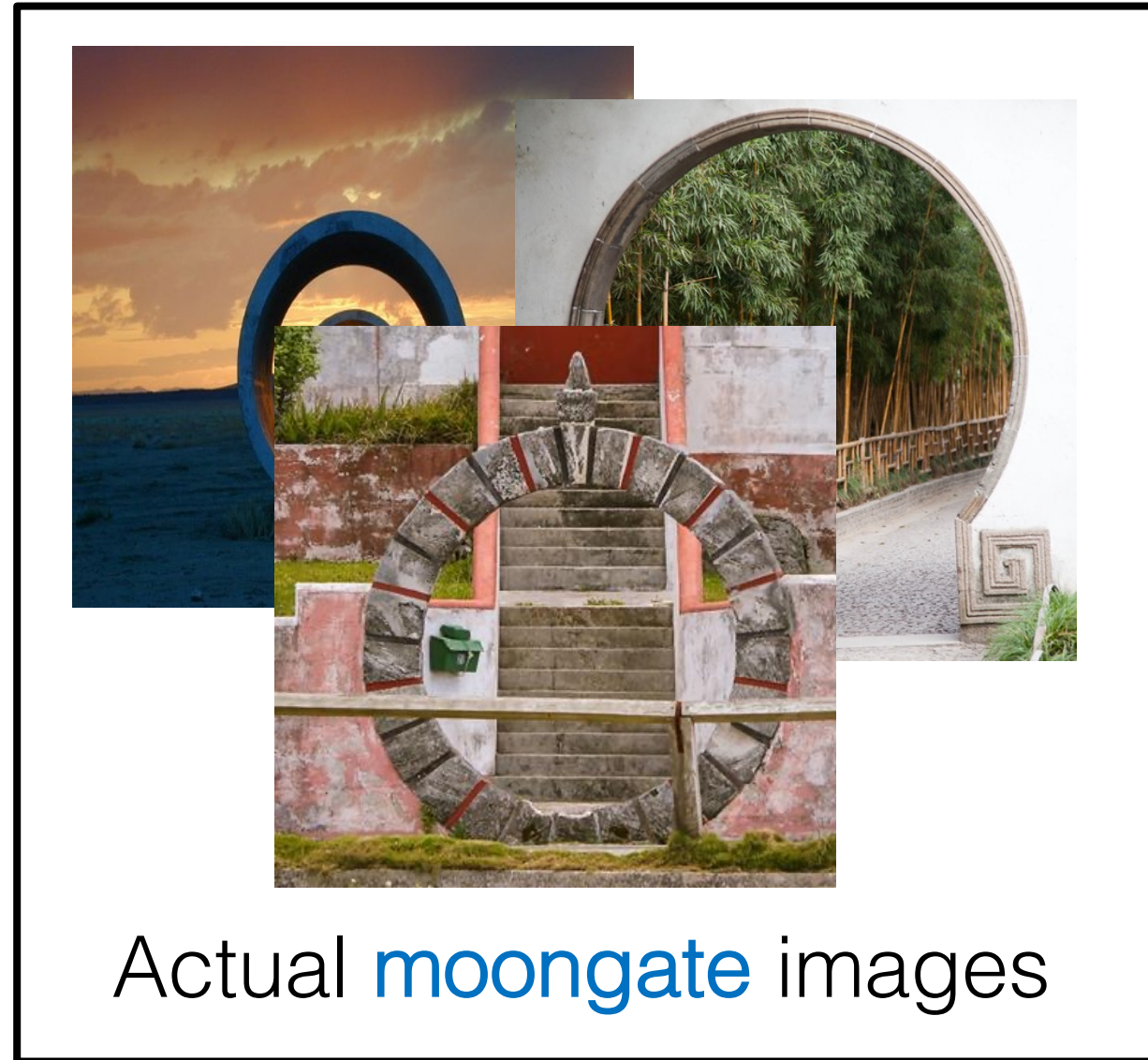


Custom
Diffusion

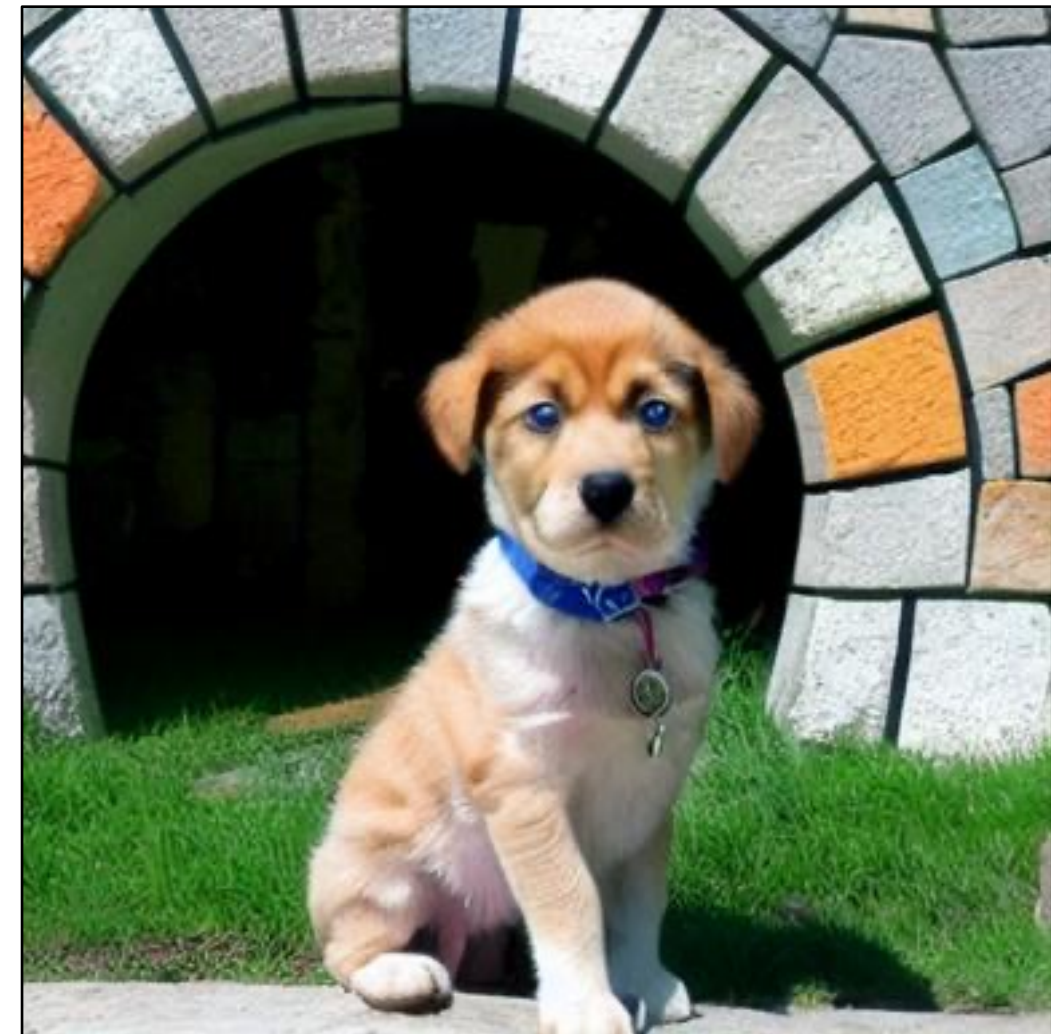


Moongate in snowy ice

Unseen contexts



Custom
Diffusion

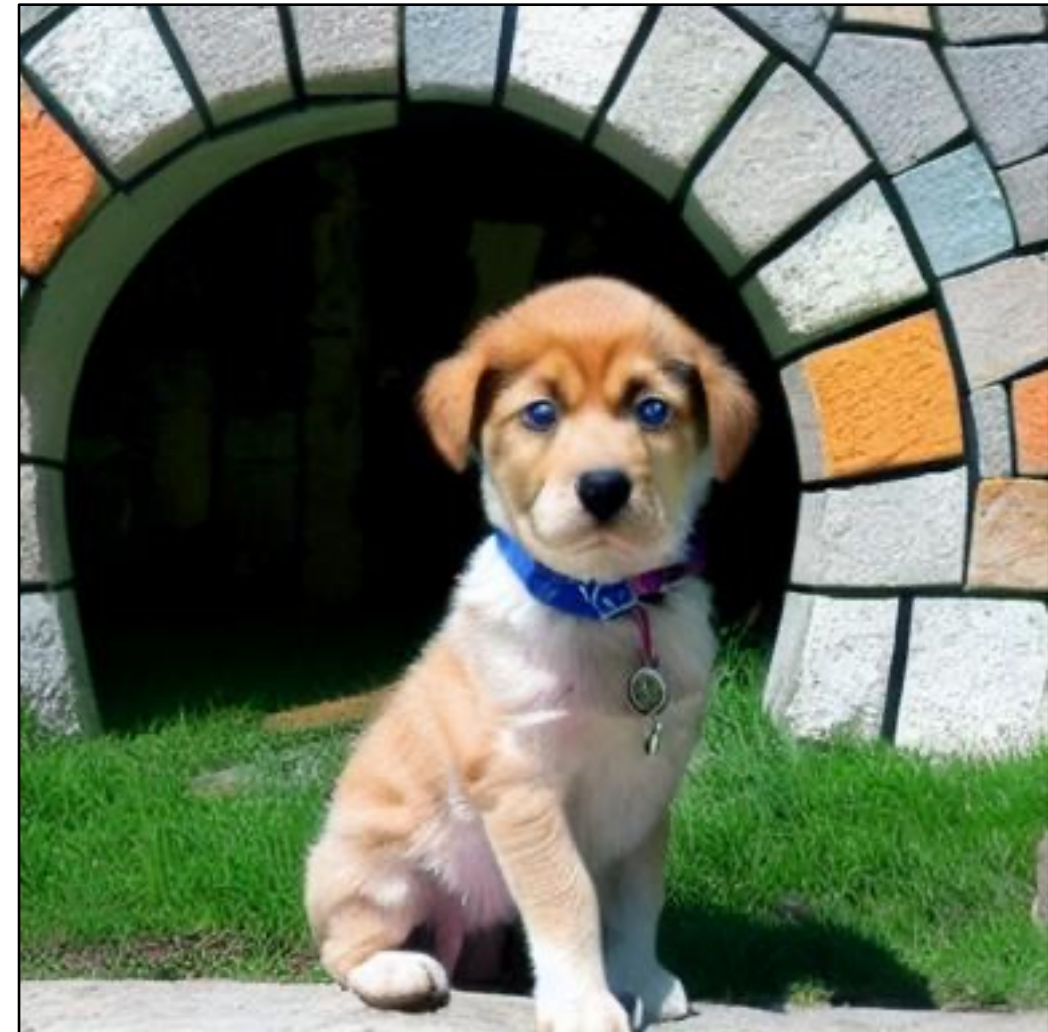


Multiple concepts



My **dog**, Stark

Custom
Diffusion



A puppy in front of **Moongate**₇₀

Multiple concepts



My **dog**, Stark

Custom
Diffusion



V* **dog** wearing sunglasses in front of **moongate**

Textual Inversion



Textual Inversion

Input samples $\xrightarrow{\text{invert}}$ " S_* "



→



"An oil painting of S_* "



"App icon of S_* "



"Elmo sitting in the same pose as S_* "



"Crochet S_* "

Input samples $\xrightarrow{\text{invert}}$ " S_* "



→



"Painting of two S_* fishing on a boat"



"A S_* backpack"

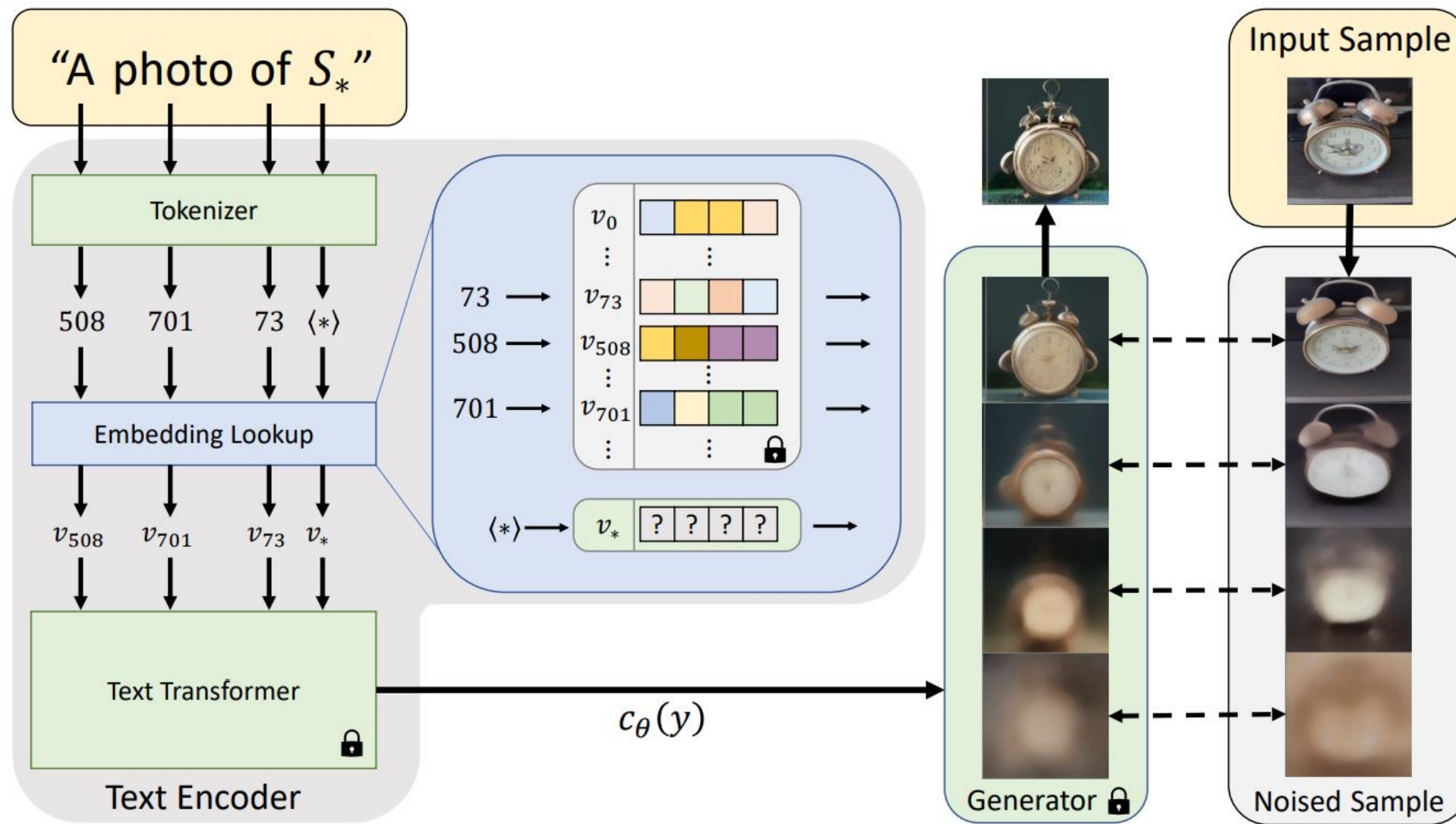


"Banksy art of S_* "



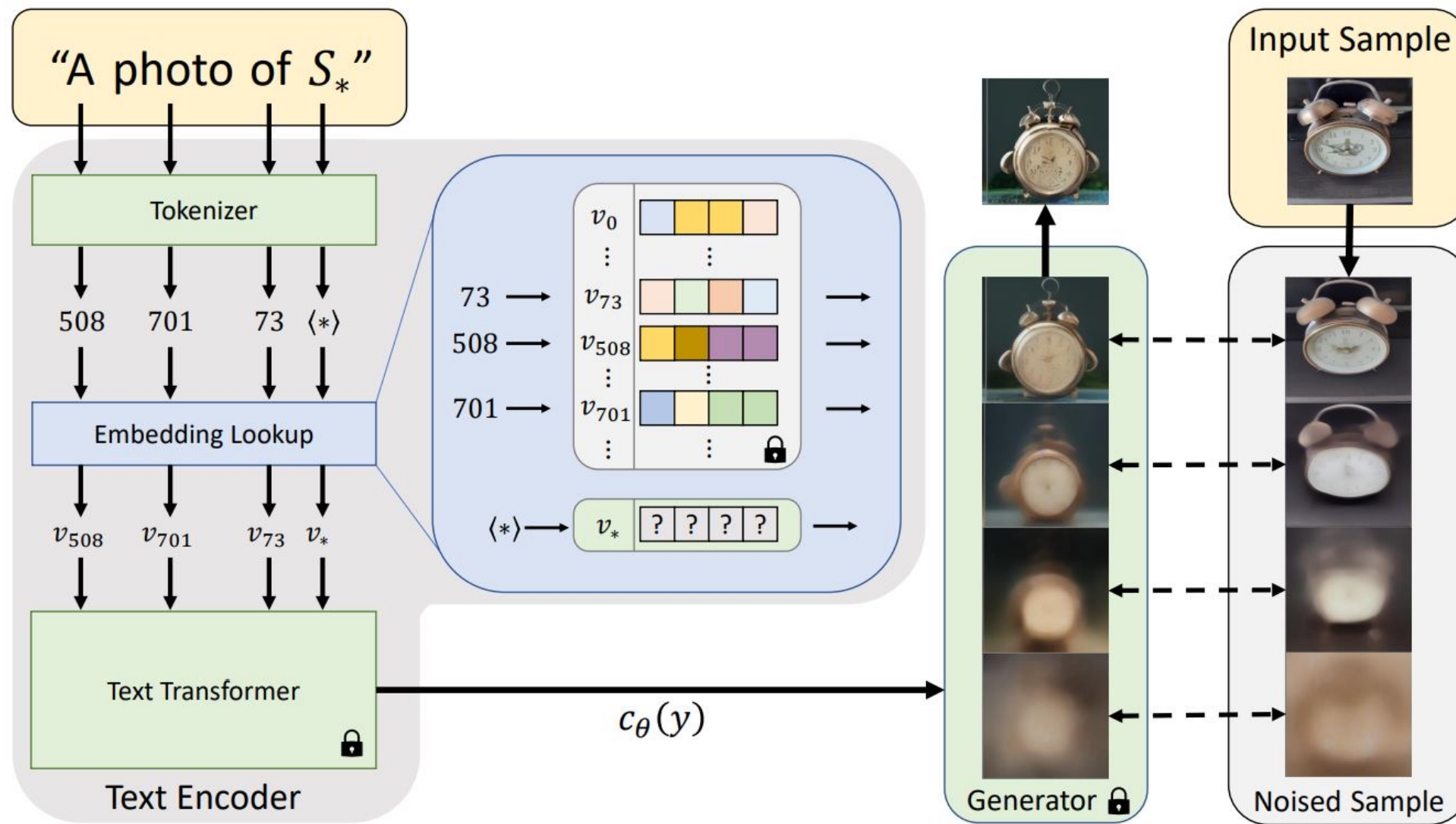
"A S_* themed lunchbox"

Textual Inversion



$$v_* = \arg \min_v \mathbb{E}_{z \sim \mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t, c_\theta(y))\|_2^2 \right]$$

Textual Inversion



$$v_* = \arg \min_v \mathbb{E}_{z \sim \mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t, c_\theta(y))\|_2^2 \right]$$

Works well for artistic styles



Input samples



“The streets of Paris
in the style of S_* ”



“Adorable corgi
in the style of S_* ”



“Painting of a black hole
in the style of S_* ”



“Times square
in the style of S_* ”

Compositional Ability



“Photo of *S_{clock}*
in the style of *S_{cat}*”

“Photo of *S_{clock}*
in the style of *S_{craft}*”

“Photo of *S_{cat}*
in the style of *S_{craft}*”

Only works for style

Reconstruction quality



Target images



S^* cat swimming in a pool

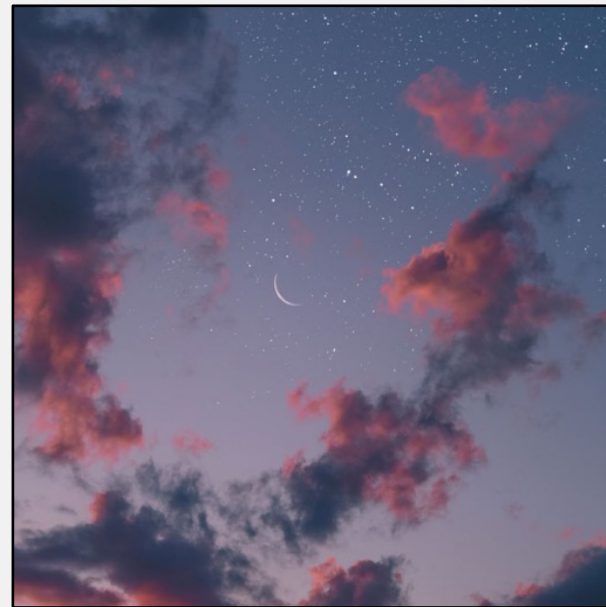
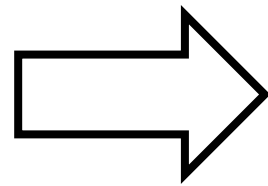
How to improve
reconstruction quality?

Optimization space
(model weights, weight subsets,
extended embedding space)

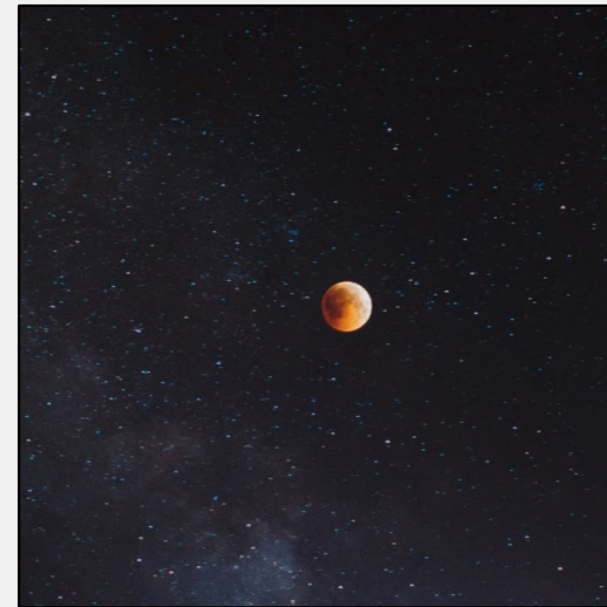
Model Training

Generate/Retrieve images with similar captions

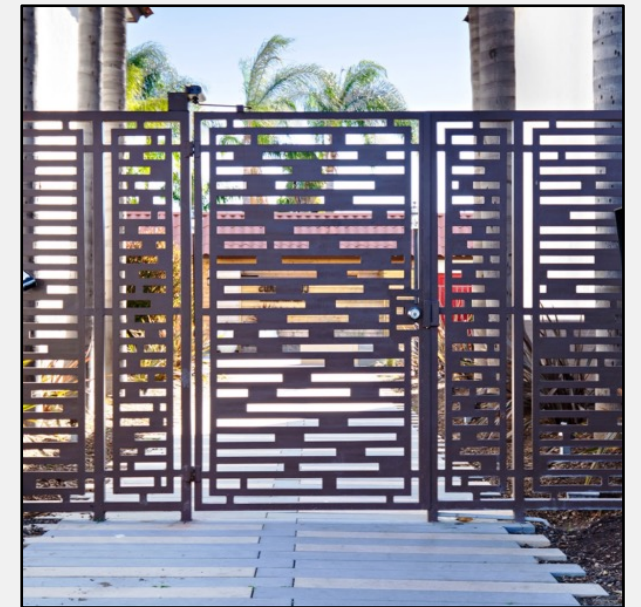
Moongate



sky full of
stars and the
moon



Blood moon



Apartment gates

...

Model Training

Training dataset



Photo of a
moongate

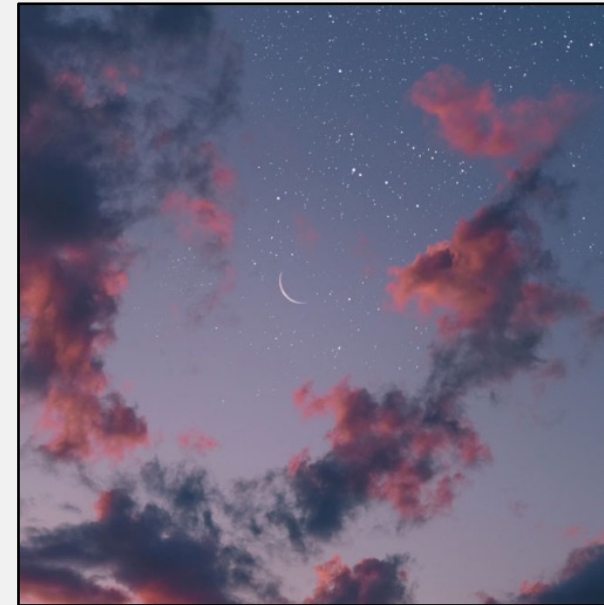


Photo of a
moongate

...

Target images

+



sky full of
stars and the
moon

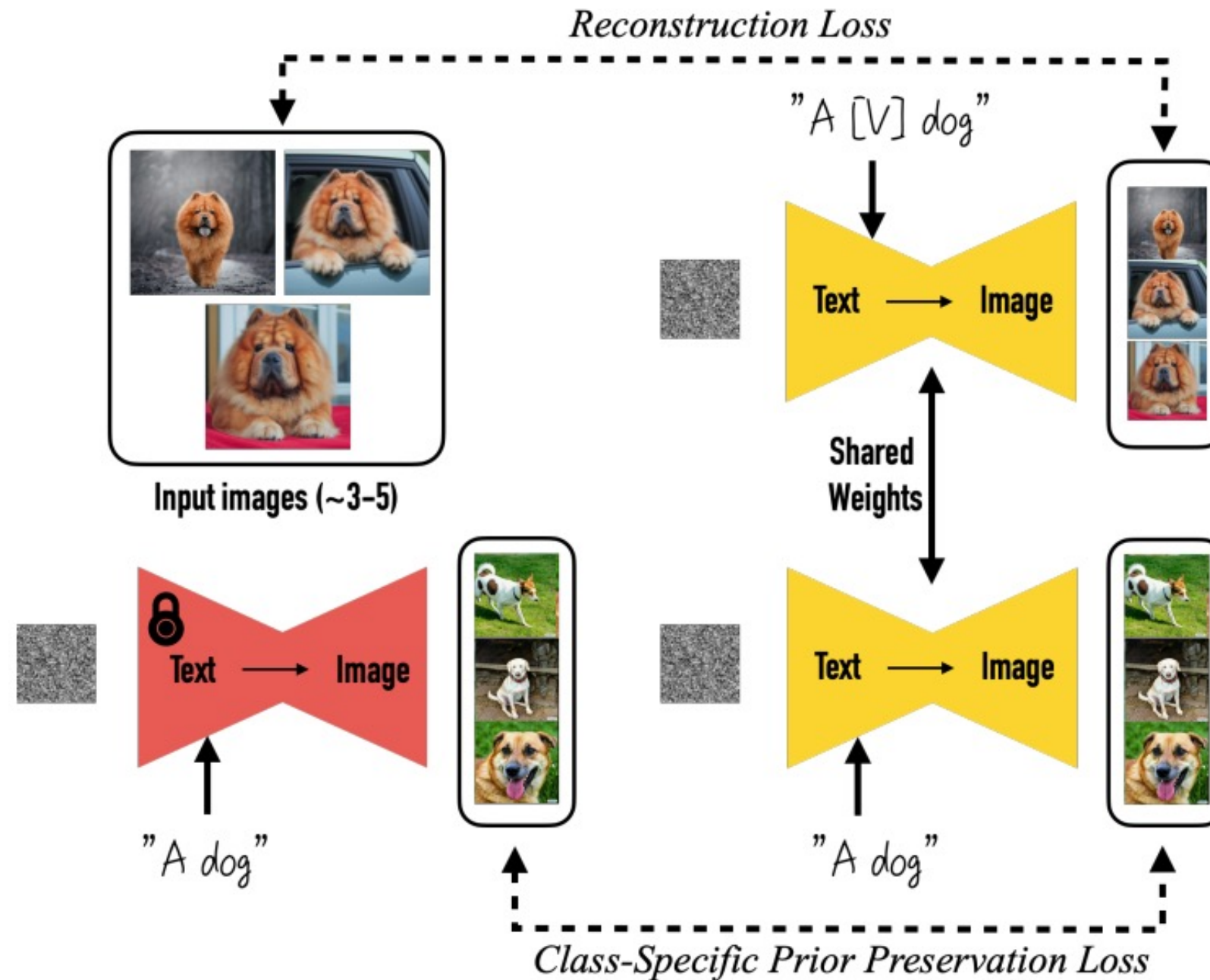


Blood moon

...

Regularization images

DreamBooth Training (Fine-tuning all the weights)



DreamBooth Results



Input images



in the Acropolis



swimming



sleeping



in a doghouse

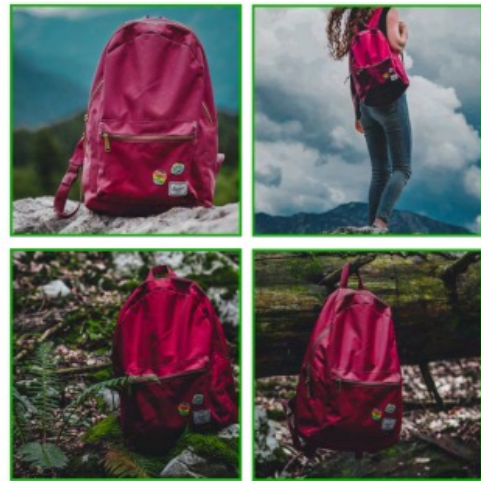


in a bucket



getting a haircut

DreamBooth Results



Input images



A [V] backpack in the Grand Canyon



A wet [V] backpack in water



A [V] backpack in Boston



A [V] backpack with the night sky



Input images



A [V] teapot floating in milk



A transparent [V] teapot with milk inside

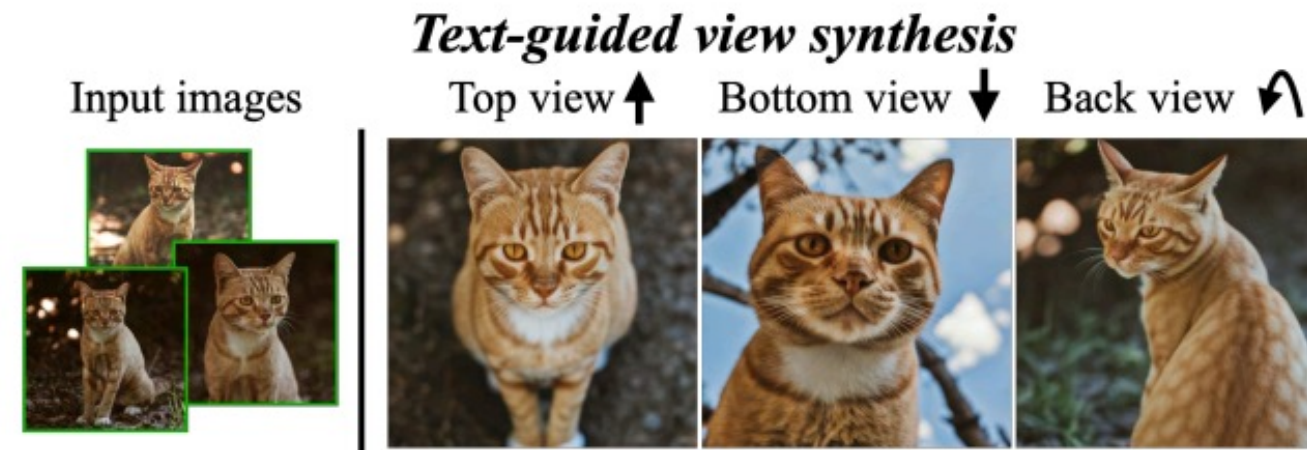


A [V] teapot pouring tea



A [V] teapot floating in the sea

DreamBooth Applications

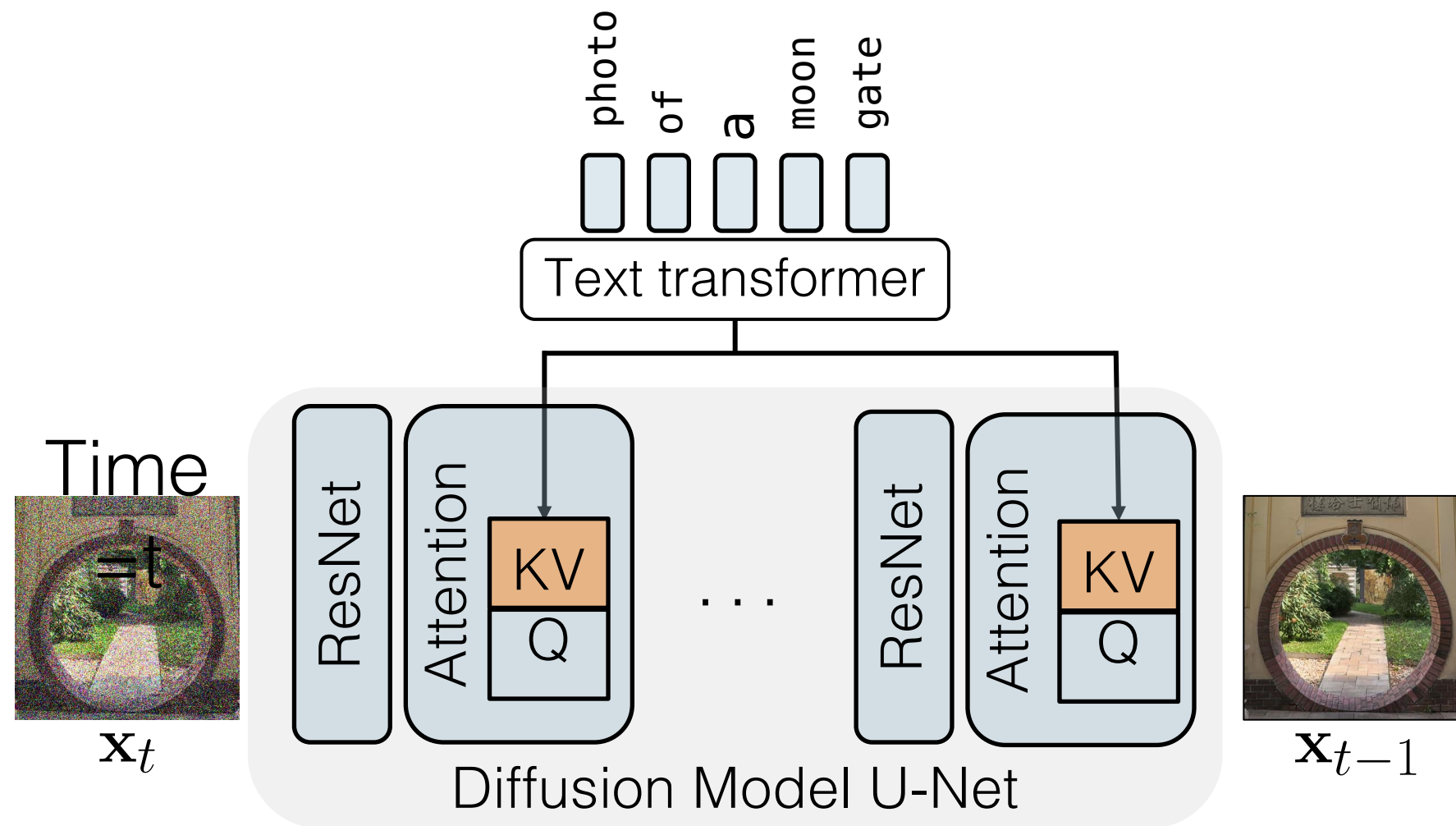


practical concerns

1. training time (30-60 min)
2. storage per concept (3-4 GB)

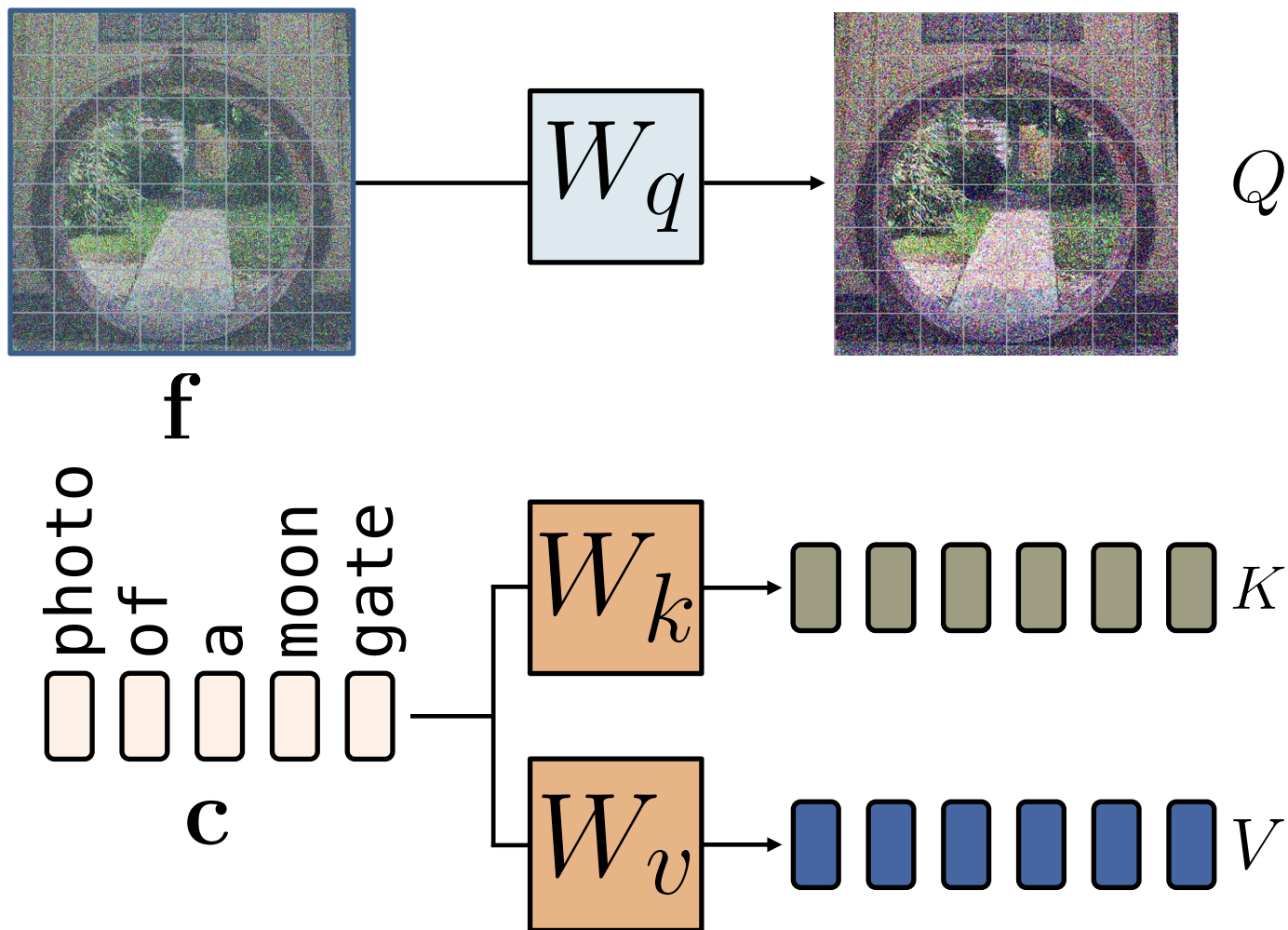
Efficient Custom Diffusion Training

Fine-tune key, value projection matrices in cross-attention layers



 Trainable  Frozen

Text-image Cross-Attention



$$Q \text{ Softmax} \left(* \right) = \text{[Attention Weights]}$$

$$= \sum \left(\text{[Attention Weights]} * \text{[Value V]} \right)$$

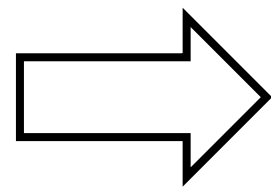
i.e.

$$\text{Output} = \text{Softmax} \left(\frac{Q \cdot K^T}{\sqrt{d'}} \right) V$$

Trainable Frozen

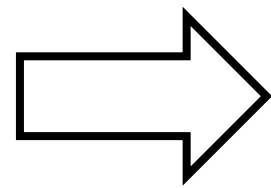
Custom Diffusion for multiple-concepts

Merging two concepts



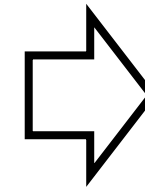
$$W_{k1} \quad W_{v1}$$

+



$$W_{k2} \quad W_{v2}$$

$$W_k \quad W_v$$



V* dog wearing
sunglasses
in front of a
moongate

Custom Diffusion for multiple-concepts

Merging two concepts

$$\hat{W} = \arg \min_W \|WC_{\text{reg}}^\top - W_0 C_{\text{reg}}^\top\|_F$$

s.t. $WC^\top = V$, where $C = [\mathbf{c}_1 \cdots \mathbf{c}_N]^\top$
and $V = [W_1 \mathbf{c}_1^\top \cdots W_N \mathbf{c}_N^\top]^\top$.

C_{reg} : a collection of random text prompts.

C : target prompts.

Custom Diffusion for multiple-concepts

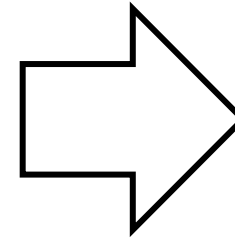
Merging two concepts

$$\hat{W} = W_0 + \mathbf{v}^\top \mathbf{d}, \text{ where } \mathbf{d} = C(C_{\text{reg}}^\top C_{\text{reg}})^{-1}$$
$$\text{and } \mathbf{v}^\top = (V - W_0 C^\top)(\mathbf{d} C^\top)^{-1}.$$

C_{reg} : a collection of random text prompts.

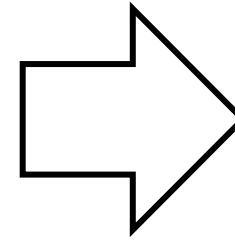
C : target prompts.

More examples



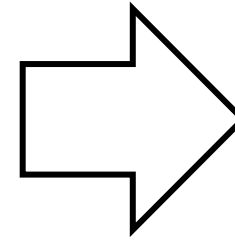
V_1^* chair with the V_2^* cat sitting on it near a beach

More examples



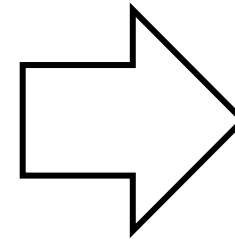
V_1^* chair with the V_2^* cat sitting on it near a beach

More examples



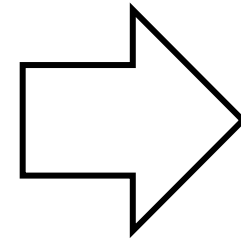
V_1^* chair with the V_2^* cat sitting on it near a beach

More examples



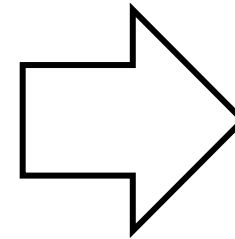
The V_1^* cat is sitting inside a V_2^* wooden pot and looking up

More examples



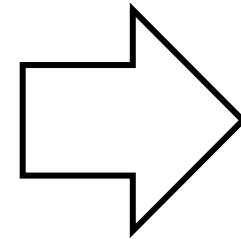
The V_1^* cat is sitting inside a V_2^* wooden pot and looking up

More examples



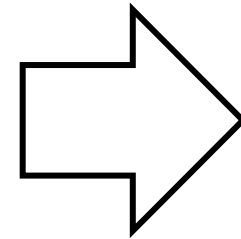
The V_1^* cat is sitting inside a V_2^* wooden pot and looking up

More examples



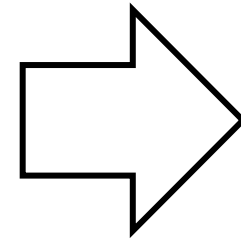
V_1^* flower in the V_2^*
wooden pot on a table

More examples



V_1^* flower in the V_2^*
wooden pot on a table

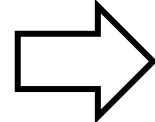
More examples



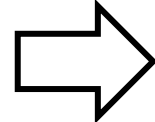
V_1^* flower in the V_2^*
wooden pot on a table



Drawings from Aaron Hertzmann



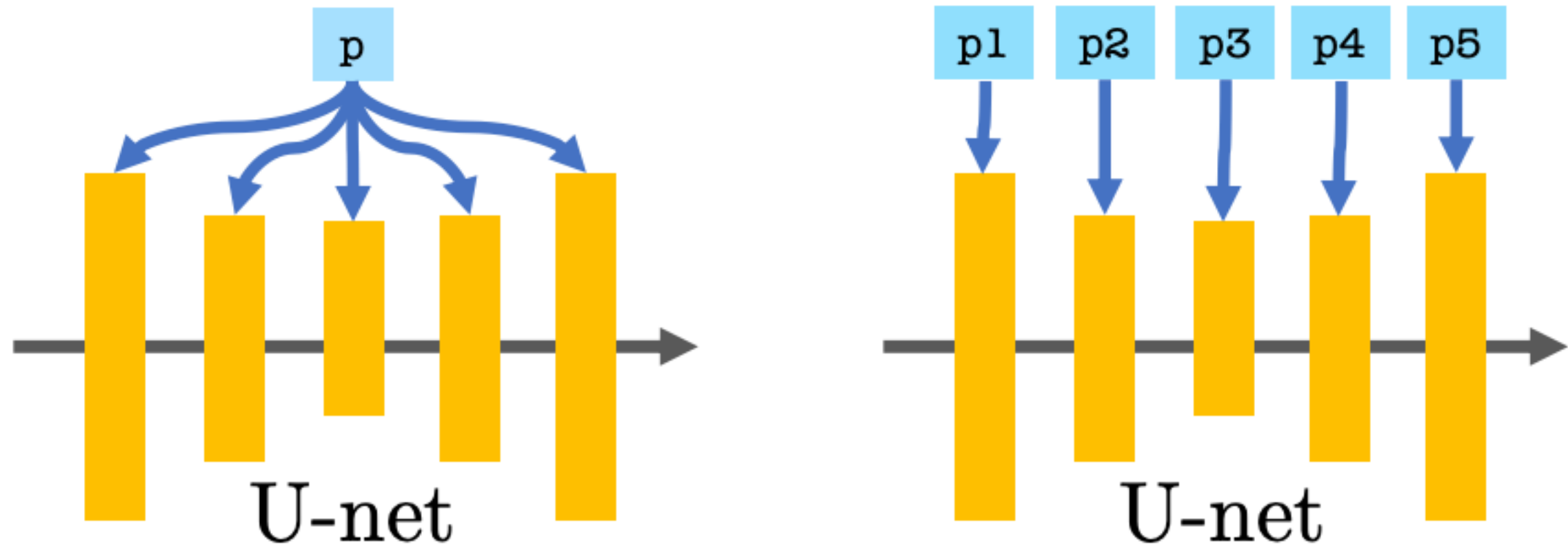
Plant painting in style of **V*** art



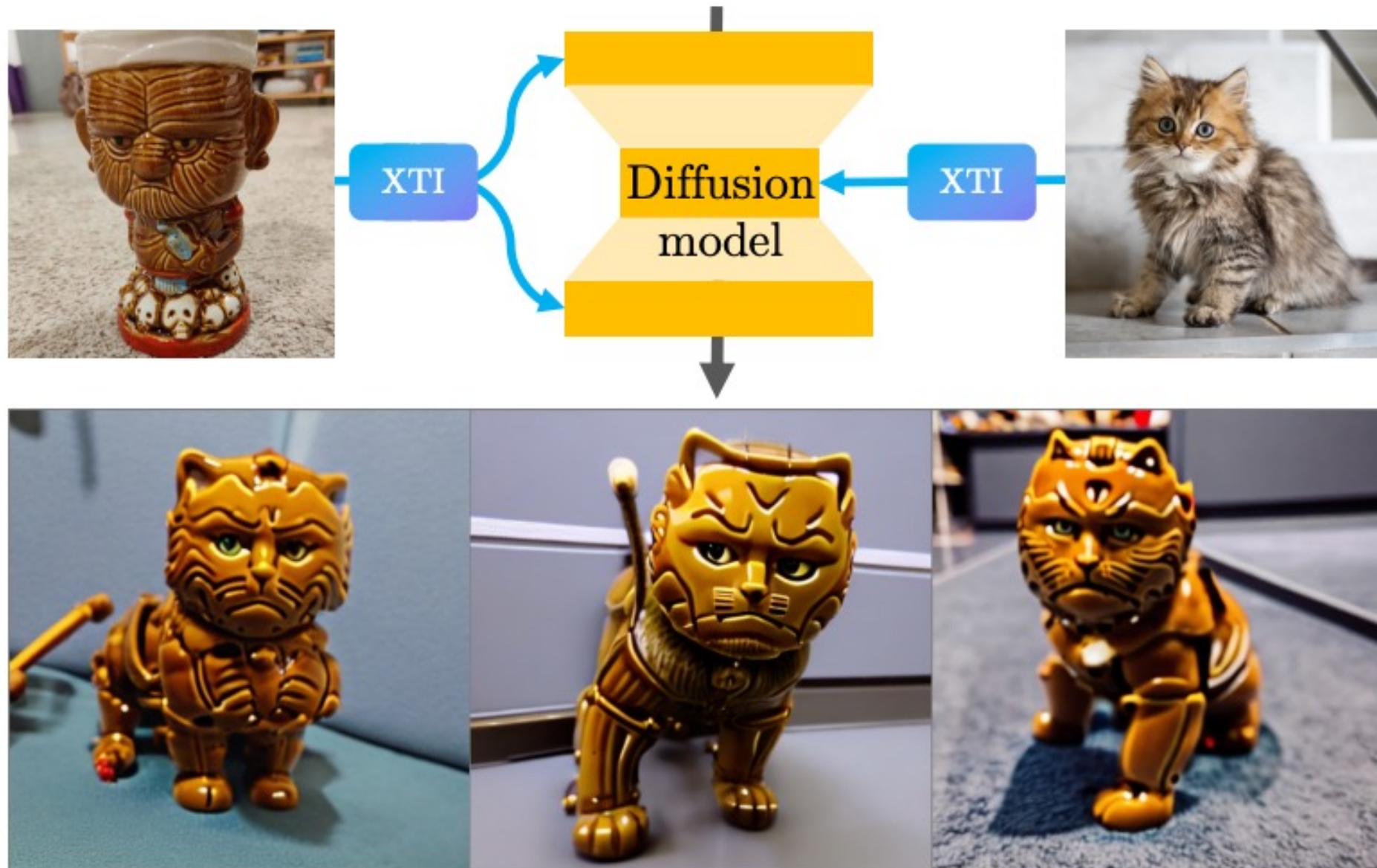
Drawings from Aaron Hertzmann

Painting of dog in style of V* art

Extended Textual Inversion



Shape-Style Mixing



Extended Textural Inversion

Real



Textual Inversion



Extended Textual Inversion



<teddy bear> in Times Square



<cat> wearing sunglasses