

Ideal models (Dream)

Pros: good sample, fast sample, Exact/fast likelihoods
good coverage, easy to training, learn low-dimensional latent representation.
Cons: no cons

Autoregressive models

Pros: Exact likelihoods, good samples, good coverage
Cons: Slow to evaluate or sample

VAEs

Pros: Cheap to sample, good coverage
Cons: Blurry samples (in practice)

GANs

Pros: Cheap to sample, fast to train
Cons: No likelihoods (density), hard to train

Diffusion models

Pros: Easy to train, good samples
Cons: slow to sample; slow/hard to compute likelihoods

No Free Lunch

Stable Training \leftrightarrow Slow Inference

Hybrid Models

VQ-VAE2: VAE + autoregressive

VQ-GAN: VAE, GANs, Perceptual loss + autoregressive

Latent Diffusion: VAE, GANs, perceptual loss + diffusion models

Base model (feature): autoregressive, diffusion

Upsampler (feature->image): VAE, GANs, Perceptual loss

Model Distillation

Teacher: Multi-step Models (e.g., Latent Diffusion Models)

Student: Single-Step Model (e.g., Conditional GANs)

Distillation Loss: GAN Loss, Perceptual Loss, ...

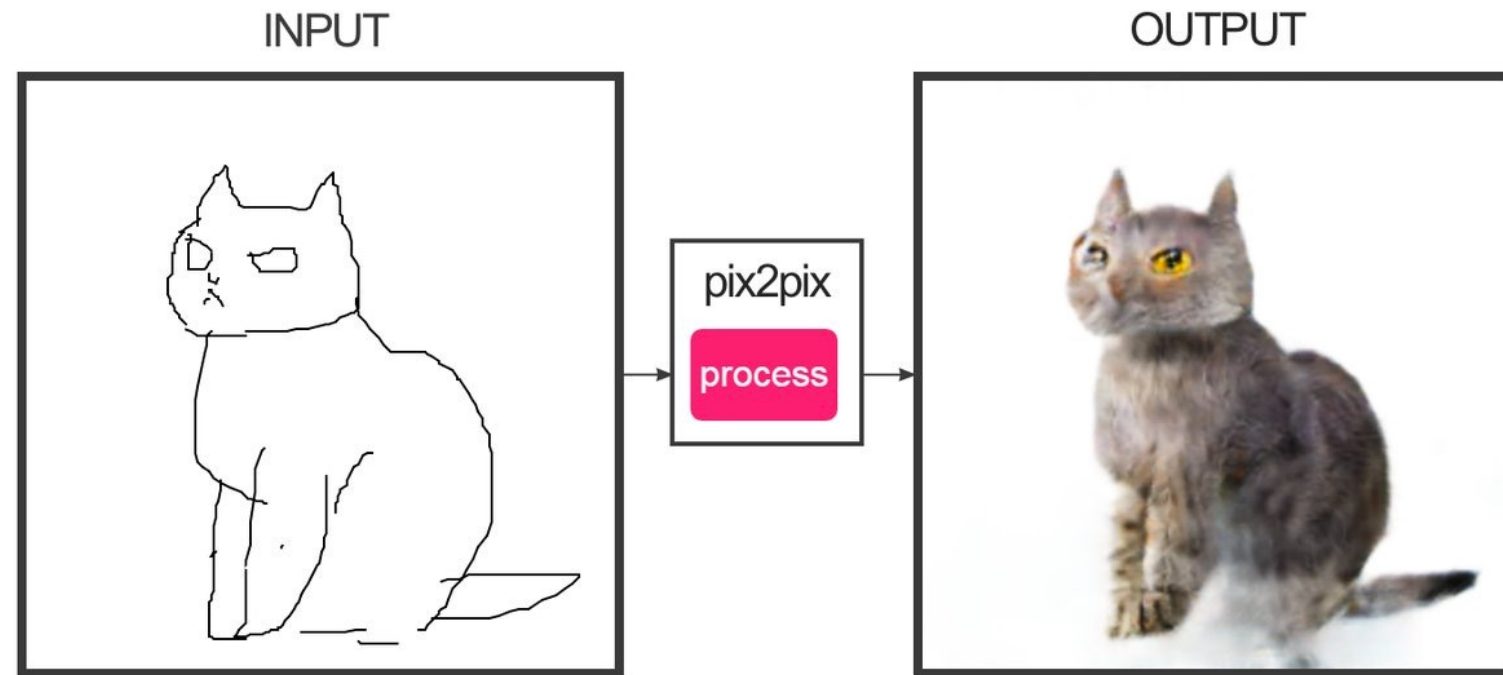
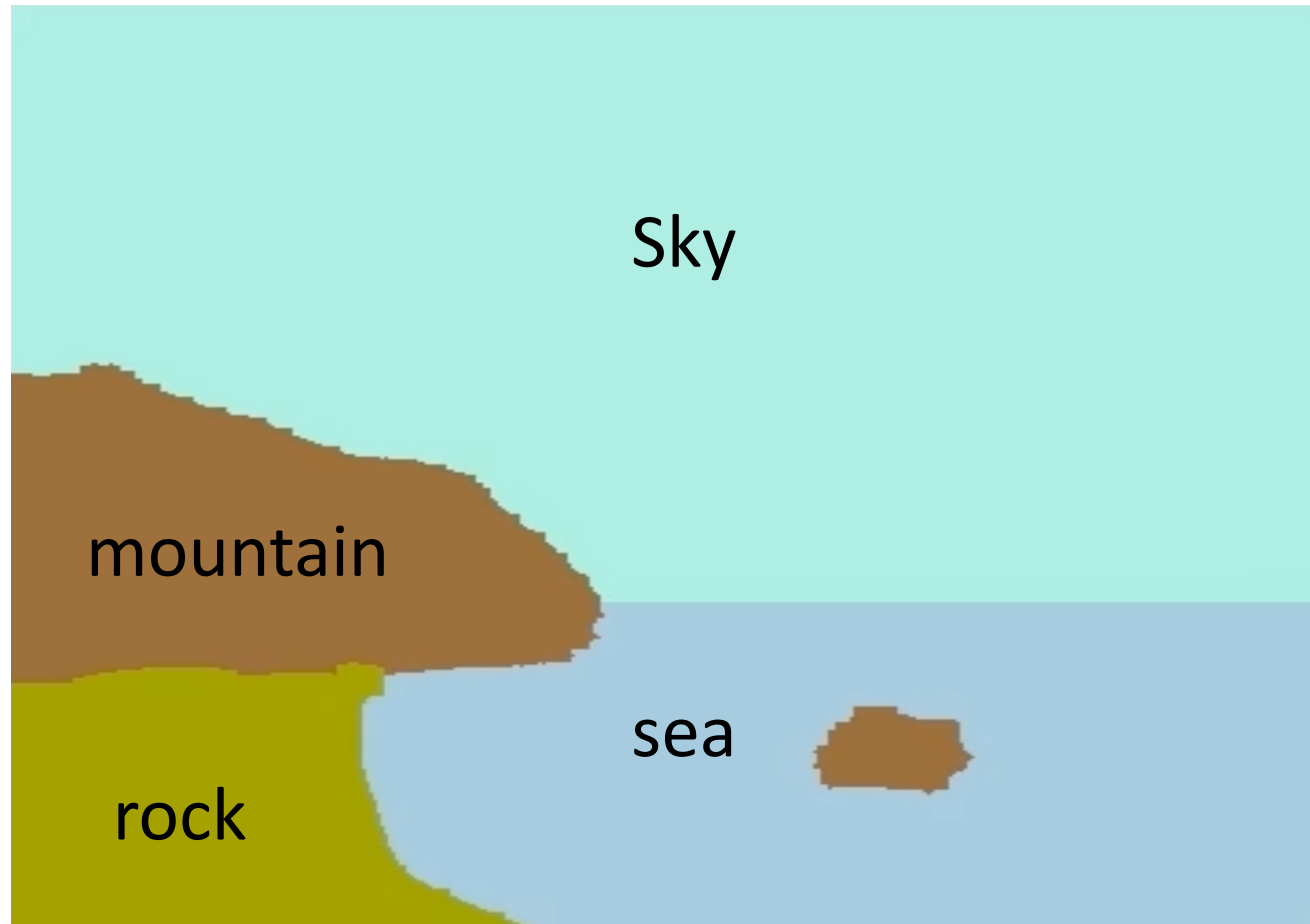


Image-to-Image Translation

Conditional Generative Models

Jun-Yan Zhu
16-726, Spring 2025

Problem Statement



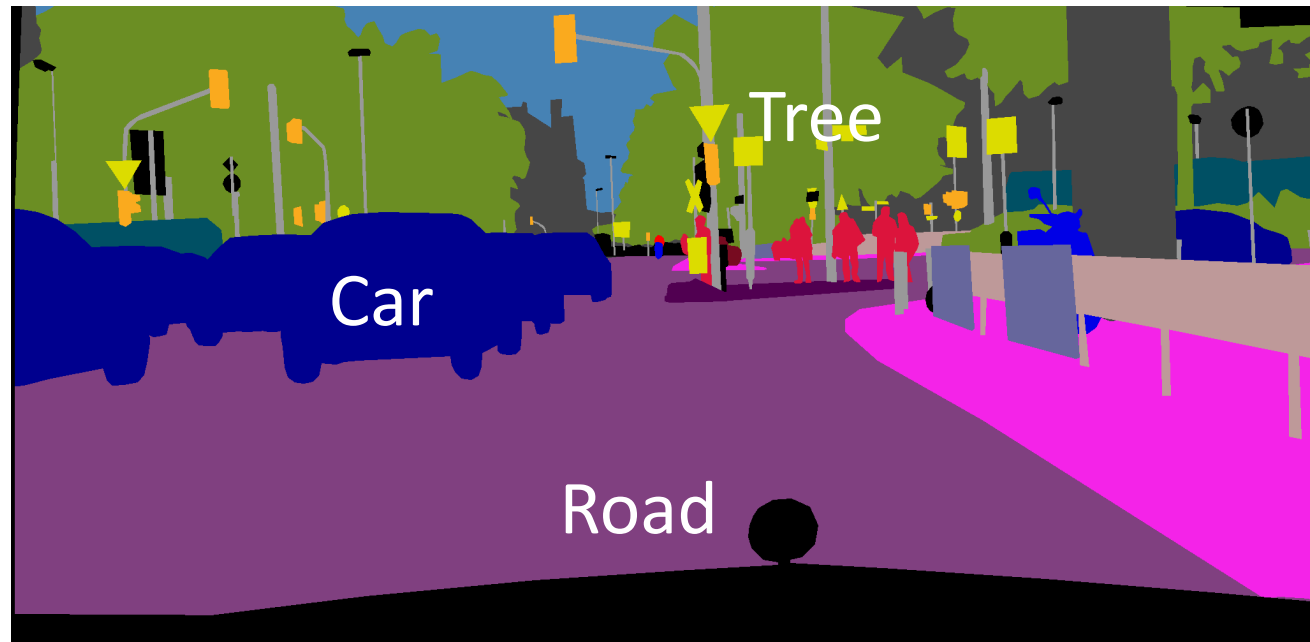
Input



Output

Goal: synthesize a photograph given an input image

Problem Statement



Input



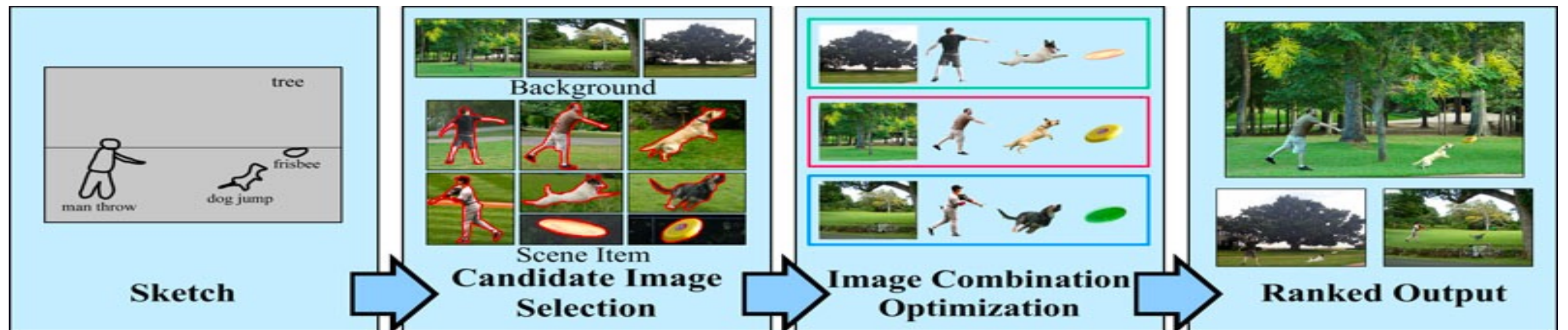
Output

Goal: synthesize a photograph given an input image

Early work (Example-based)

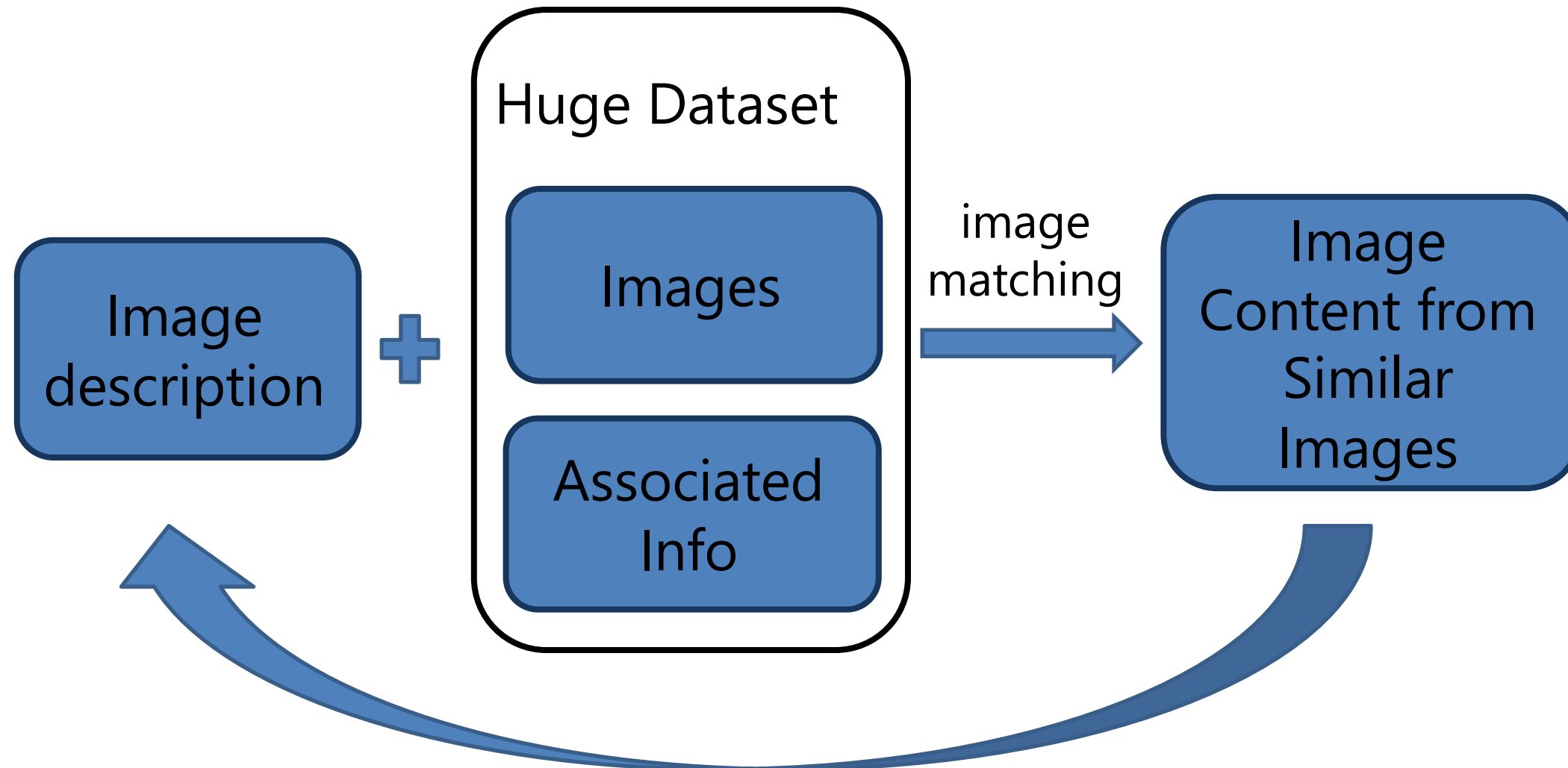


Semantic Photo Synthesis [Johnson et al., Eurographics 2006]



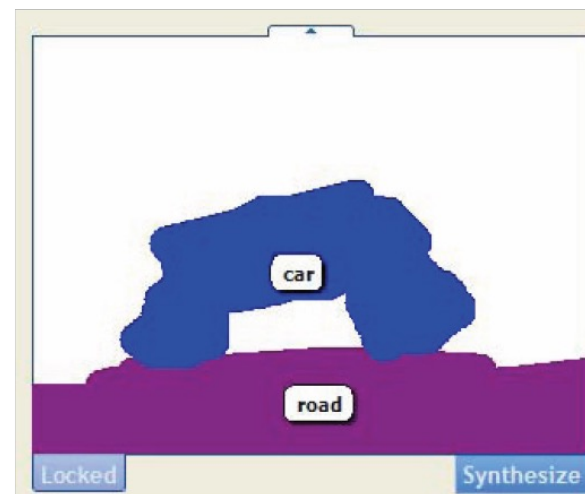
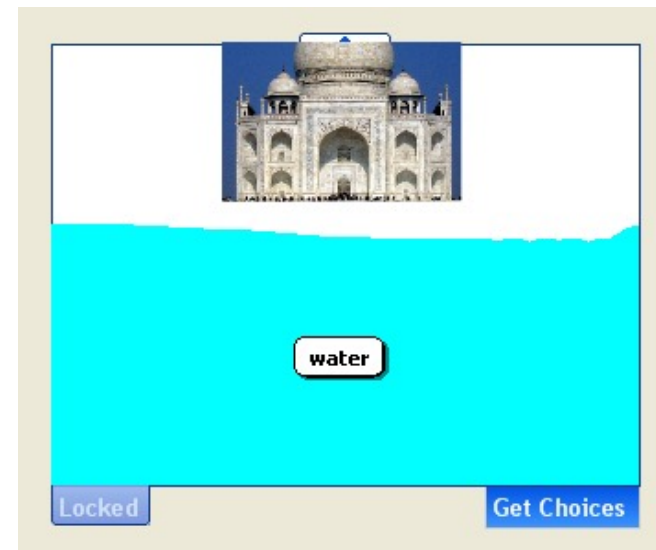
Sketch2Photo [Tao et al., SIGGRAPH Asia 2009]

Semantic Photo Synthesis



M. Johnson, G. Brostow, J. Shotton, O. A. C., and R. Cipolla, "Semantic Photo Synthesis," Eurographics 2006

Semantic Photo Synthesis [EG'06]



M. Johnson, G. Brostow, J. Shotton, O. A. C., and R. Cipolla, "Semantic Photo Synthesis," Eurographics 2006

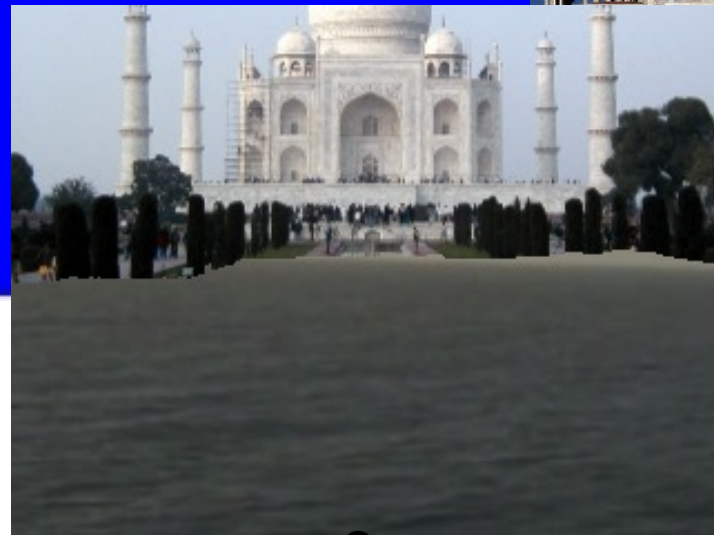
Semantic Photo Synthesis



1



2

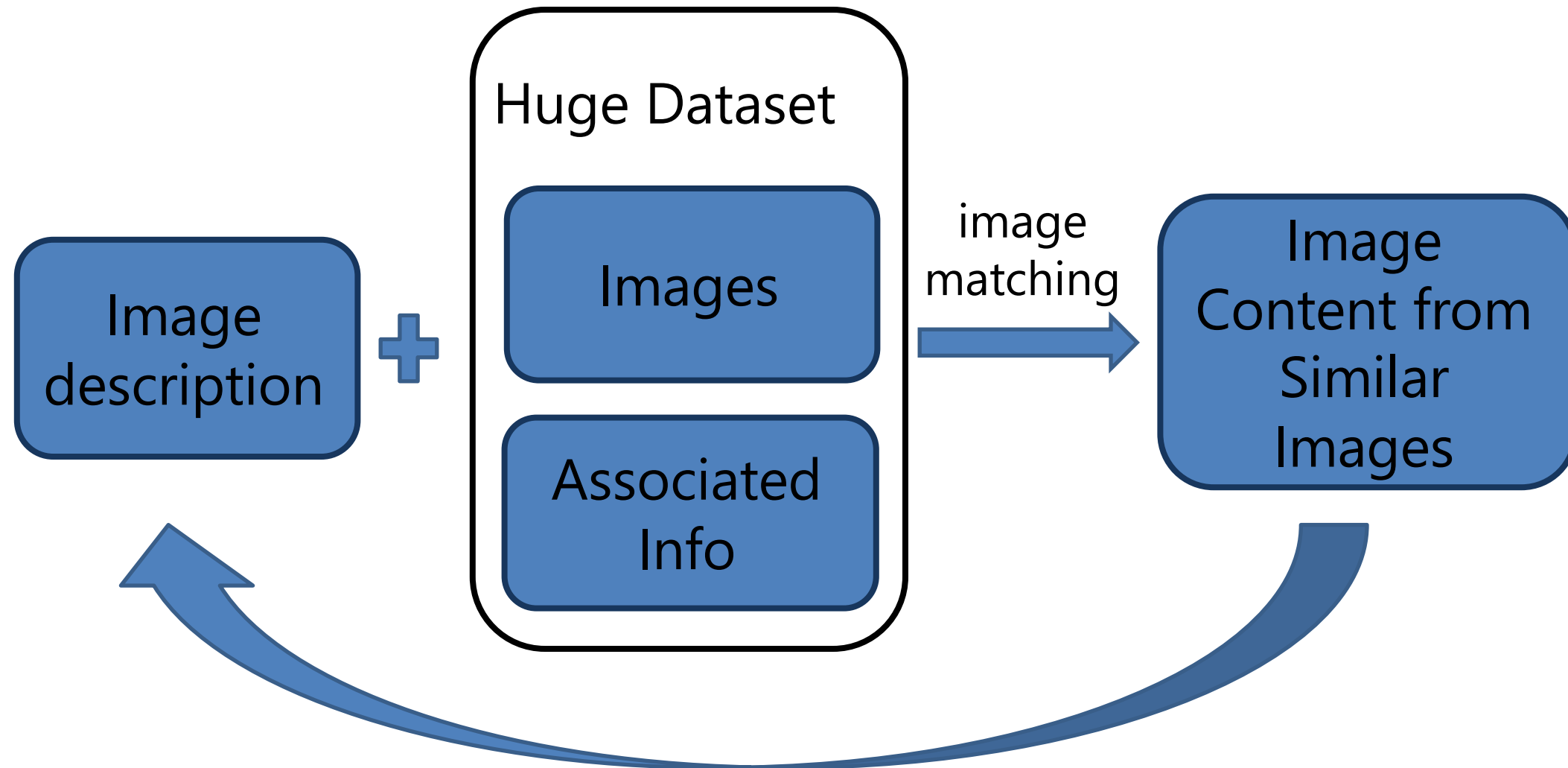


3



4

Semantic Photo Synthesis



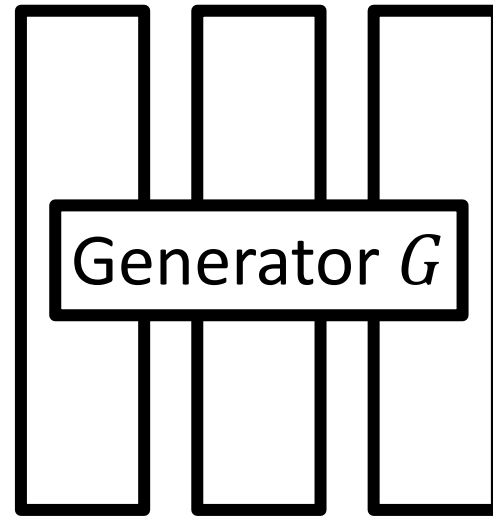
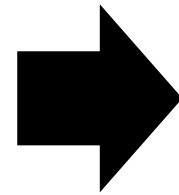
M. Johnson, G. Brostow, J. Shotton, O. A. C., and R. Cipolla, "Semantic Photo Synthesis," Computer Graphics Forum Journal (Eurographics 2006), vol. 25, no. 3, 2006.

Learning-based methods

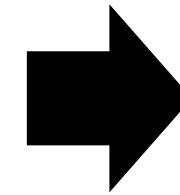
Loss functions for Image Synthesis



Input x



Learnable rendering



Output Image $G(x)$

What is a good objective \mathcal{L} ?

- What is a good loss?
- How to calculate it efficiently?
- How to collect data (x, y) ?

Problem Statement

Loss function



$$\arg \min_G \mathcal{L}(G(x), y)$$

G

Generator

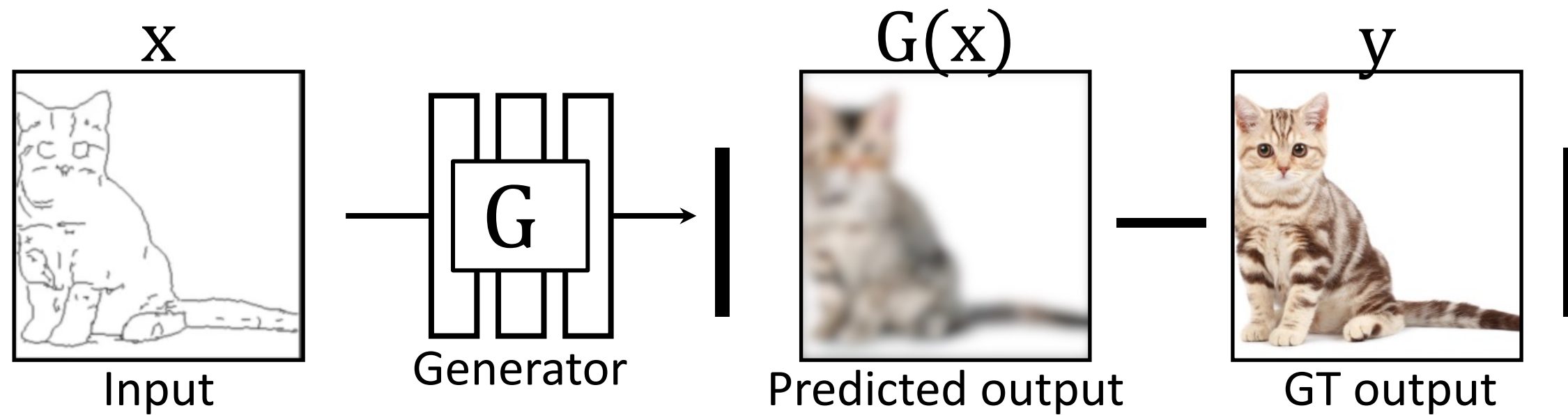
x

Input Info

y

Output image

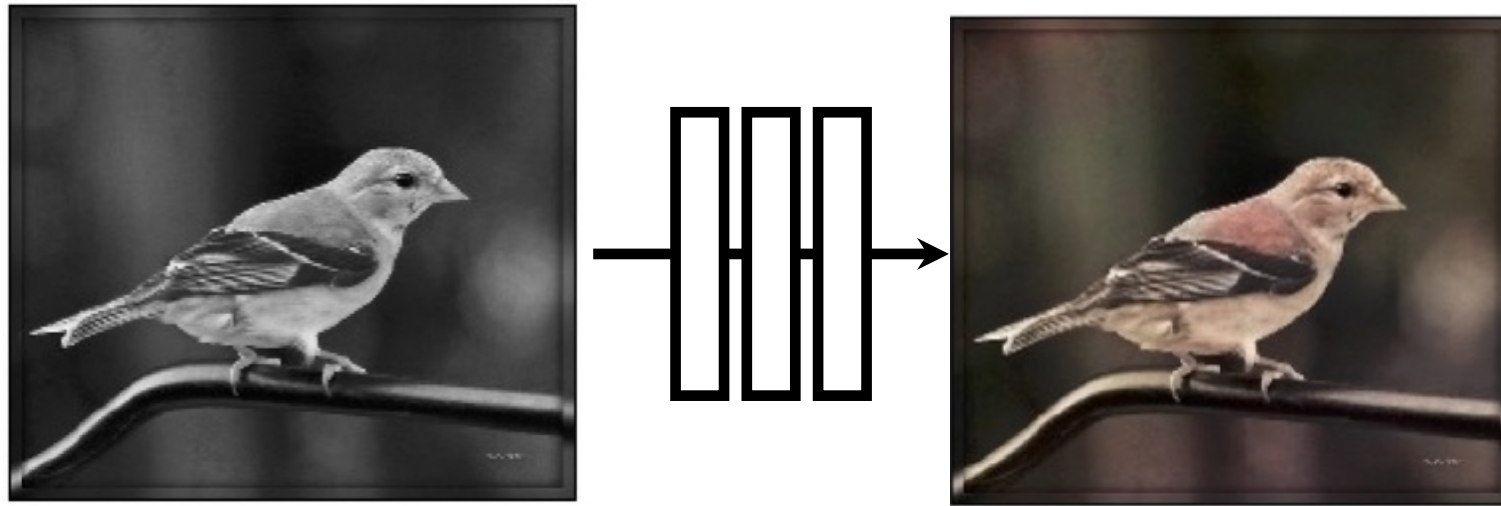
Designing Loss Functions



L2 regression $\arg \min_G \mathbb{E} [||G(x) - y||]$

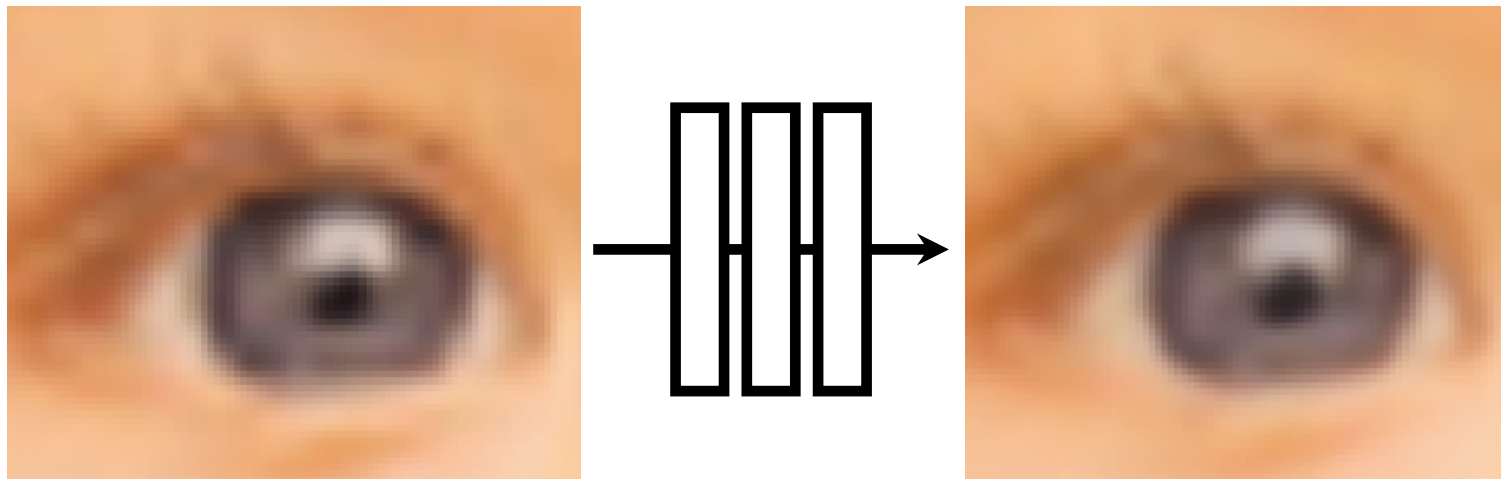
Designing Loss Functions

Image colorization



L2 regression

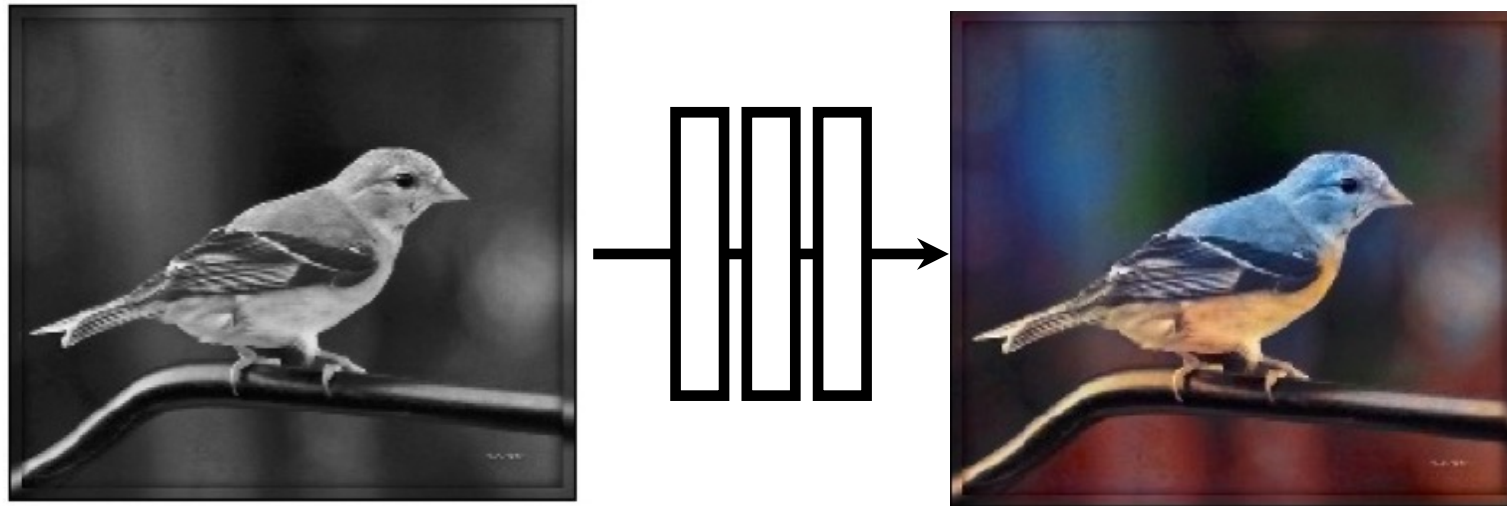
Super-resolution



L2 regression

Designing Loss Functions

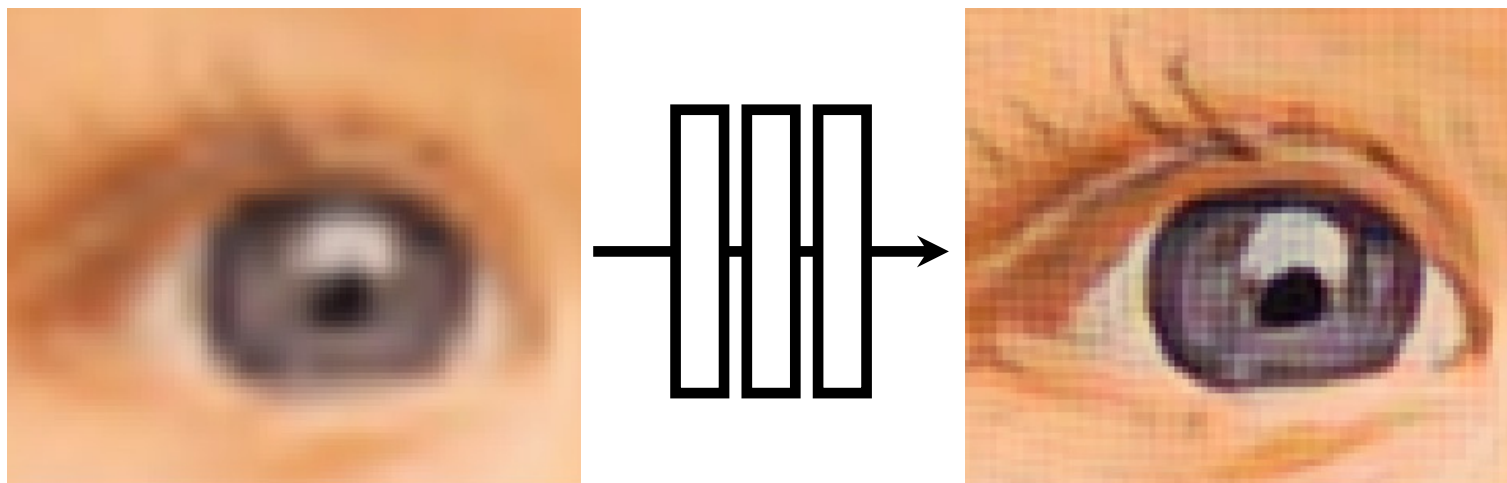
Image colorization



[Zhang et al. 2016]

Classification Loss:
Cross entropy objective,
with colorfulness term

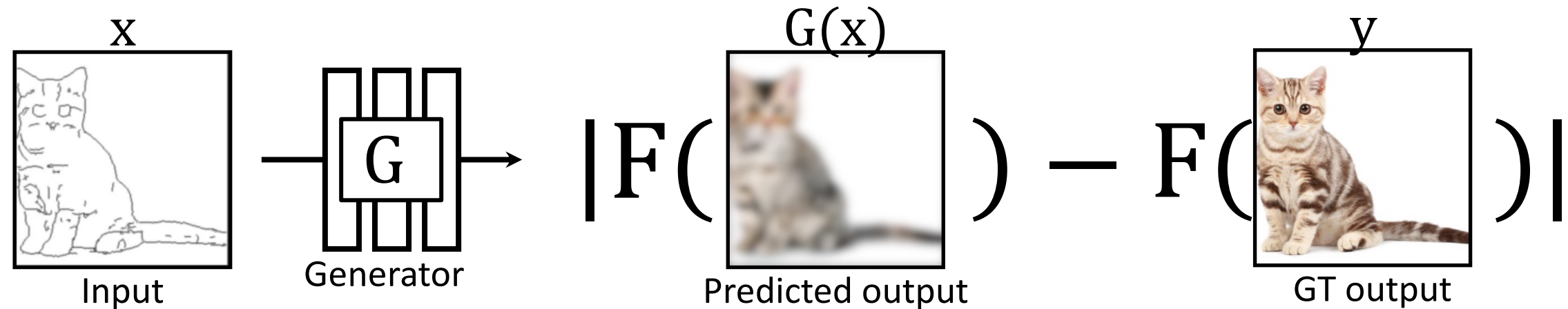
Super-resolution



Feature/Perceptual loss
Deep feature matching
objective

[Gatys et al., 2016], [Johnson et al. 2016], [Dosovitskiy and Brox. 2016]

CNNs as a Perceptual Metric



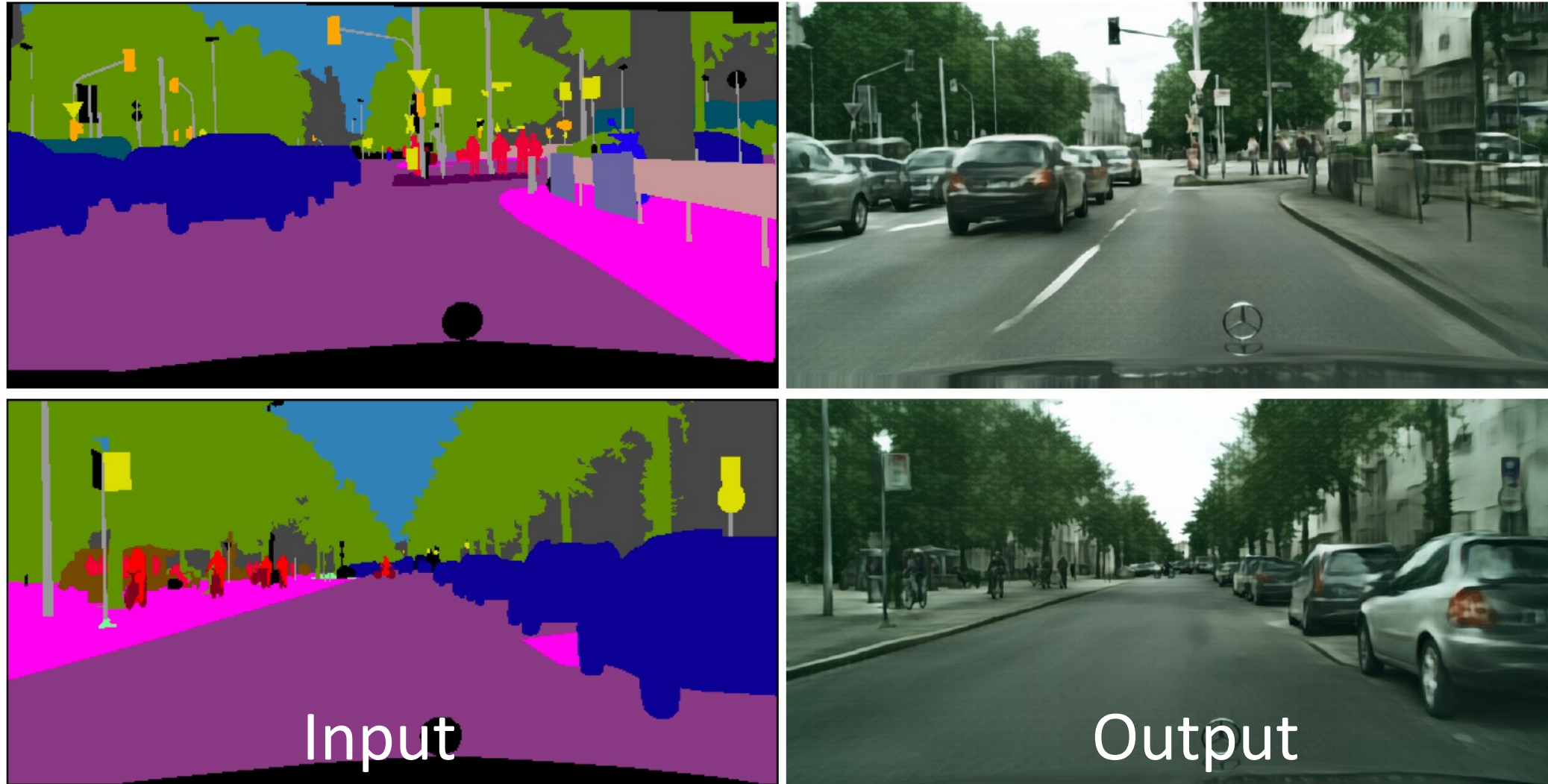
F is a deep network (e.g., ImageNet classifier)

Perceptual Loss

$$\arg \min_G \mathbb{E}_{(x,y)} \sum_{i=1}^N \lambda_i \frac{1}{M_i} \left\| F^{(i)}(G(x)) - F^{(i)}(y) \right\|_2^2$$

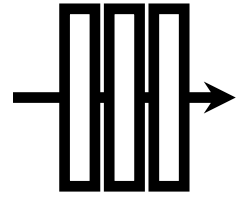
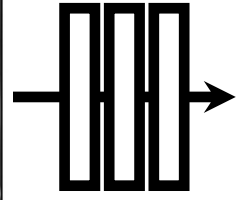
The number of elements in the (i)-th layer

Learning with Perceptual Loss



Training objective: $\arg \min_G \mathbb{E}_{(x,y)} \sum_{i=1}^N \lambda_i \frac{1}{M_i} \|F^{(i)}(G(x)) - F^{(i)}(y)\|_2^2$

Generated images



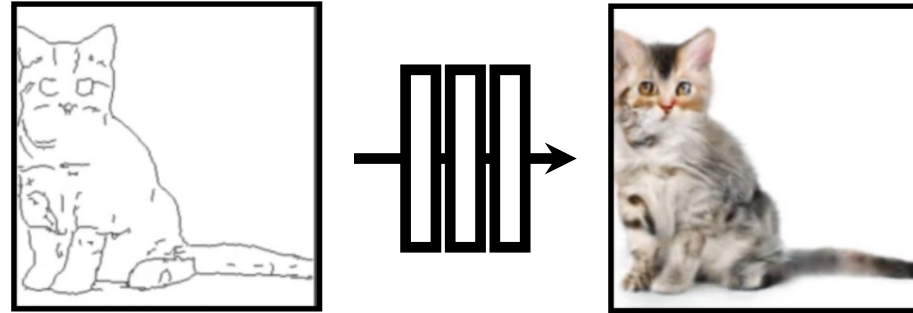
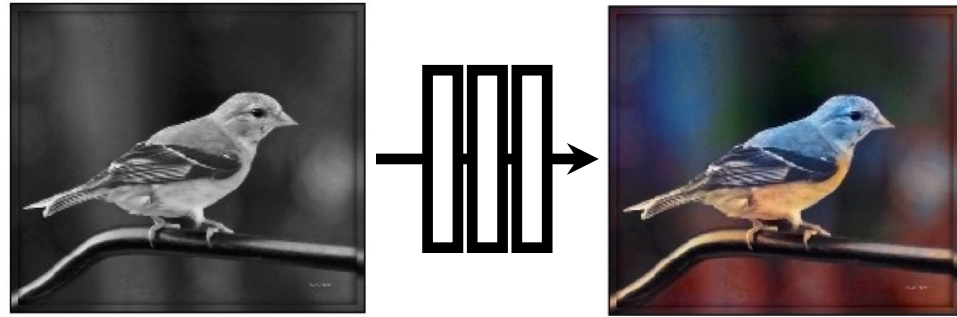
⋮

⋮



Universal loss?

Generated images



⋮

⋮

Generative Adversarial Network (GANs)

Classifier

Real vs. Fake



Real photos



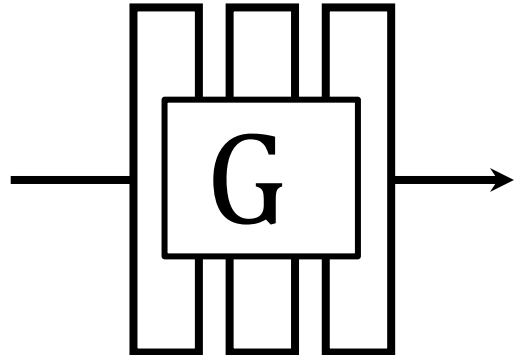
⋮

[Goodfellow, Pouget-Abadie, Mirza, Xu, Warde-Farley, Ozair, Courville, Bengio 2014]

x



Input image

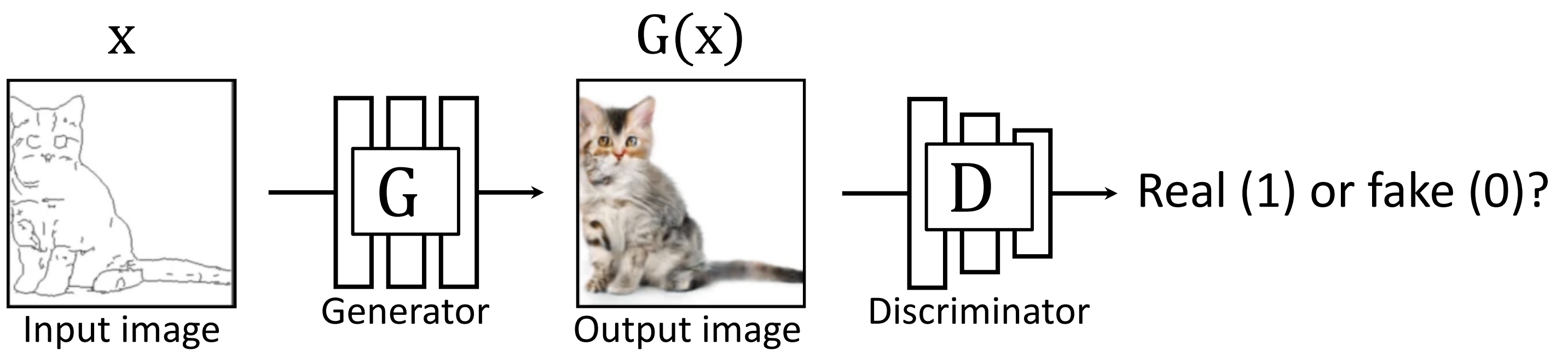


Generator

$G(x)$

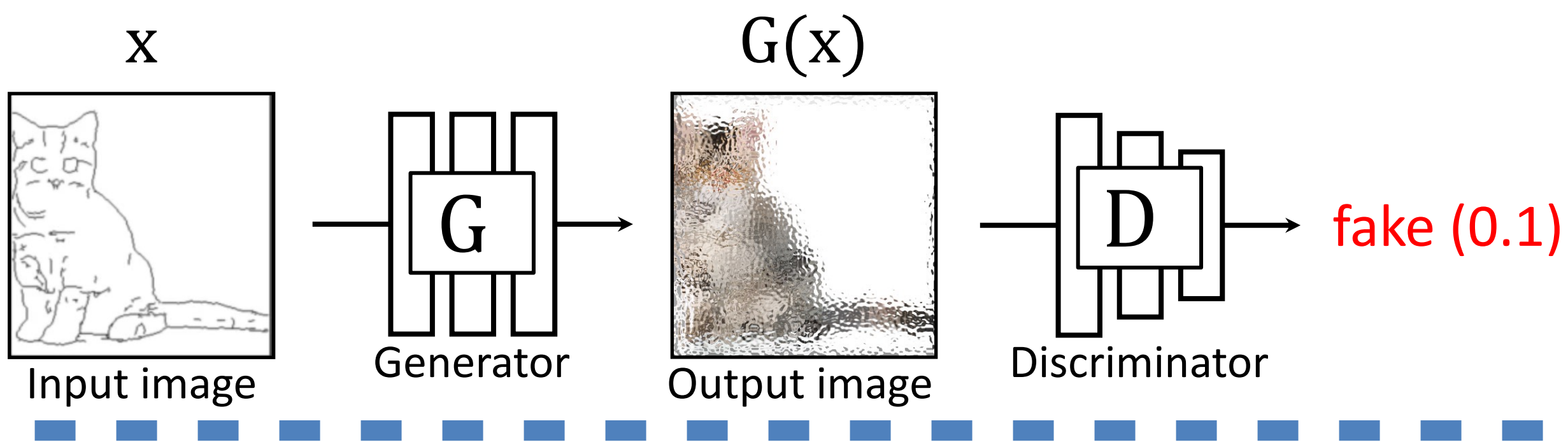


Output image



A two-player game:

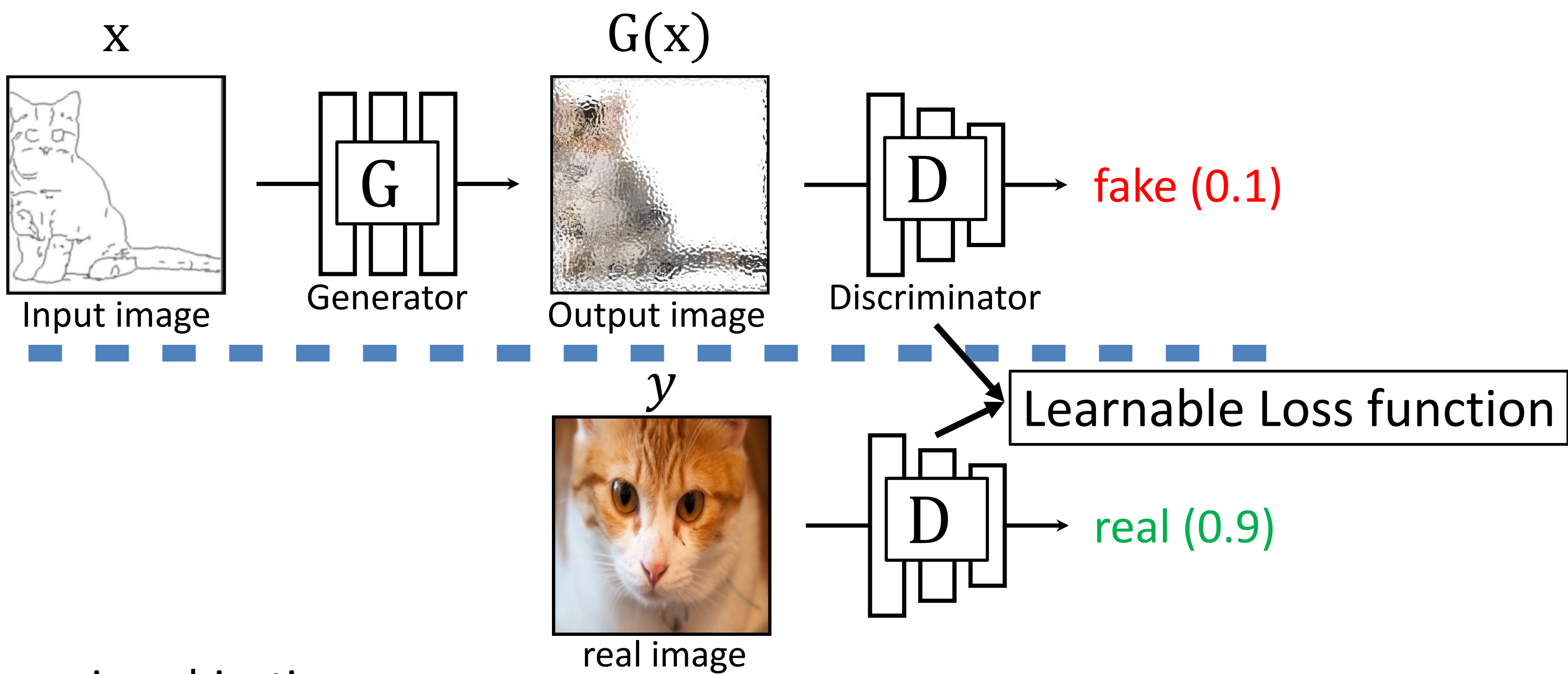
- G tries to generate fake images that can fool D .
- D tries to detect fake images.



Learning objective

$$\min_G \max_D \mathbb{E}_x [\log(1 - D(G(x)))] + \mathbb{E}_y [\log D(y)]$$

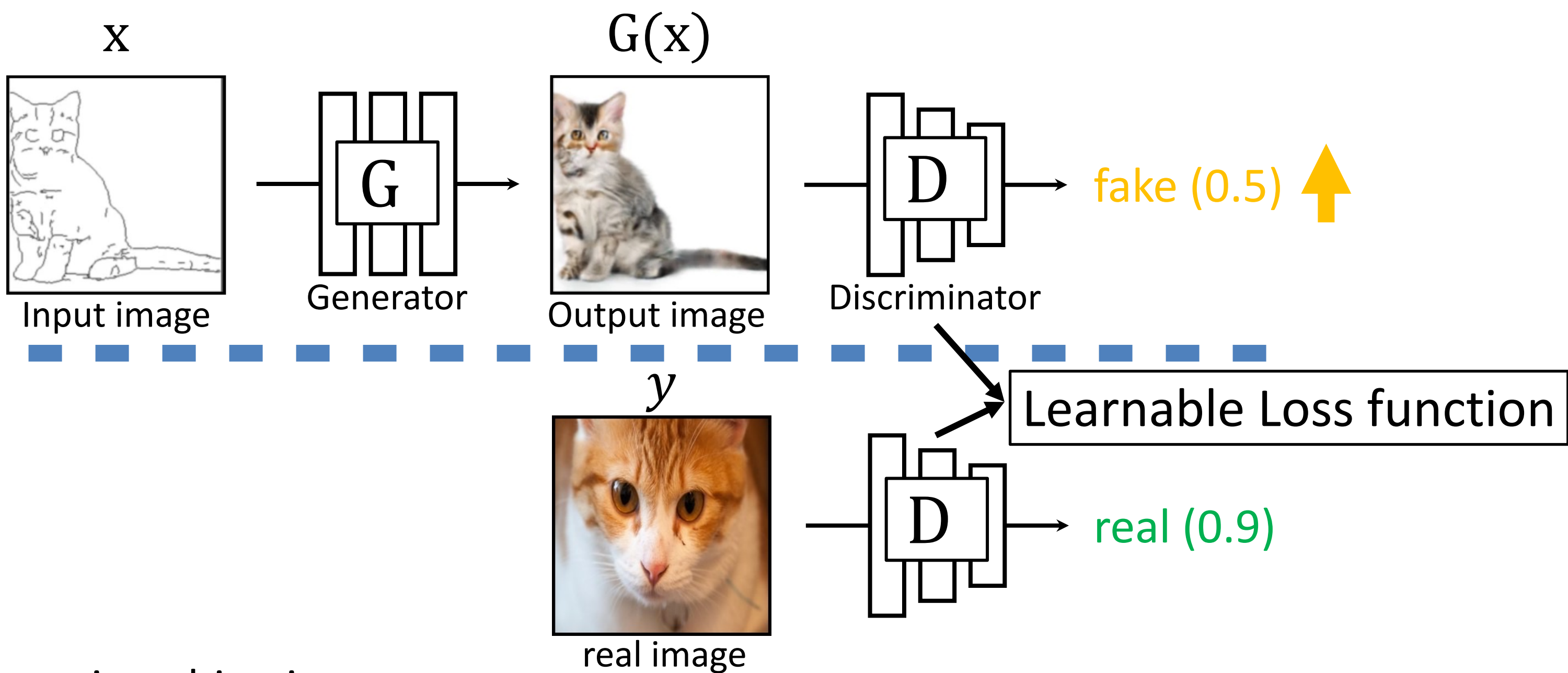
[Goodfellow et al. 2014]



Learning objective

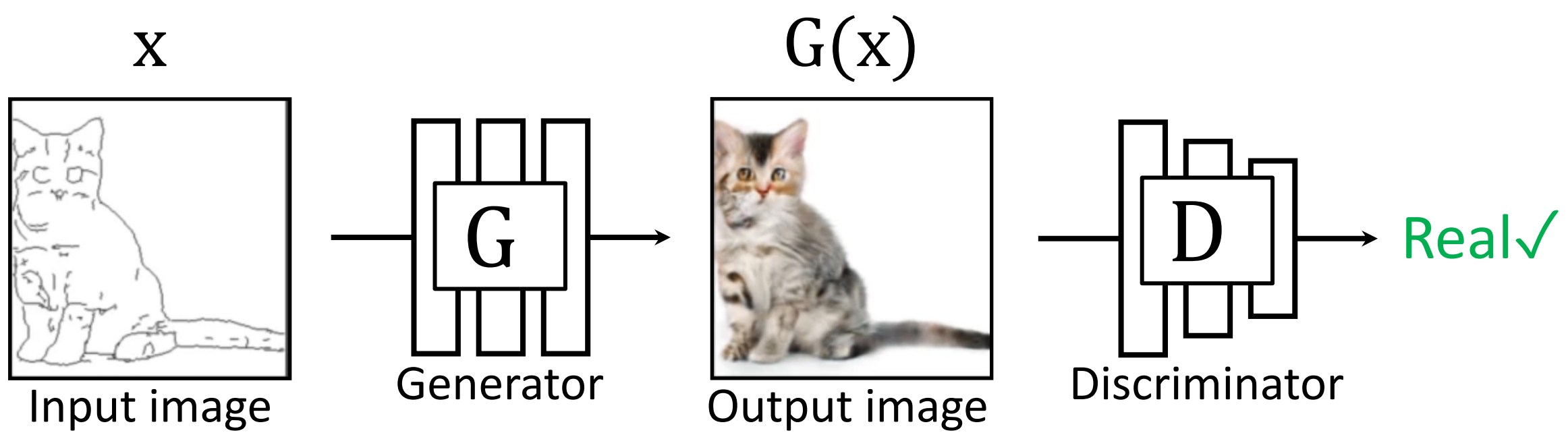
$$\min_G \max_D \mathbb{E}_x [\log(1 - D(G(x)))] + \mathbb{E}_y [\log D(y)]$$

[Goodfellow et al. 2014]



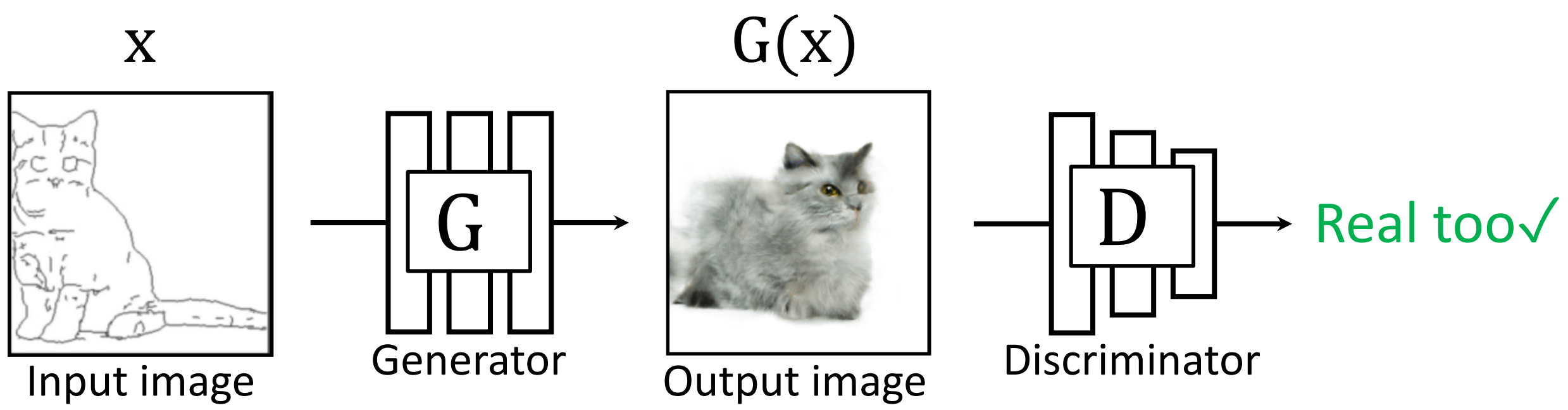
Learning objective

$$\min_G \max_D \mathbb{E}_x [\log(1 - D(G(x)))] + \mathbb{E}_y [\log D(y)]$$



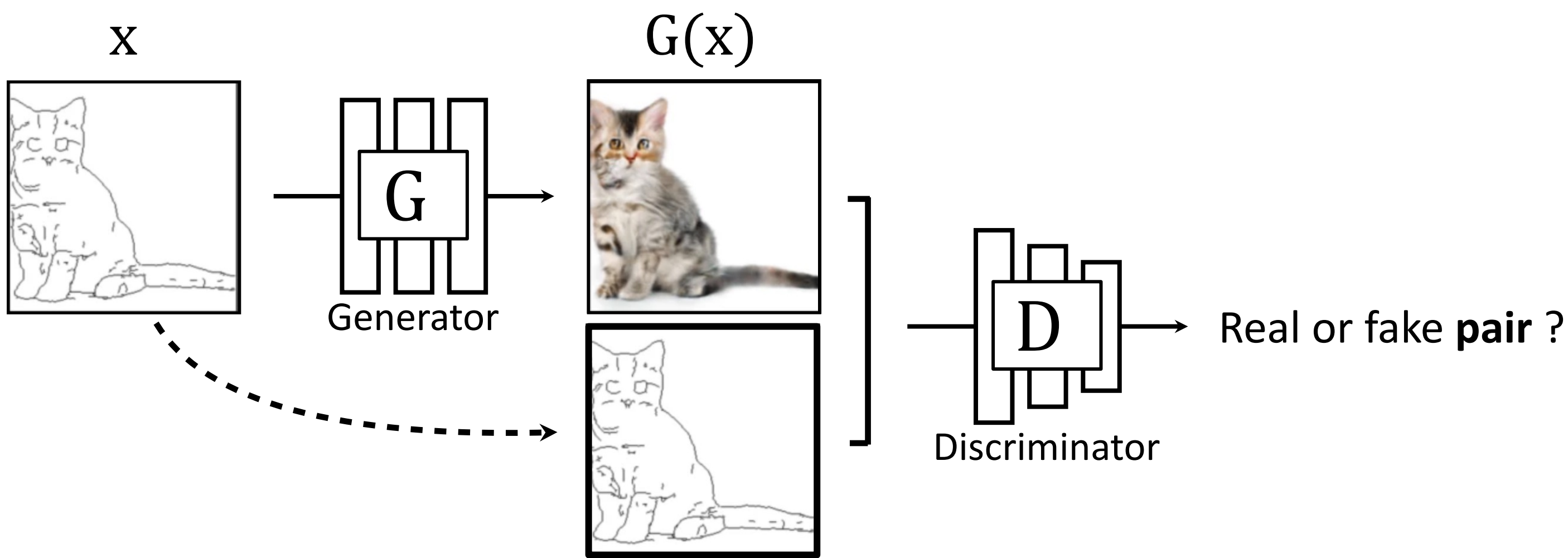
Learning objective

$$\min_G \max_D \mathbb{E}_x [\log(1 - D(G(x)))] + \mathbb{E}_y [\log D(y)]$$



Learning objective

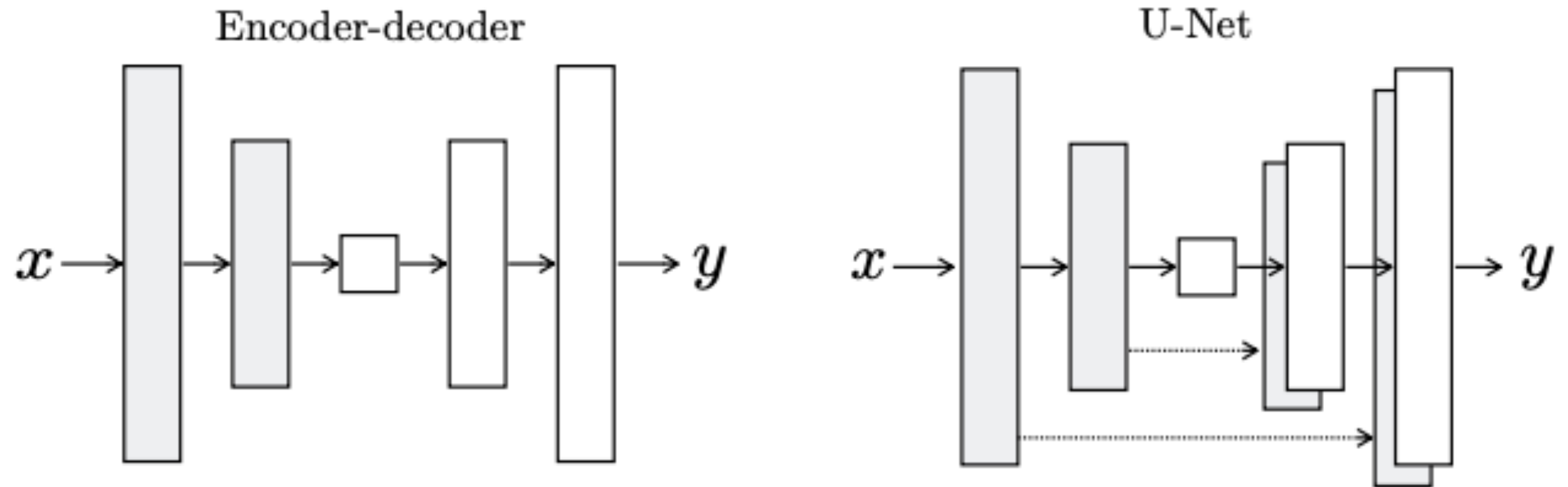
$$\min_G \max_D \mathbb{E}_x [\log(1 - D(G(x)))] + \mathbb{E}_y [\log D(y)]$$



Learning objective

$$\min_G \max_D \mathbb{E}_x [\log(1 - D(x, G(x)))] + \mathbb{E}_{x,y} [\log D(x, y)]$$

pix2pix Generator (U-Net)



U-Net [Ronneberger et al.]: popular CNN backbone for biomedical image segmentation

U-Net: preserve high-frequency information (e.g., edge) of the input image.

Encoder-decoder: lose high-frequency details due to the information bottleneck

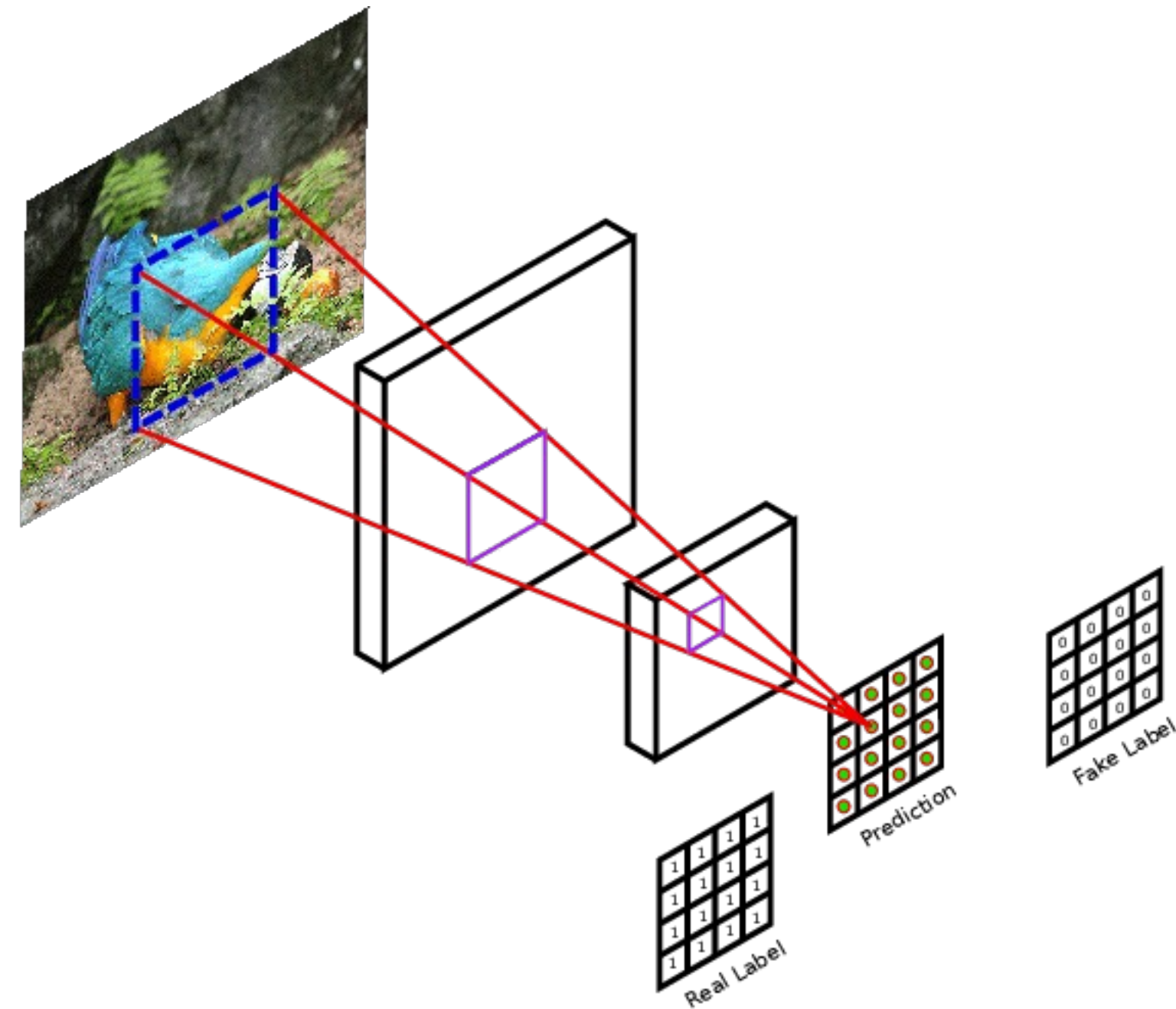
pix2pix Generator (U-Net)



Generator design is critical for image quality.

cGAN (conditional GANs) loss: capture realism. L1 loss stabilizes training (faster convergence)

pix2pix Discriminator (PatchGAN)

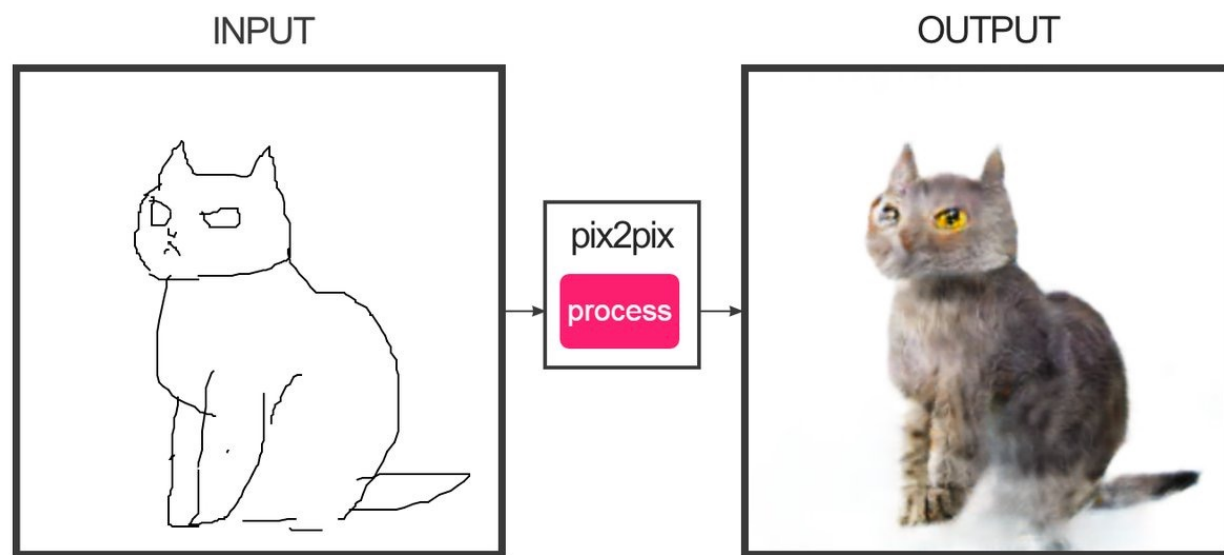


- Rather than penalizing if output *image* looks fake, penalize if each overlapping *patches* looks fake
- Focus on local visual cues (color, textures).
- Global structure: the input image has already encoded global structure. L1 loss can help as well.

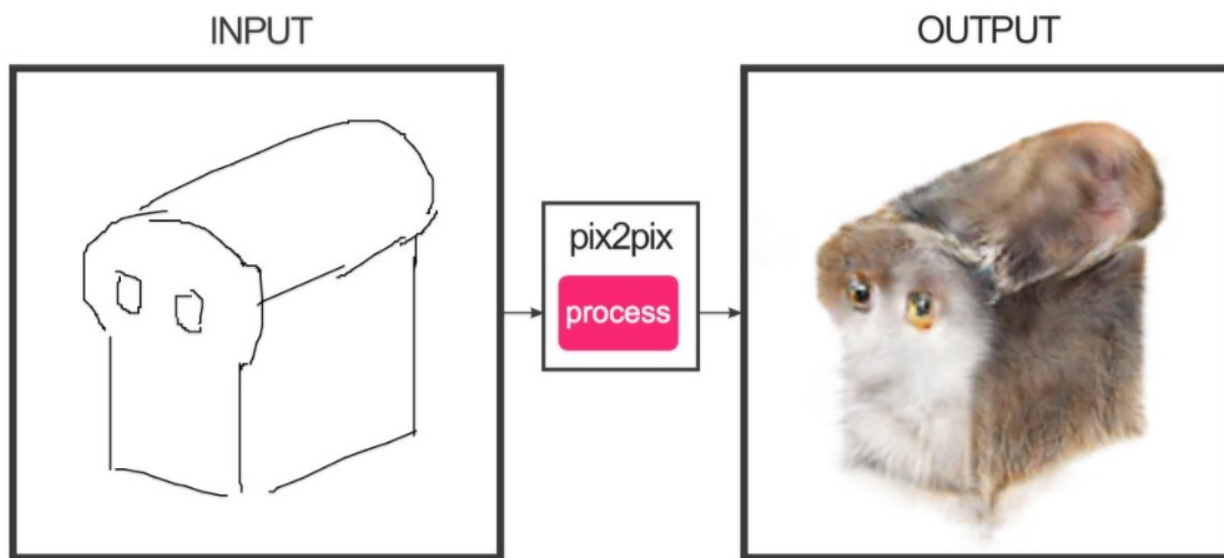
Advantages

- Faster, fewer parameters
- More supervised observations
- Applies to arbitrarily large images

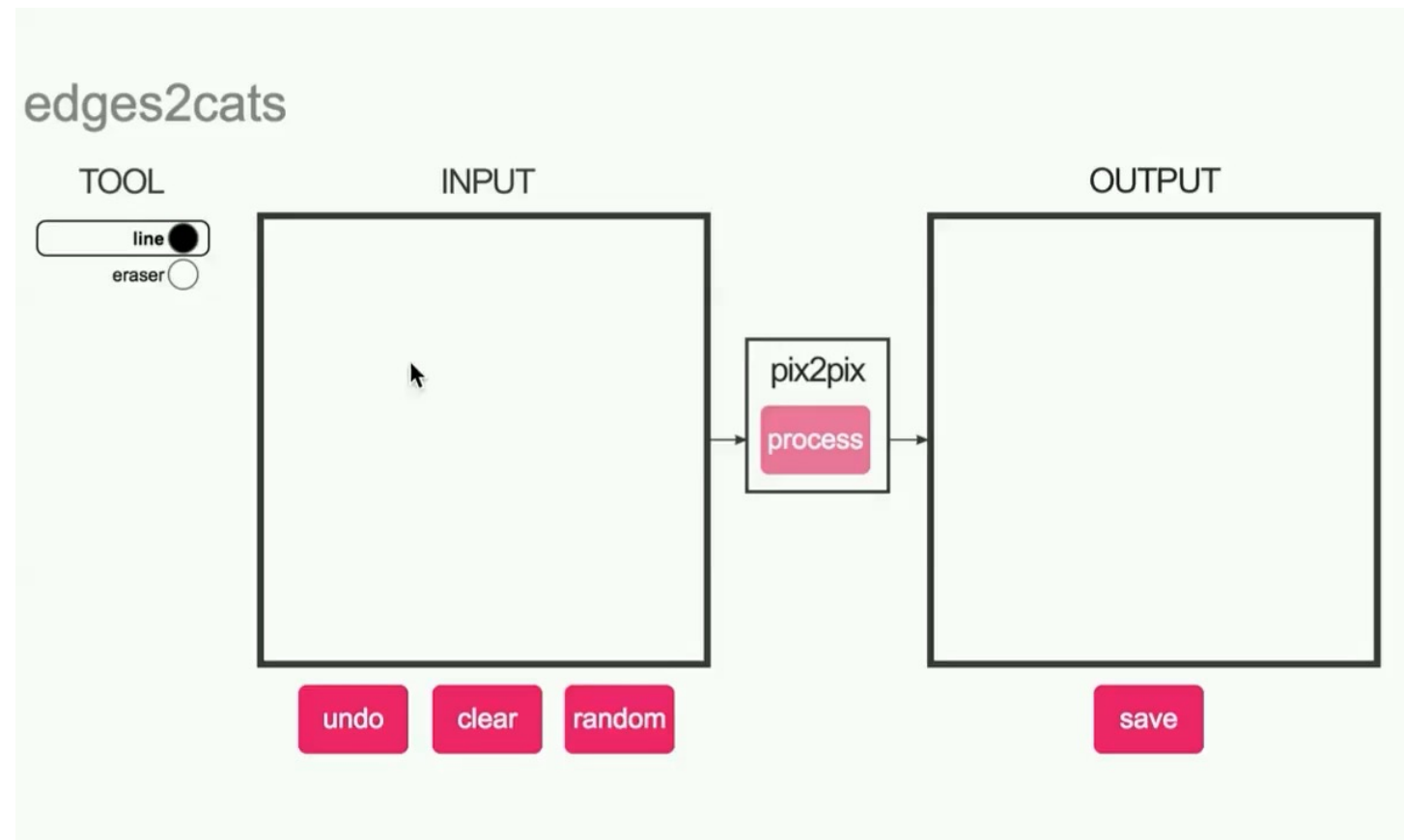
#edges2cats [Christopher Hesse]



@gods_tail



Ivy Tasi @ivymyt

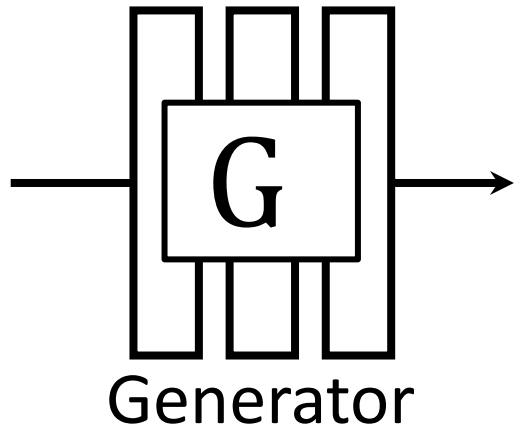


@matthematician

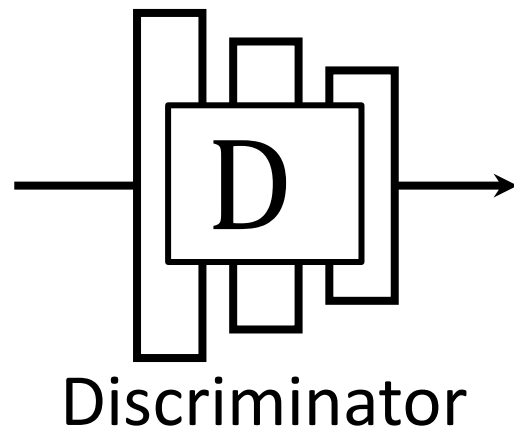


Vitaly Vidmirov @vvid

X



G(x)



Real or fake pair ?

Input: ~~Grayscale~~ ~~Output: Photo~~ → Output: ~~Color~~

Automatic Colorization with pix2pix

Input

Output

Input

Output

Input

Output



Automatic Colorization with pix2pix

Input

Output

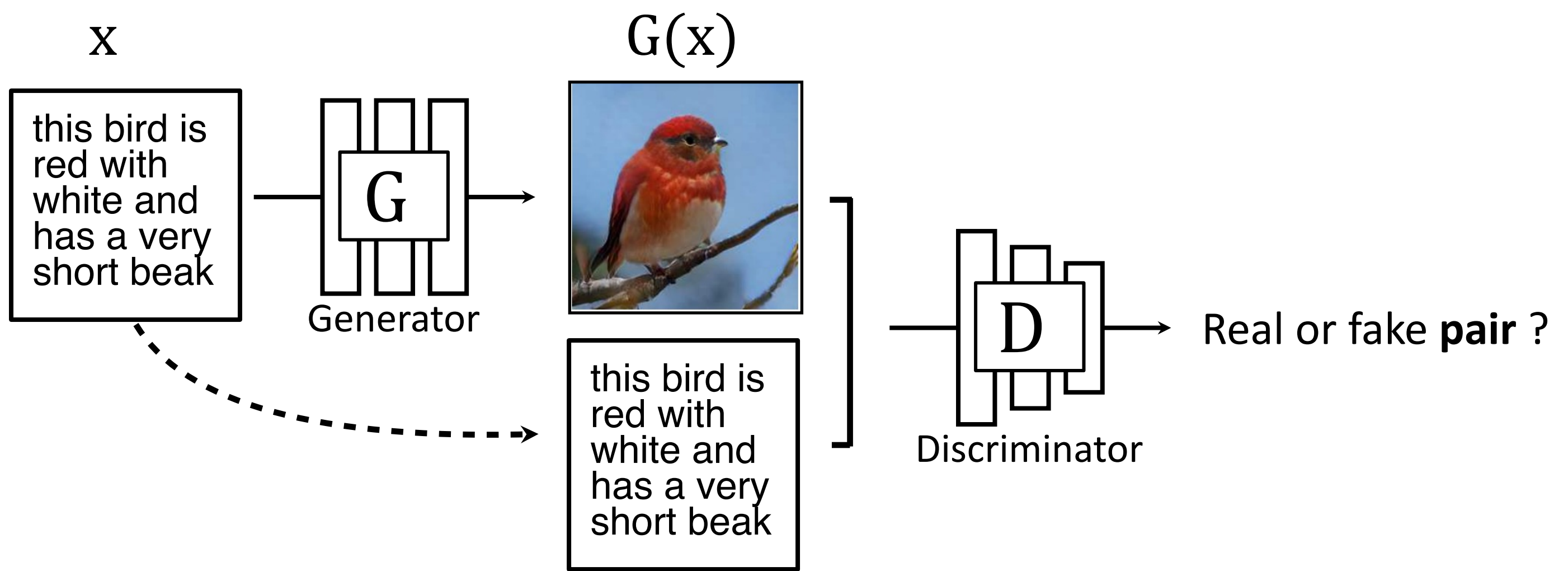
Input

Output

Input

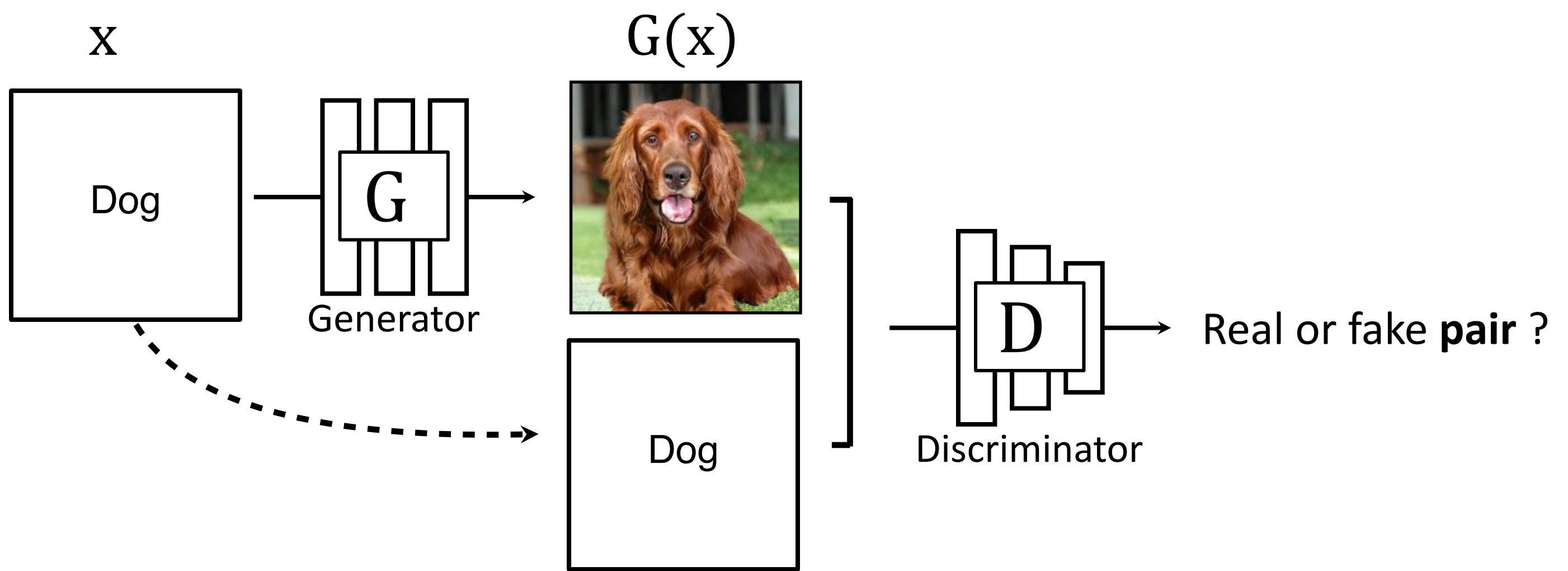
Output





Input: **Text** → Output: **Photo**

Text-to-Image Synthesis



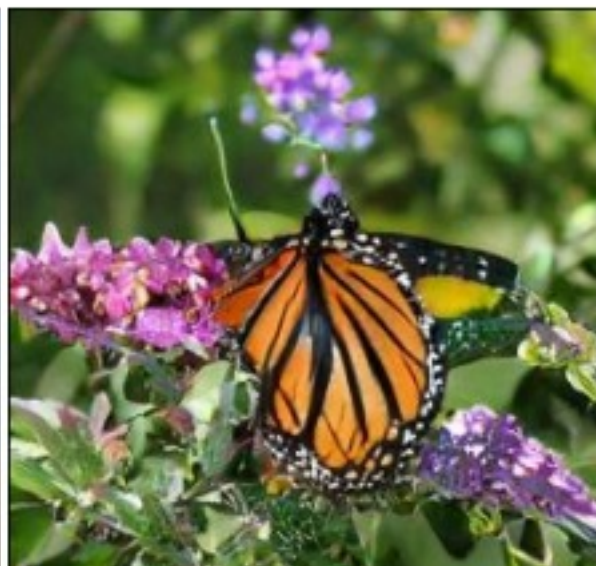
Input: **Class** → Output: **Photo**

Class-conditional GANs

cGANs [Mirza and Osindero. 2014], SAGAN [Zhang et al., 2018], BigGAN [Brock et al., 2019]

StyleGAN-XL [Sauer et al., 2022]

BigGAN

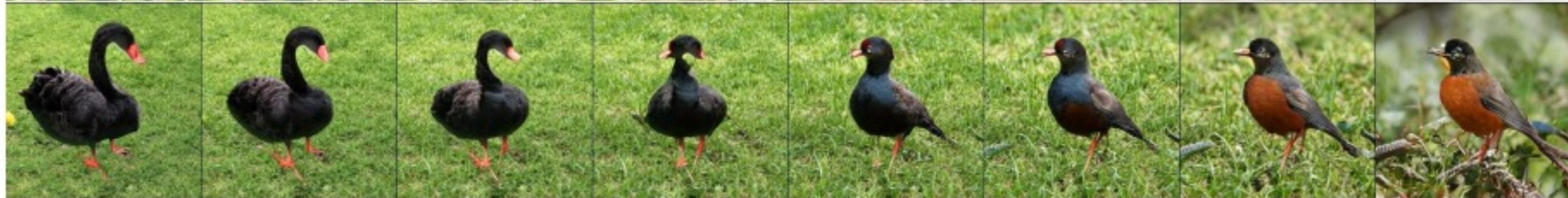


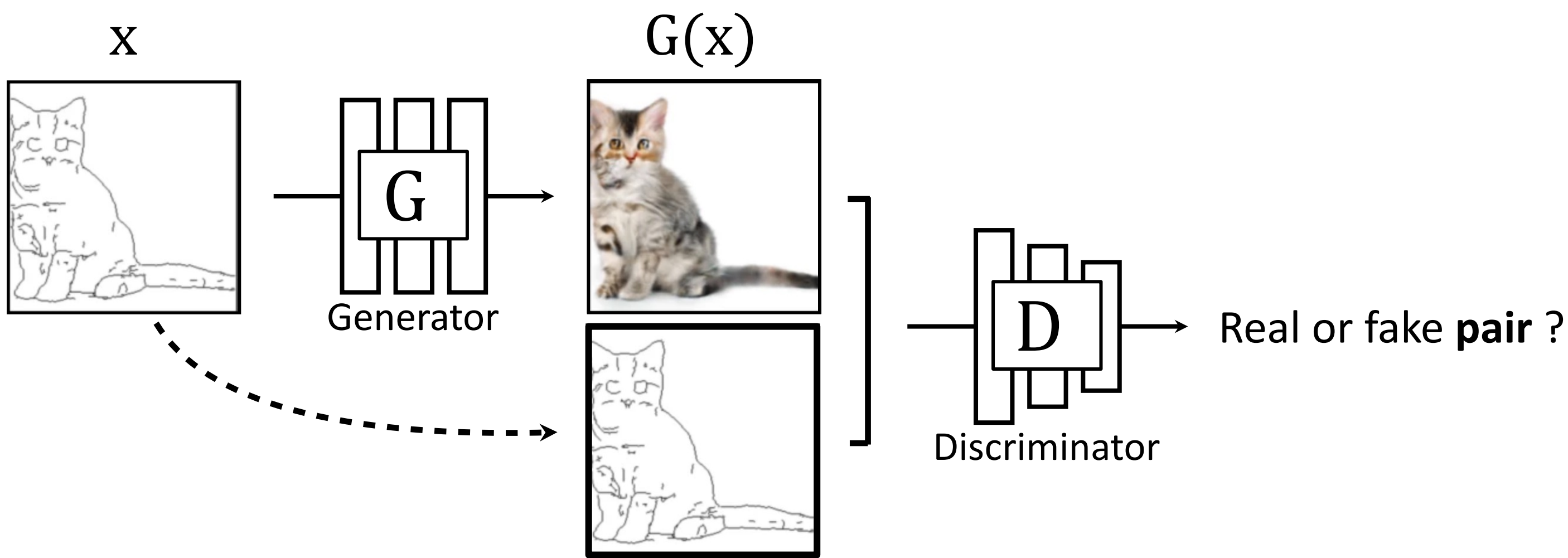
(a) 128×128

(b) 256×256

(c) 512×512

(d)





Learning objective

$$\min_G \max_D \mathbb{E}_x [\log(1 - D(x, G(x)))] + \mathbb{E}_{x,y} [\log D(x, y)]$$

Limitations

- One-to-one mapping.
- Low-resolution output.
- Requires paired training data

Improving Conditional GANs

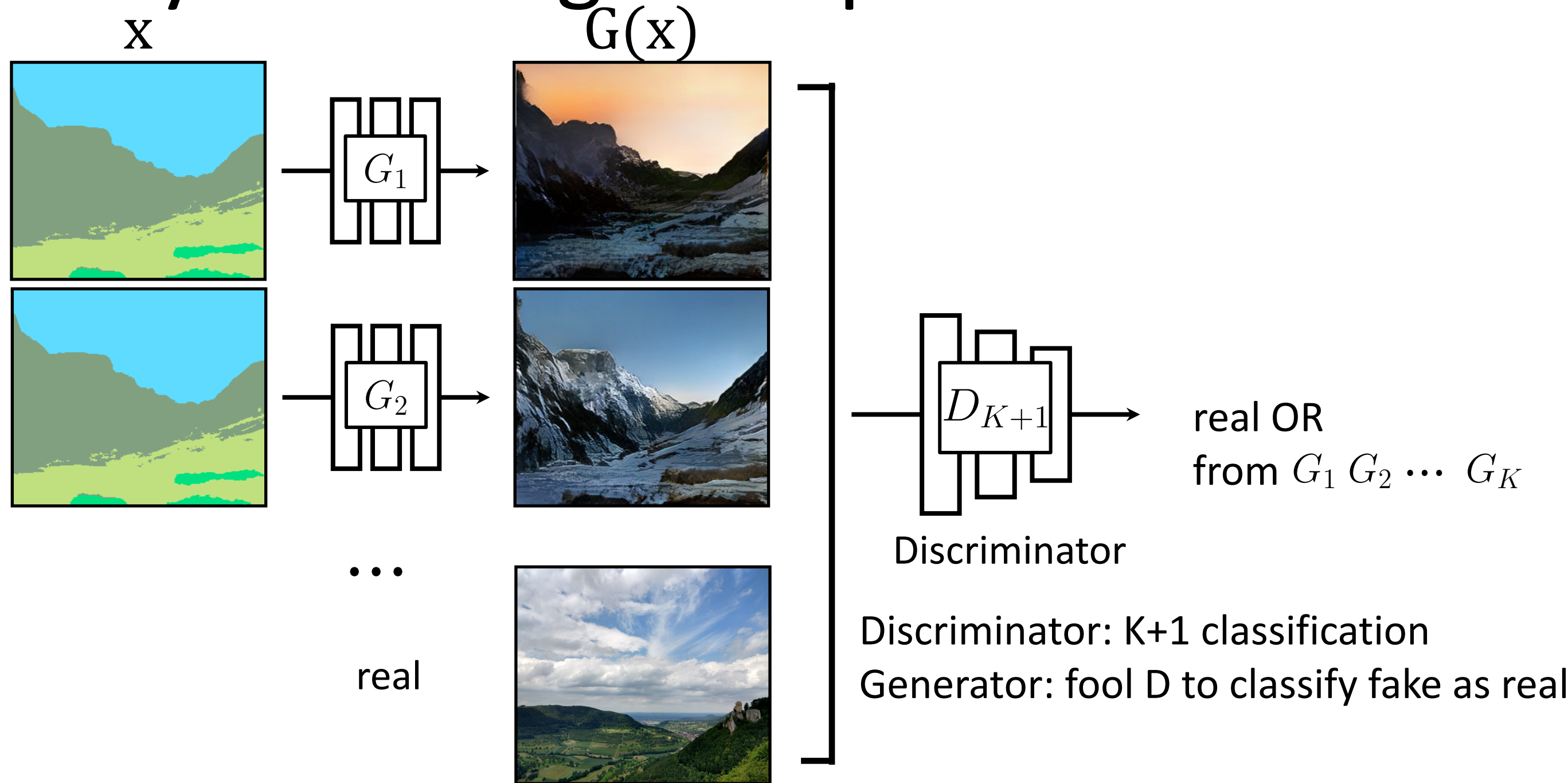
- **Multimodal synthesis.**
- High-resolution synthesis.
- Model training without pairs (next lecture)

Group Discussion



- <https://docs.google.com/forms/d/e/1FAIpQLSdfSXMRLddytfNaDOFAORBaOTxLPQTLWTELxpyAr8NJrFcZhA/viewform?usp=sharing>

Synthesizing Multiple Results



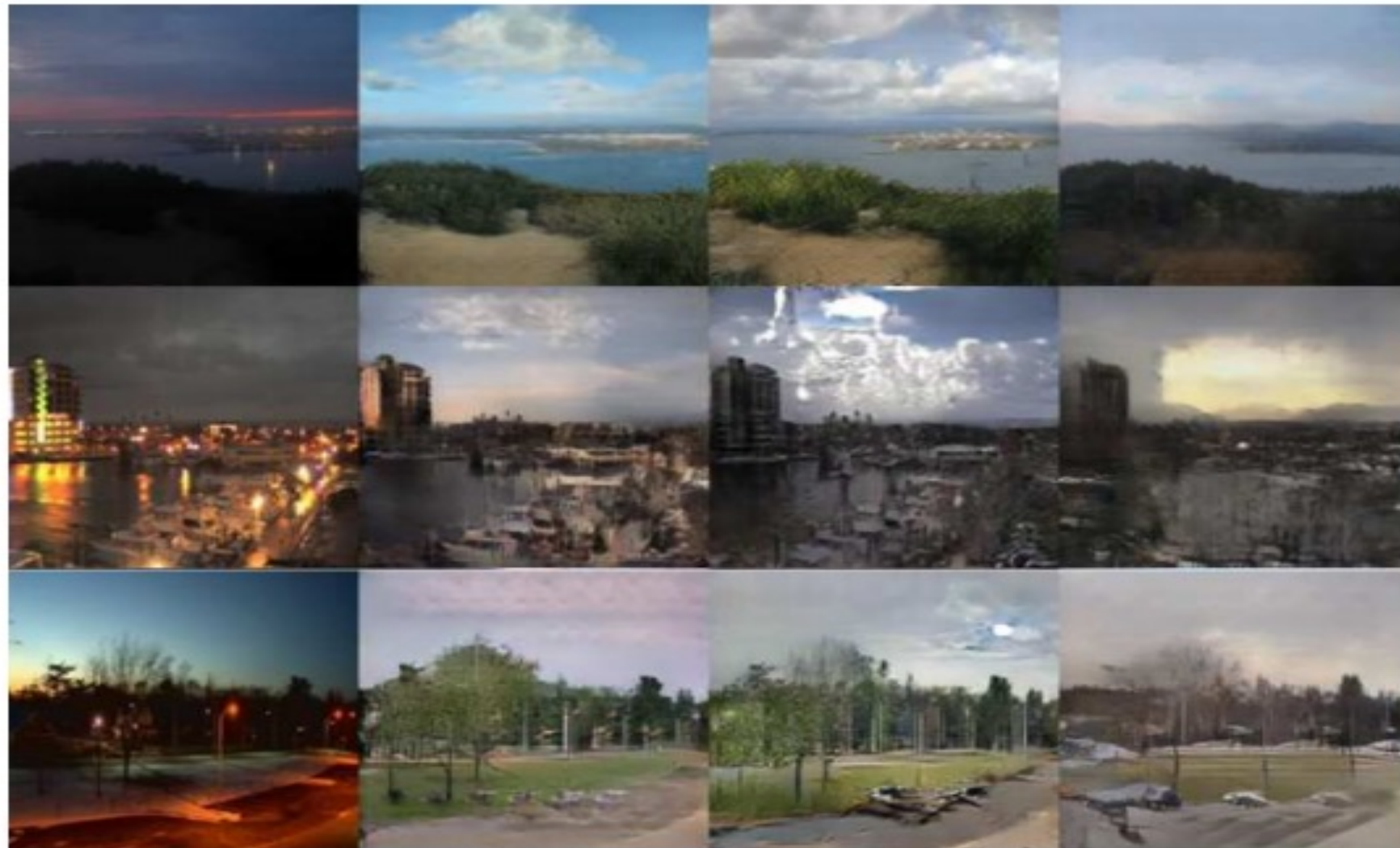
Synthesizing Multiple Results

Night input

Day output 1

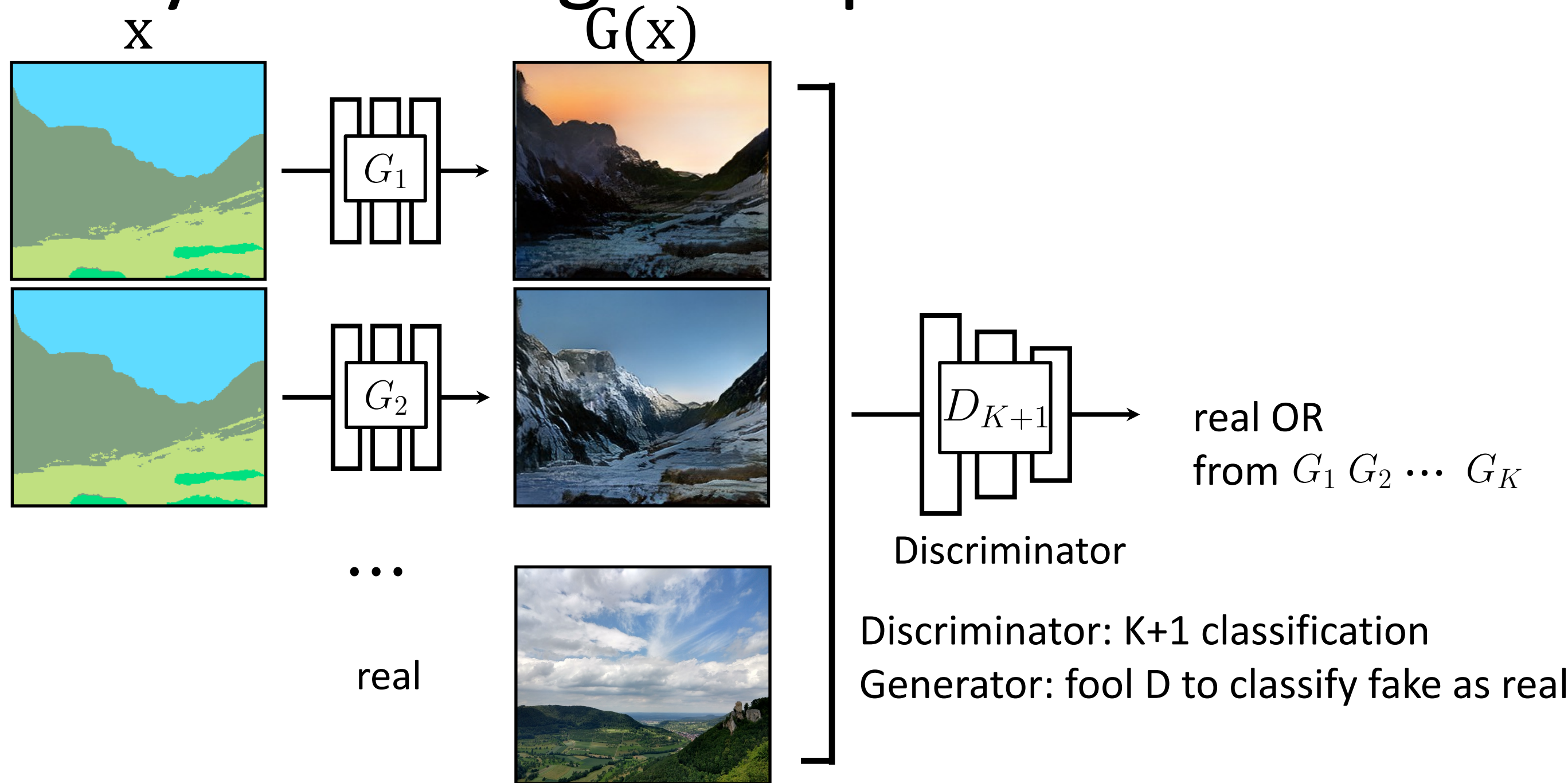
Day output 2

Day output 3



Multi-agent Diverse GANs [Ghosh et al., CVPR 2018]

Synthesizing Multiple Results



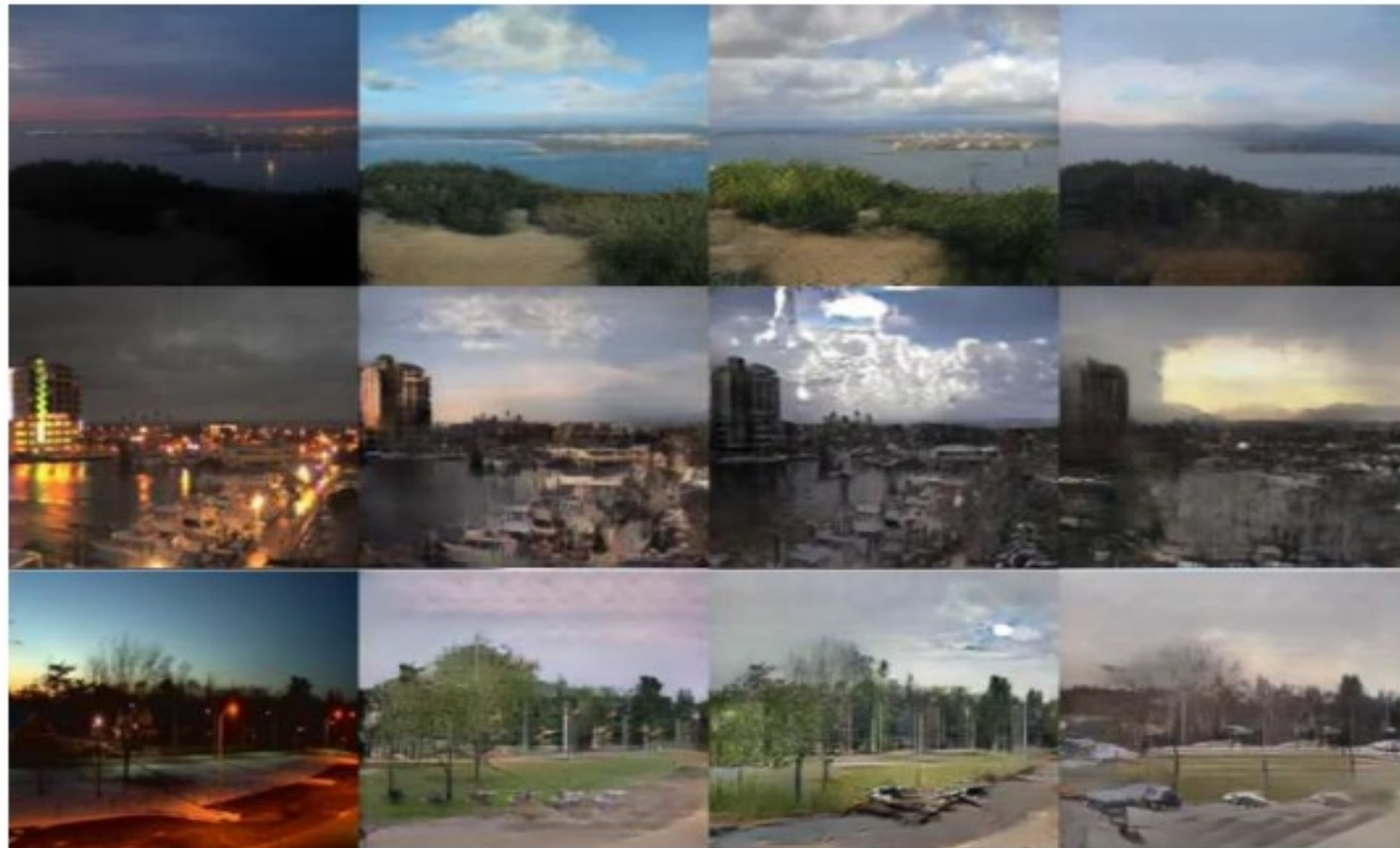
Synthesizing Multiple Results

Night input

Day output 1

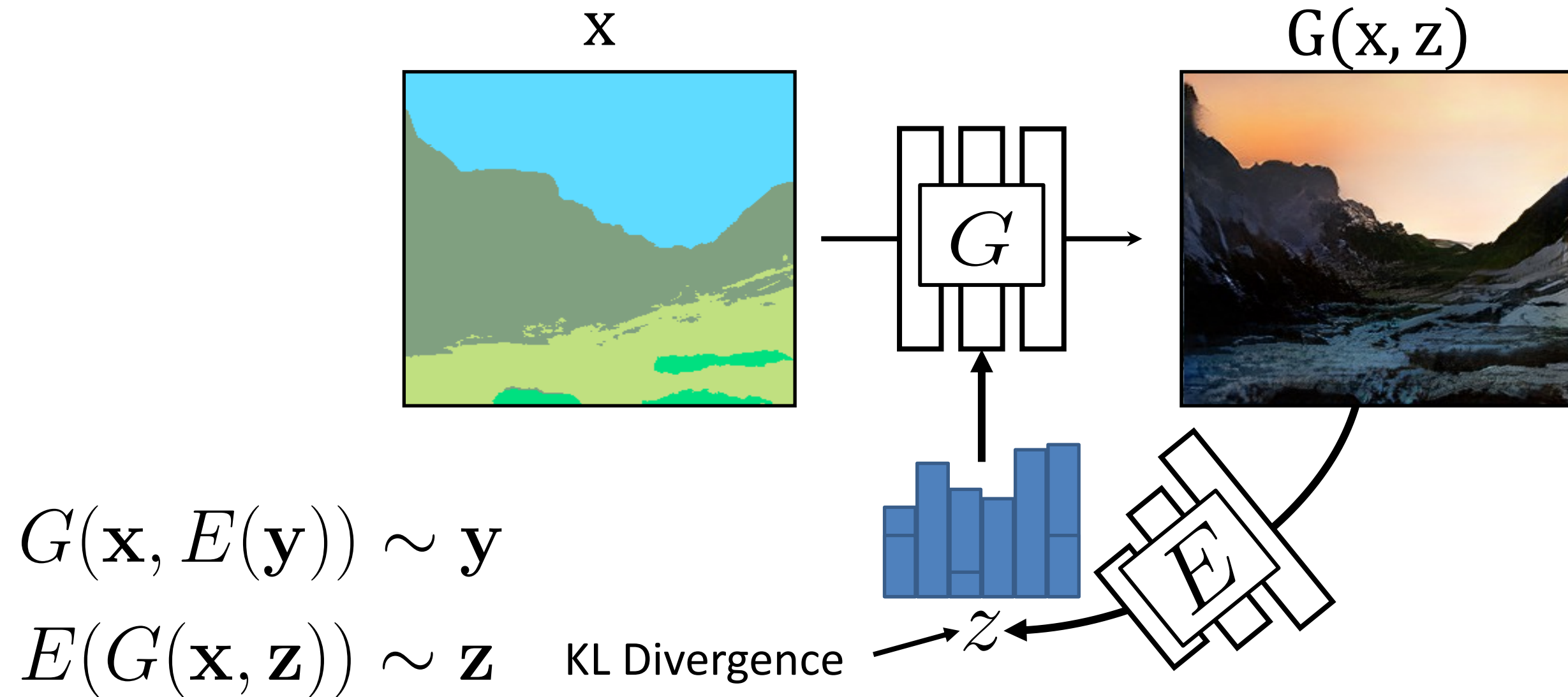
Day output 2

Day output 3



Multi-agent Diverse GANs [Ghosh et al., CVPR 2018]

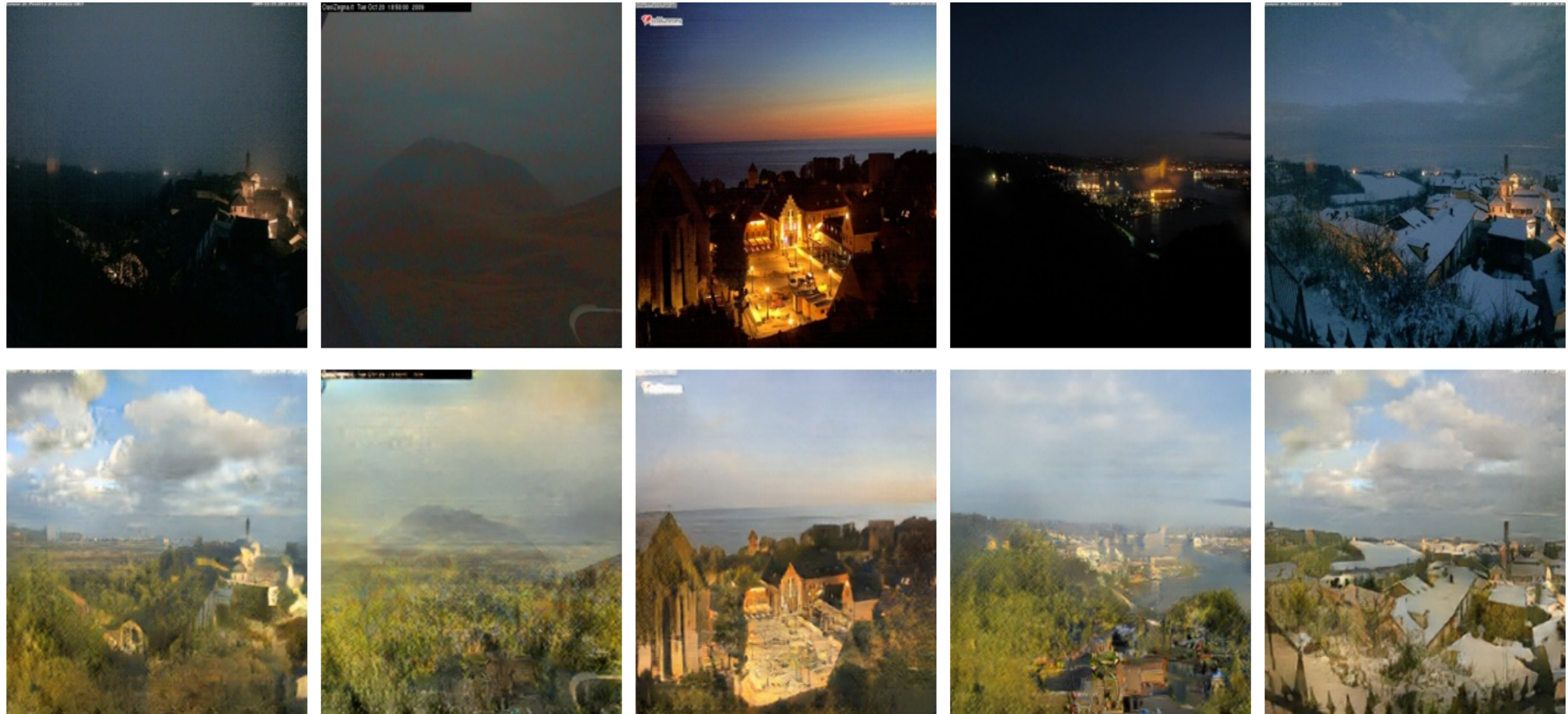
Synthesizing Multiple Results



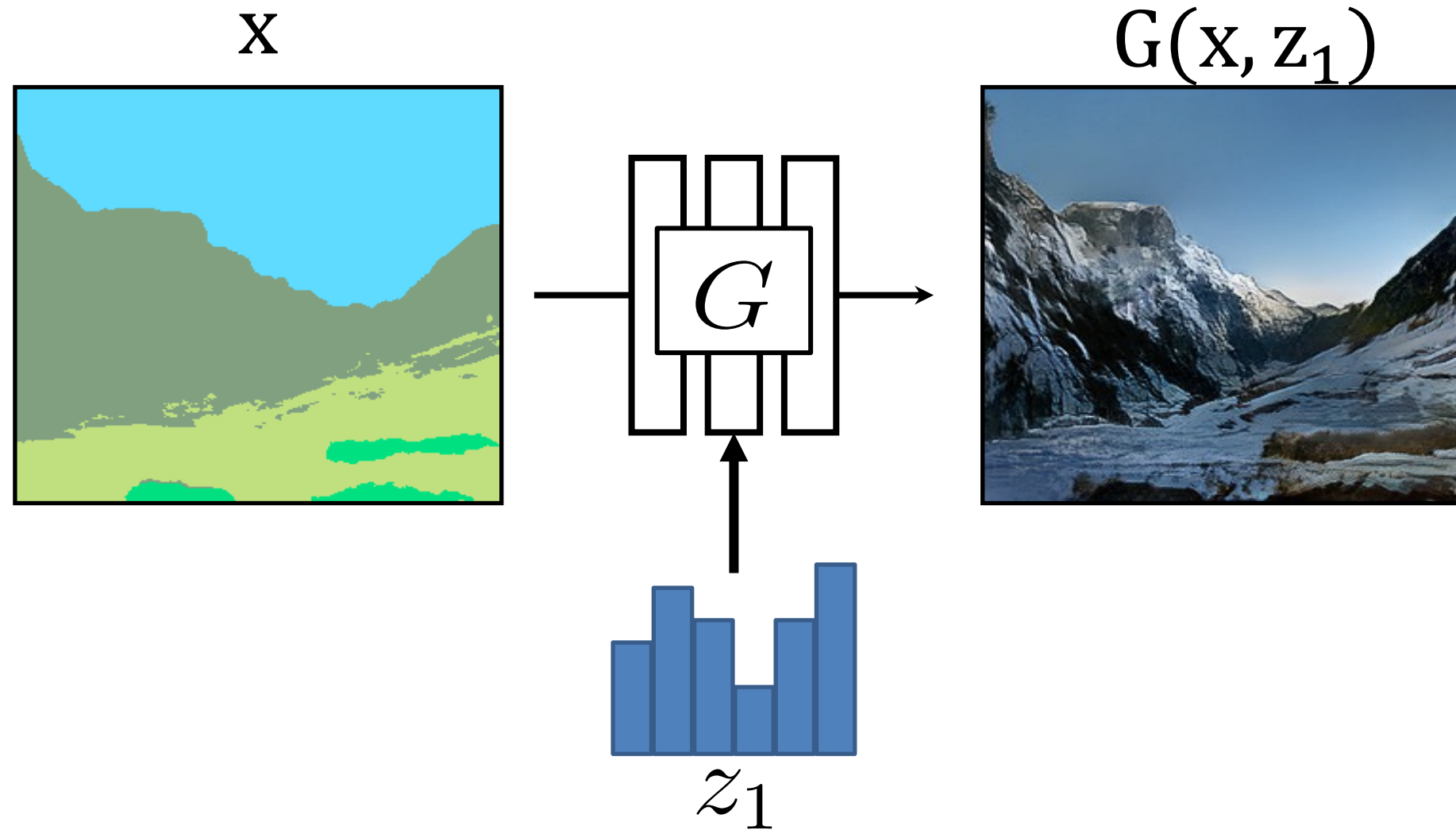
Synthesizing Multiple Results



Synthesizing Multiple Results



Synthesizing Multiple Results

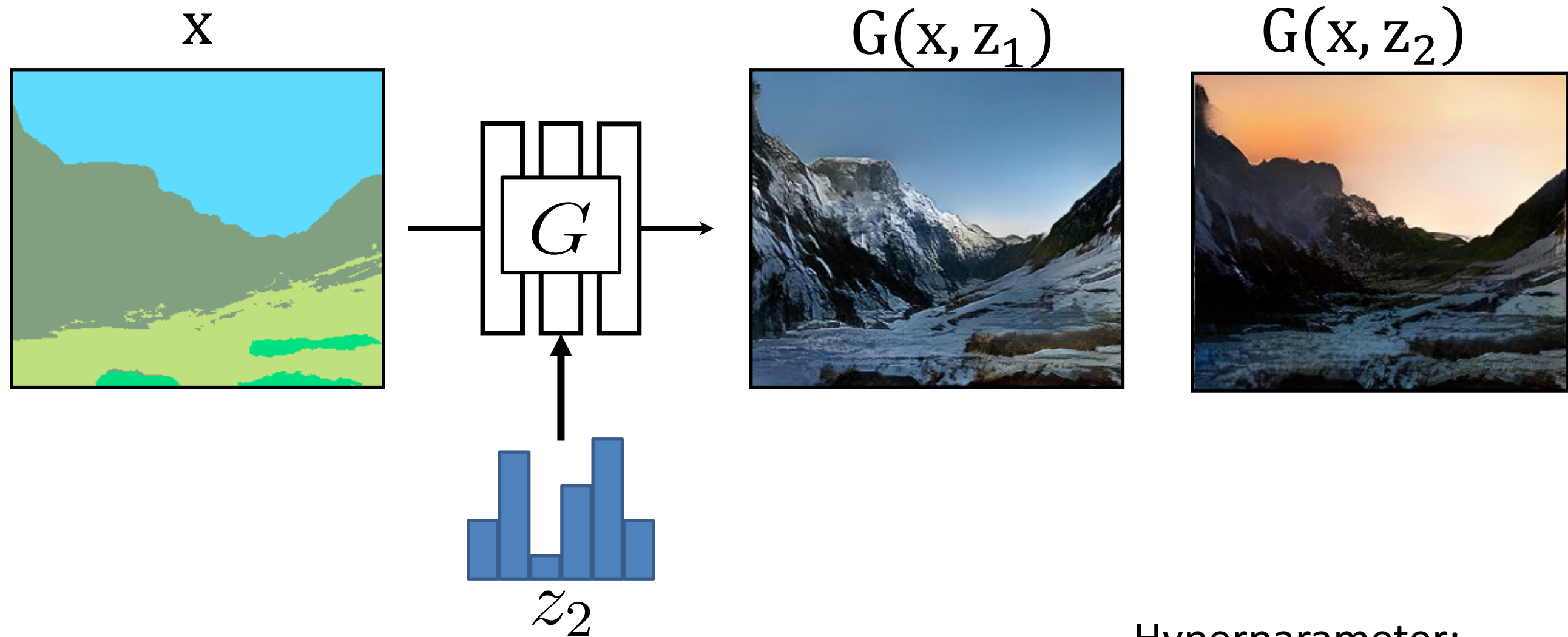


$$\max_G \mathcal{L}_z(G) = \mathbb{E}_{z_1, z_2} \left[\min \left(\frac{\|G(x, z_1) - G(x, z_2)\|}{\|z_1 - z_2\|}, \tau \right) \right],$$

Hyperparameter:
Degree of diversity

Diversity-Sensitive GAN [Yang et al., 2019]

Synthesizing Multiple Results



$$\max_G \mathcal{L}_z(G) = \mathbb{E}_{z_1, z_2} \left[\min \left(\frac{\|G(x, z_1) - G(x, z_2)\|}{\|z_1 - z_2\|}, \tau \right) \right]$$

Hyperparameter:
Degree of diversity

Diversity-Sensitive GAN [Yang et al., 2019]

Synthesizing Multiple Results



$$\max_G \mathcal{L}_z(G) = \mathbb{E}_{z_1, z_2} \left[\min \left(\frac{\|G(\mathbf{x}, z_1) - G(\mathbf{x}, z_2)\|}{\|z_1 - z_2\|}, \tau \right) \right]$$

Hyperparameter:
Degree of diversity

Diversity-Sensitive GAN [Yang et al., 2019]

Improving Conditional GANs

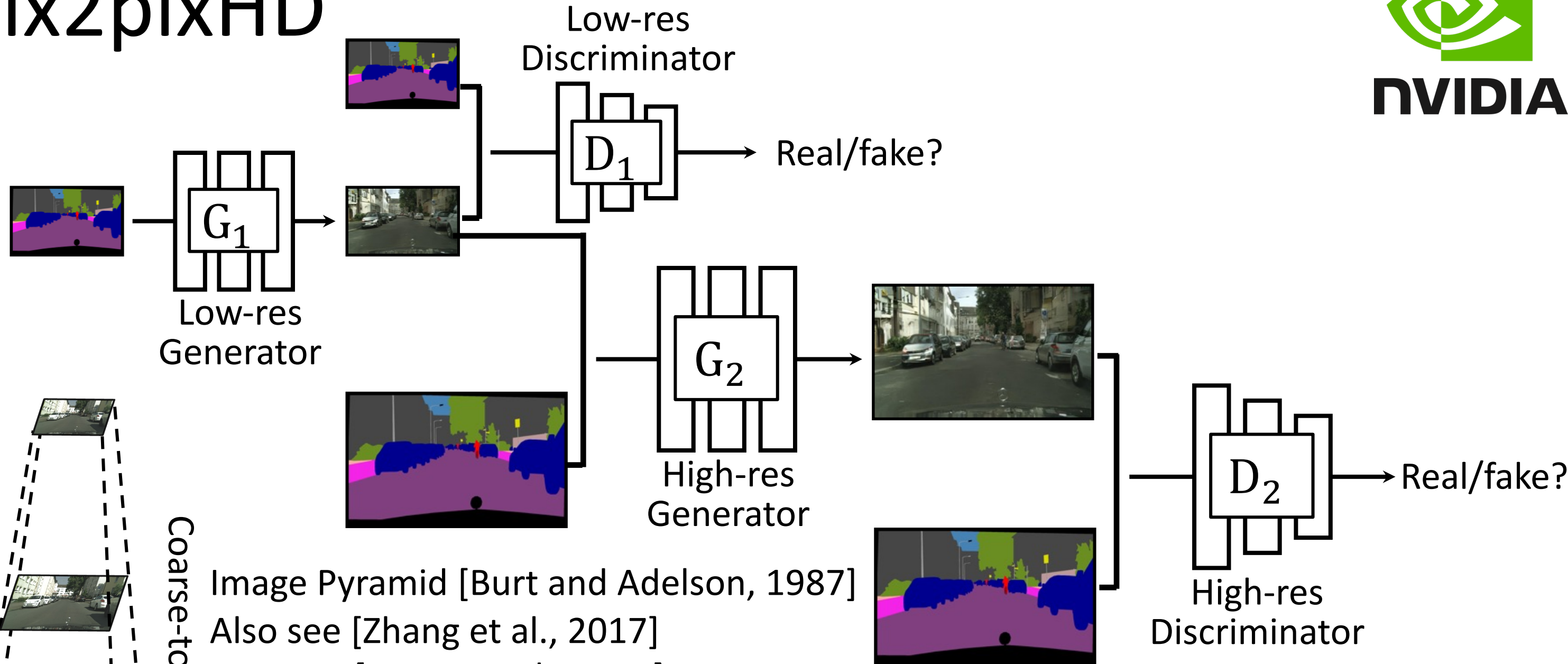
- Multimodal synthesis.
- **High-resolution synthesis.**
- Model training without pairs (next lecture)

The Curse of Dimensionality



Pix2pix output

pix2pixHD

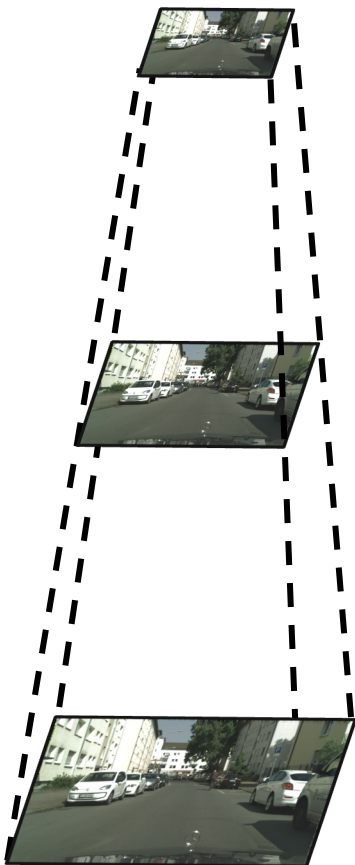


Coarse-to-fine

Image Pyramid [Burt and Adelson, 1987]
Also see [Zhang et al., 2017]
[Karras et al., 2018]

Objective: Multi-scale GANs loss + Perceptual Loss
+ Feature Matching Loss (with Discriminator's features)

pix2pixHD [Wang et al., 2018]



pix2pixHD: 2048×1024



Style

Label

Stroke

Possible Styles



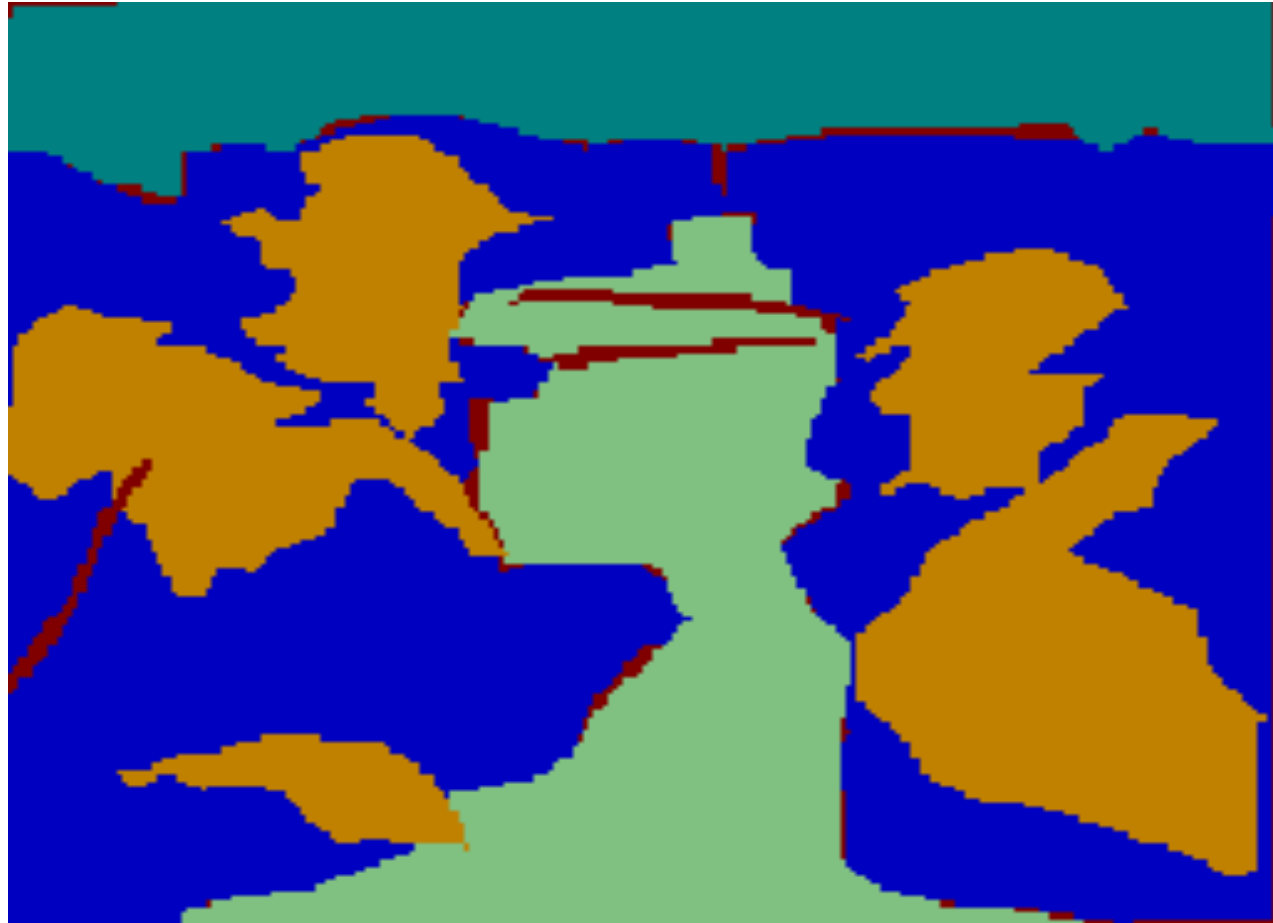
Label Map



Synthesized Result

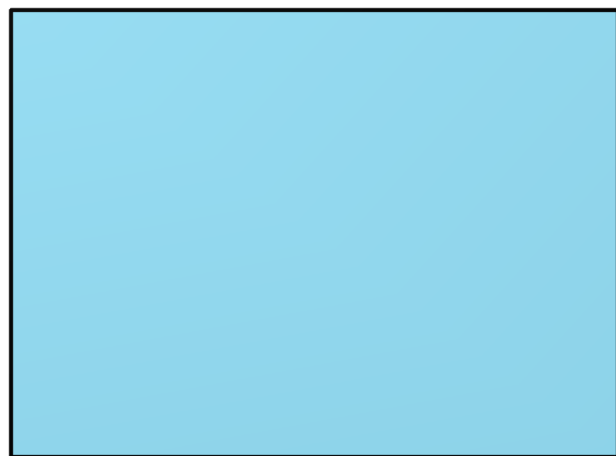


Conditional Image Synthesis in the Wild

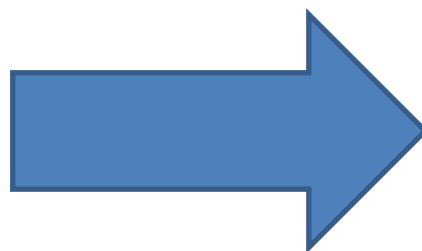


pix2pixHD [Wang et al., 2018]

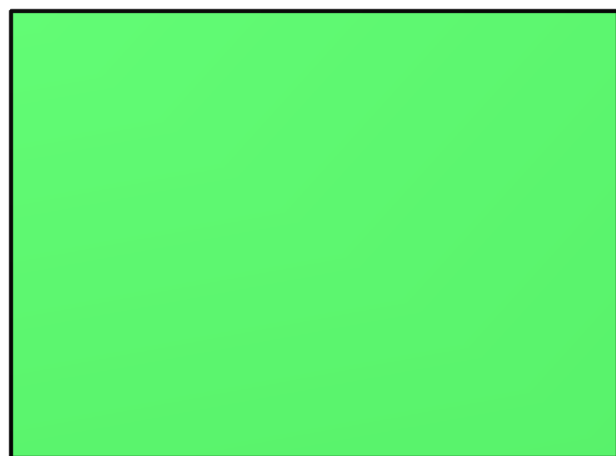
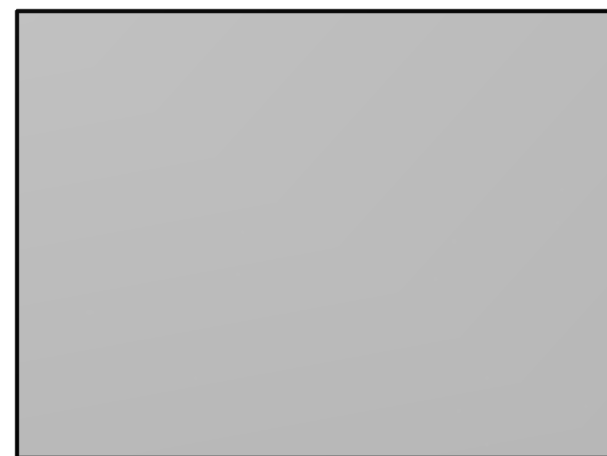
input



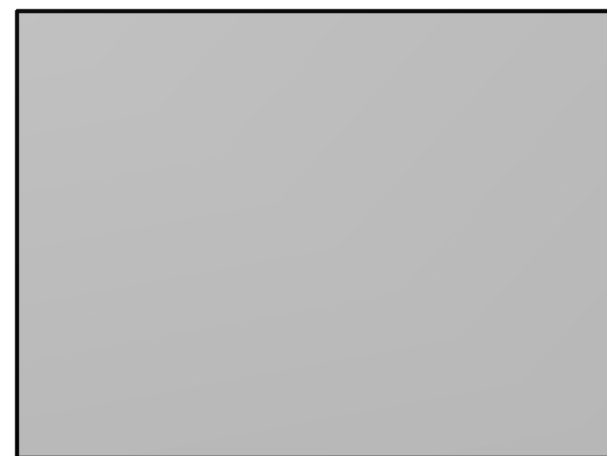
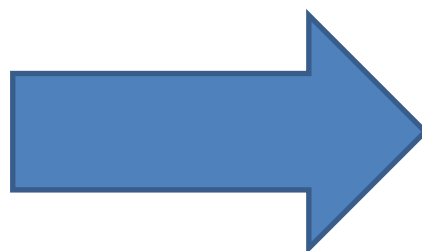
sky



output

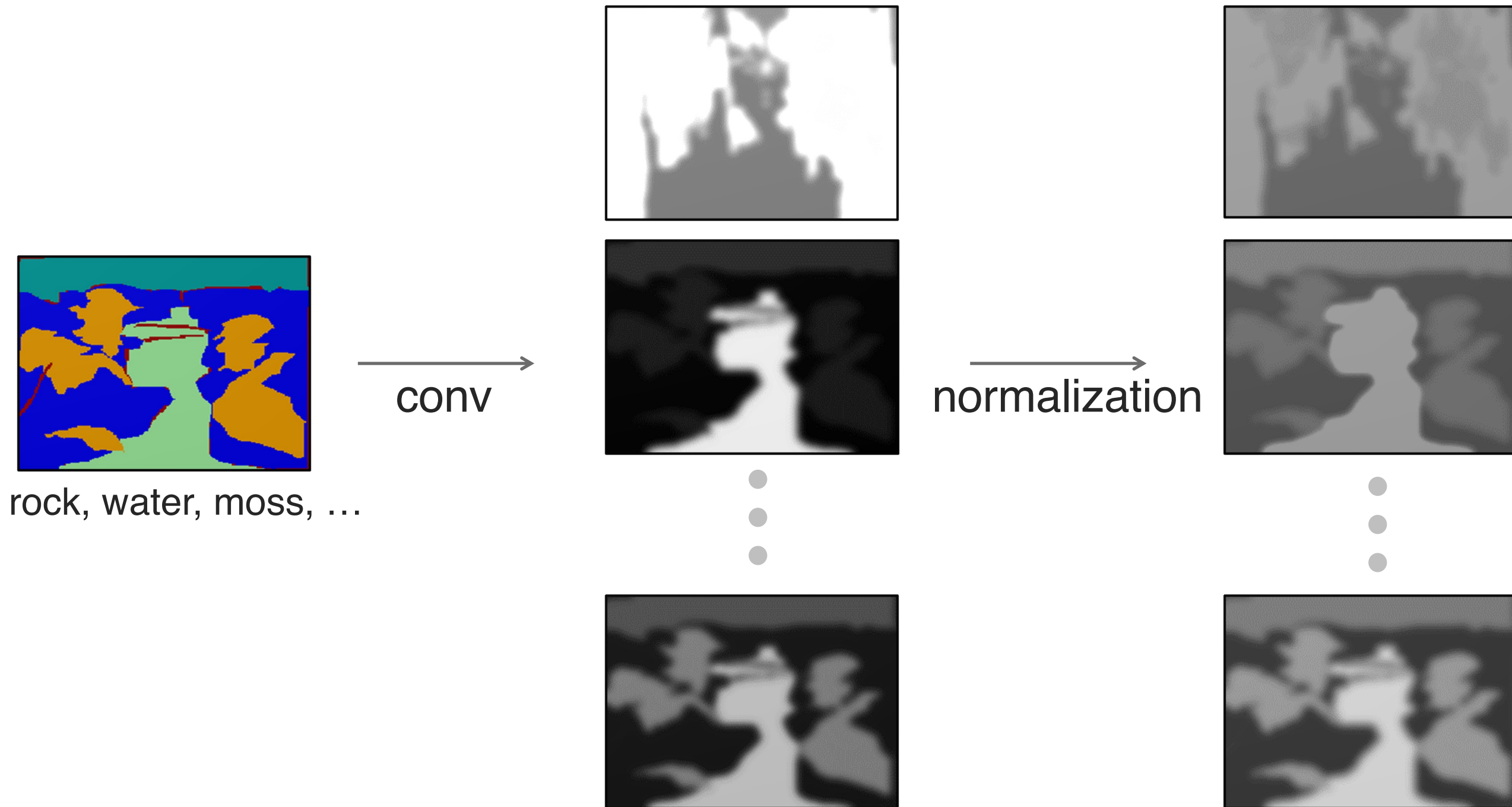


grass

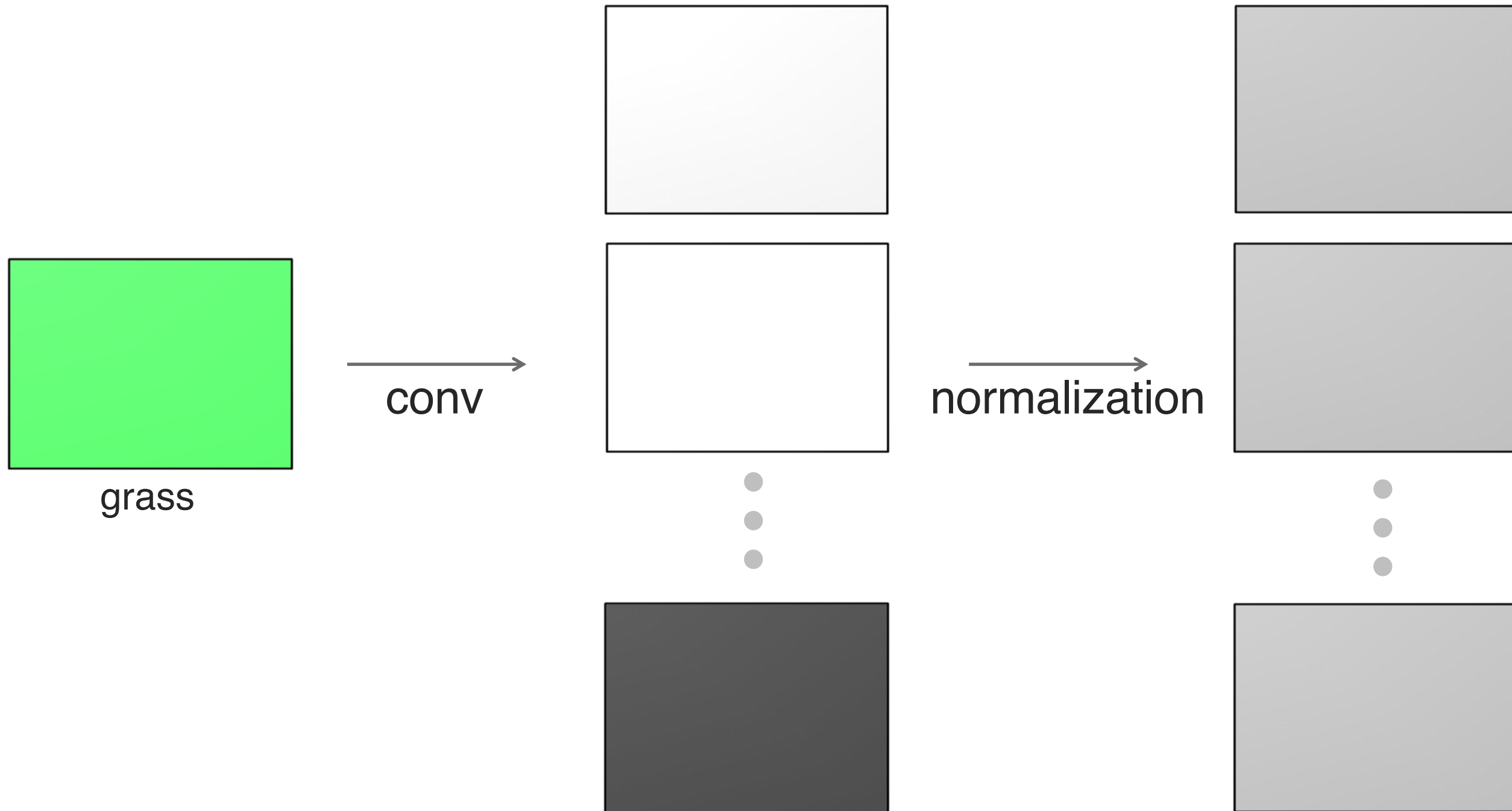


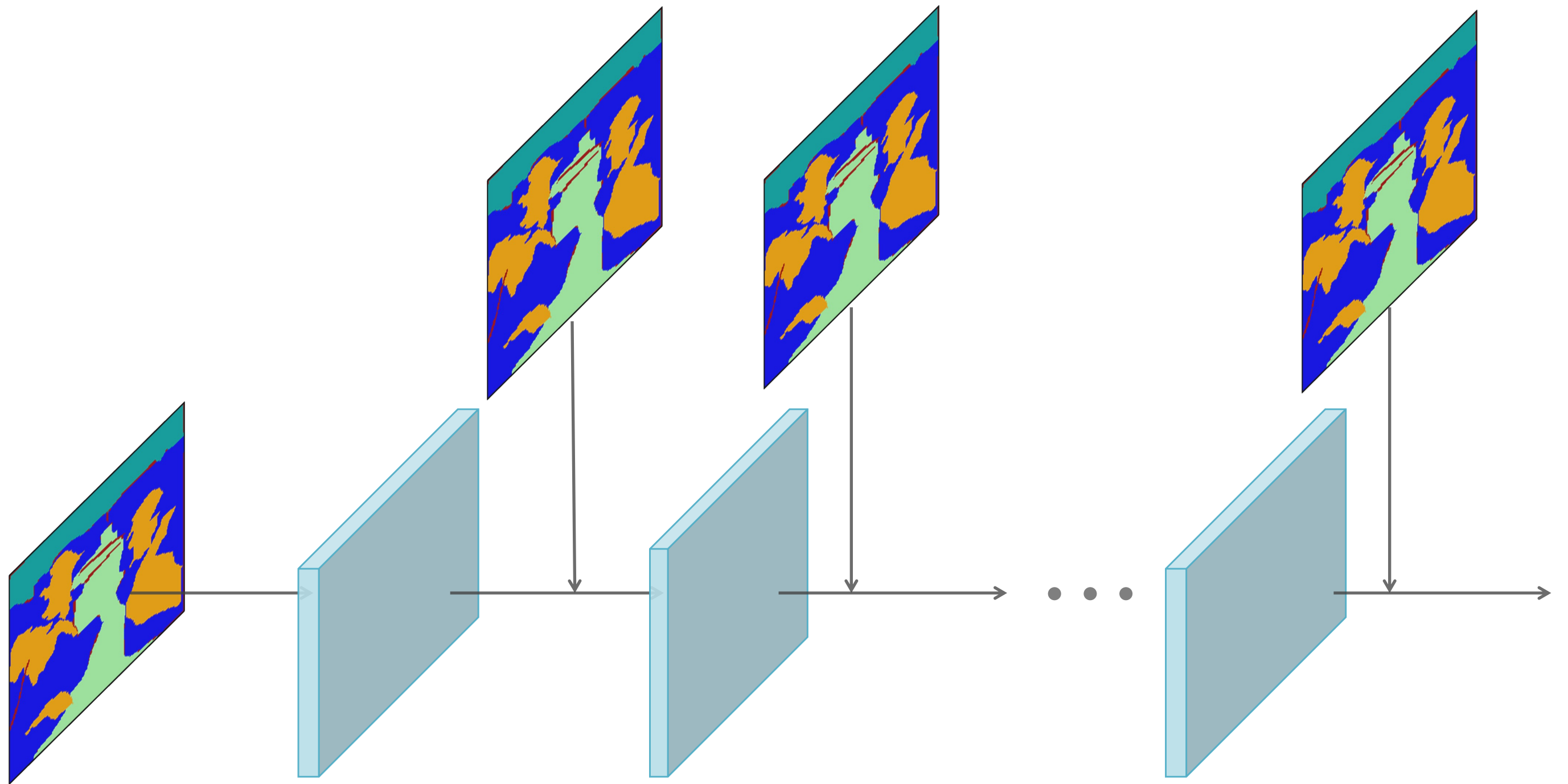
pix2pixHD [Wang et al., 2018]

Problem with standard networks



Problem with standard networks



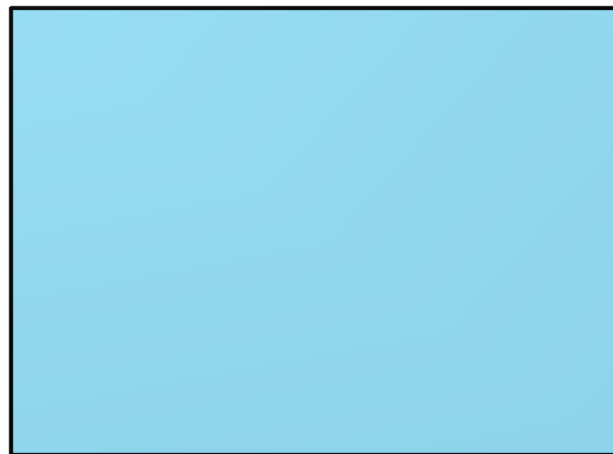


normalization

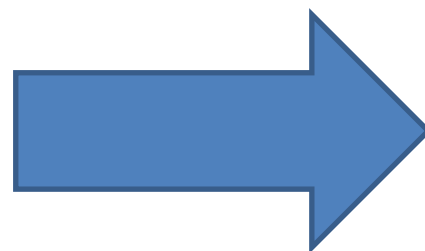
normalization

normalization

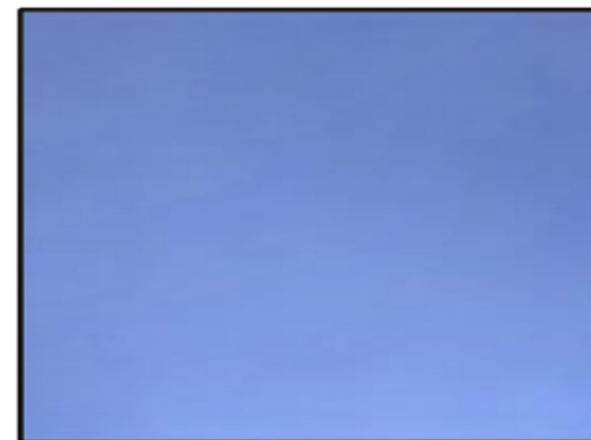
input



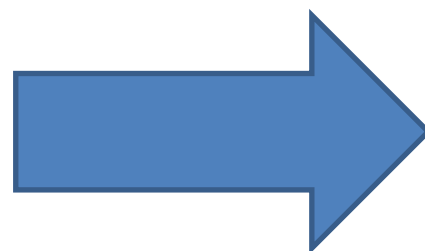
sky



output

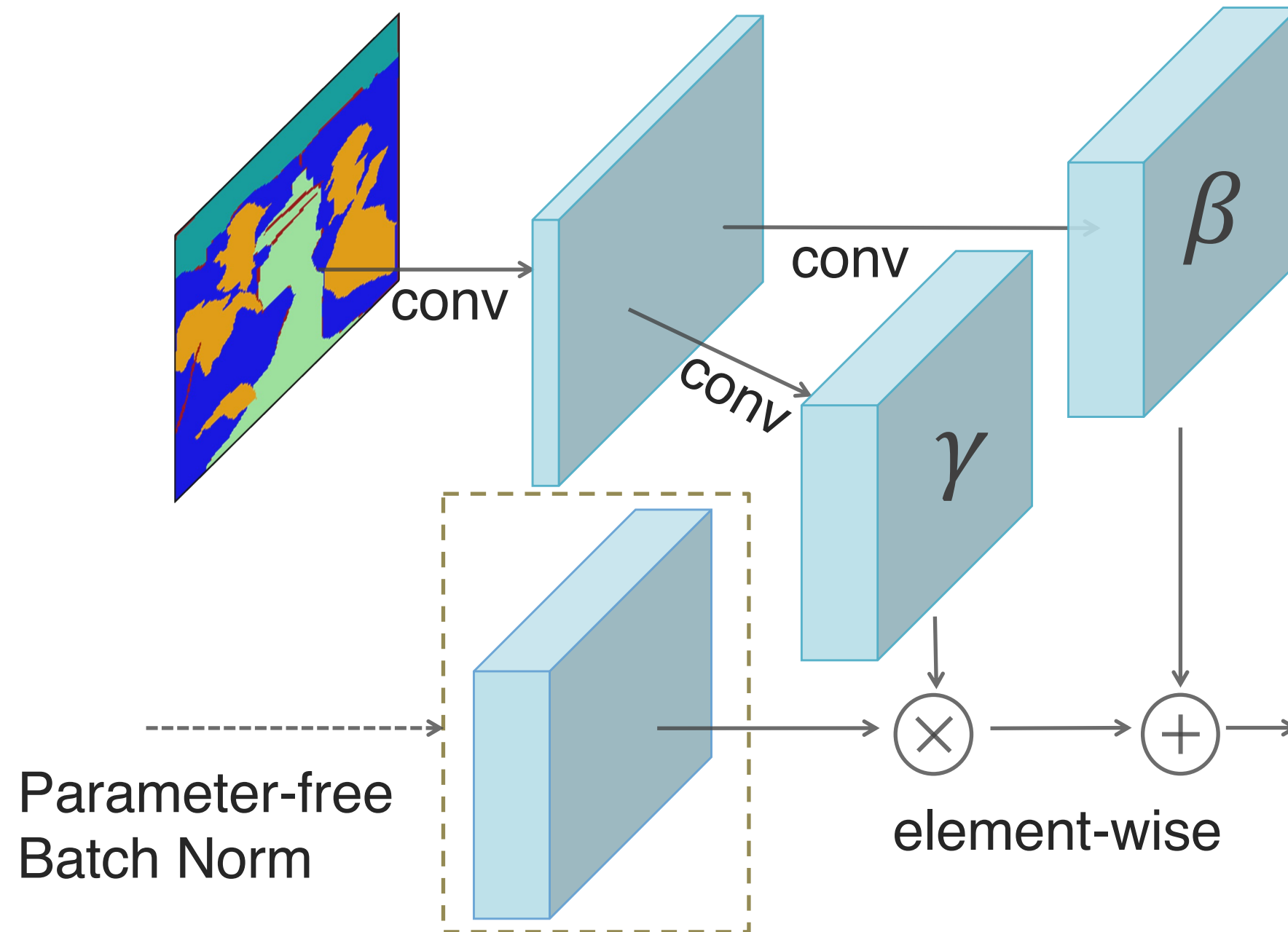


grass



SPADE (ours)

SPADE (SPAtially ADaptive DEnormalization)



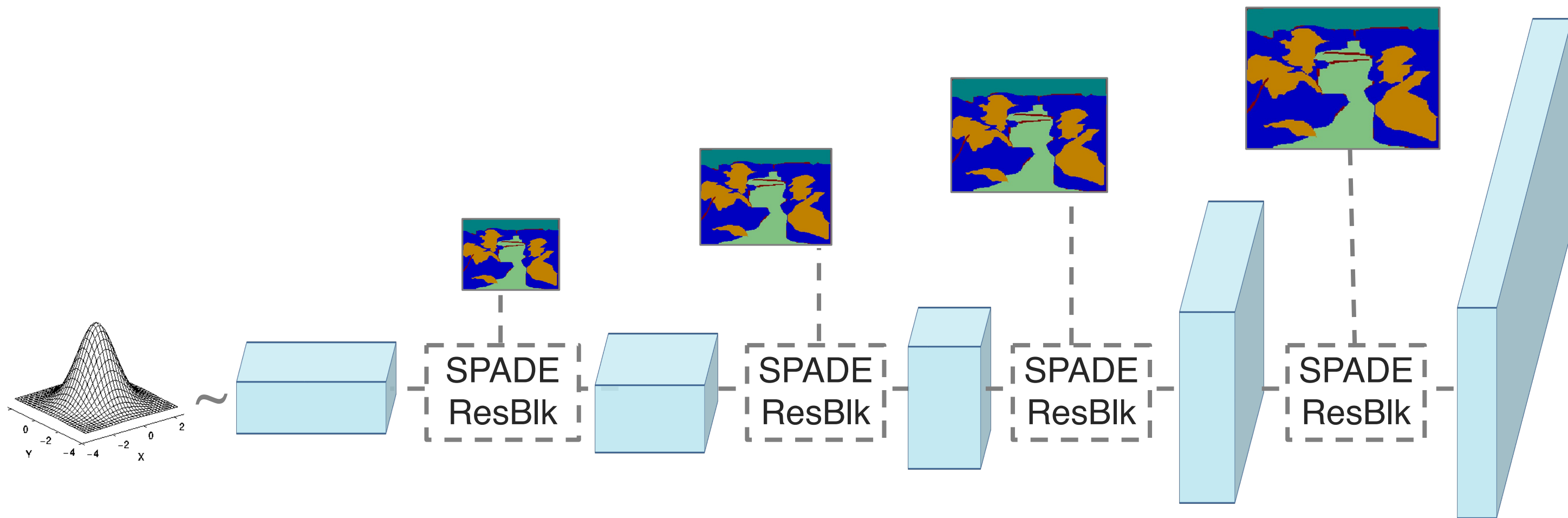
SPADE (SPAtially ADaptive DENormalization)

Batch Norm (Ioffe et al. 2015)

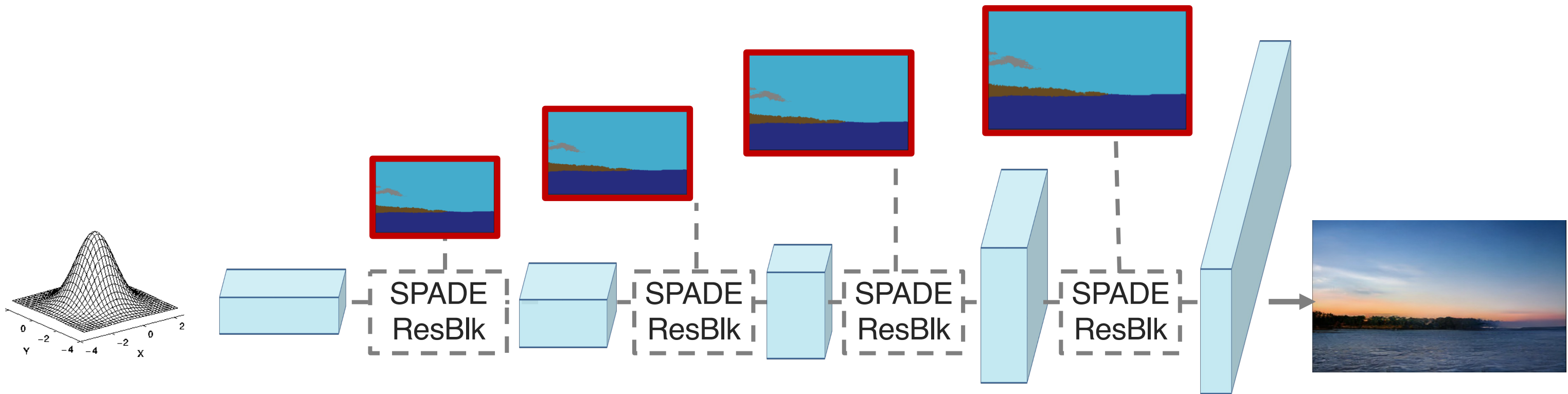
$$y = \underbrace{\frac{x - \mu}{\sigma}}_{\text{normalization}} \cdot \underbrace{\gamma + \beta}_{\text{affine transform}}$$

See other adaptive/conditional normalization: conditional BN (Dumoulin et al.), AdaIN (Huang and Belongie), SFT (Wang et al.)

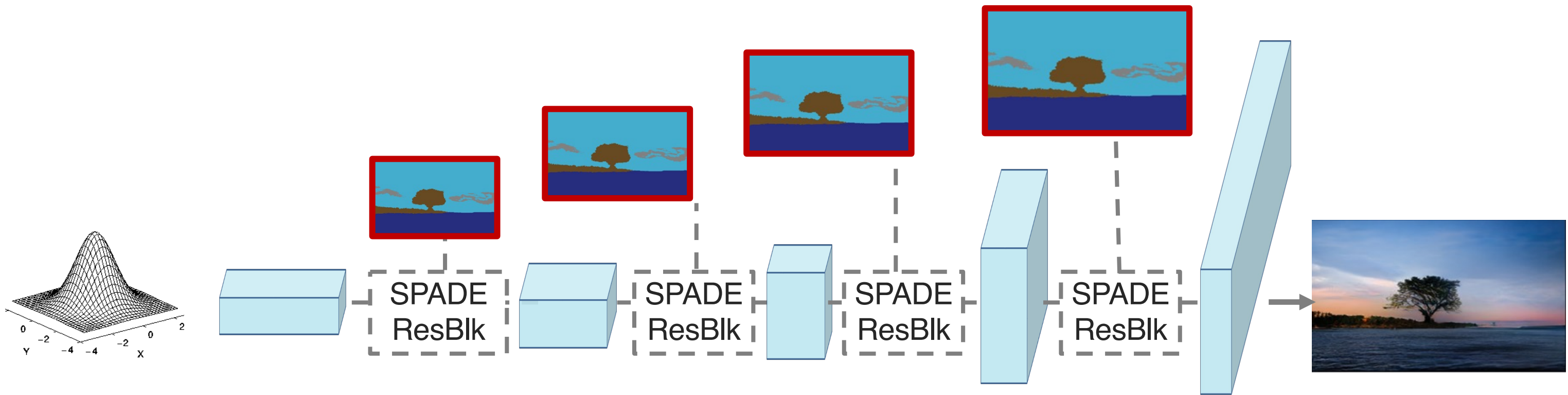
Generator



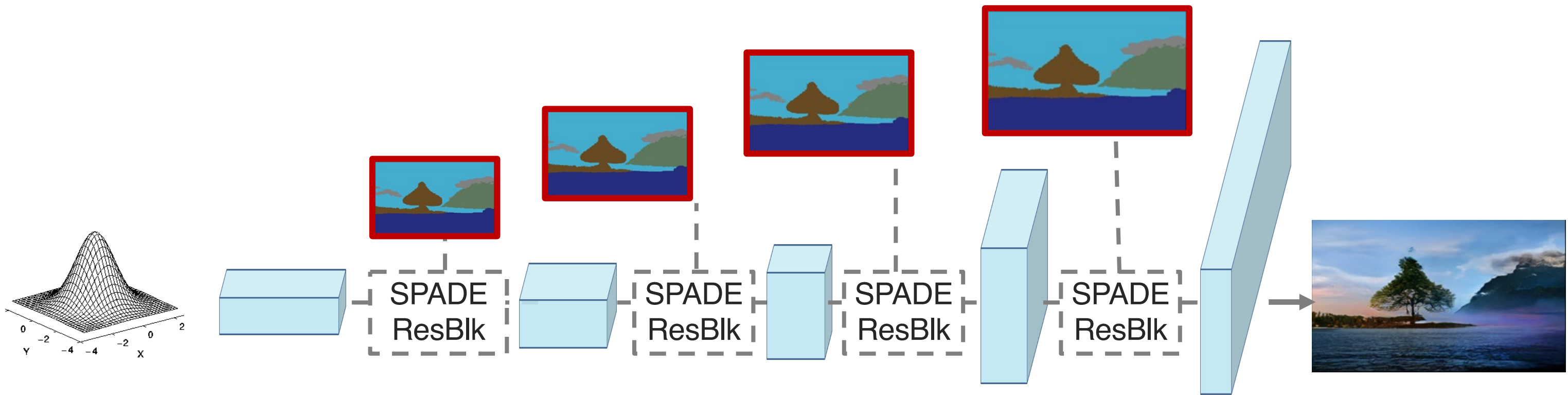
Semantic Control



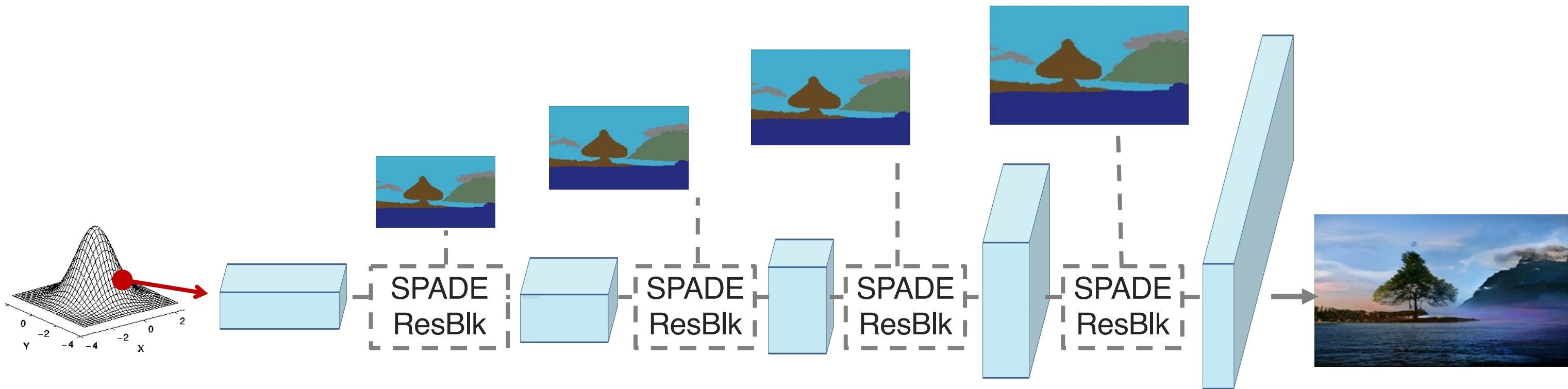
Semantic Control



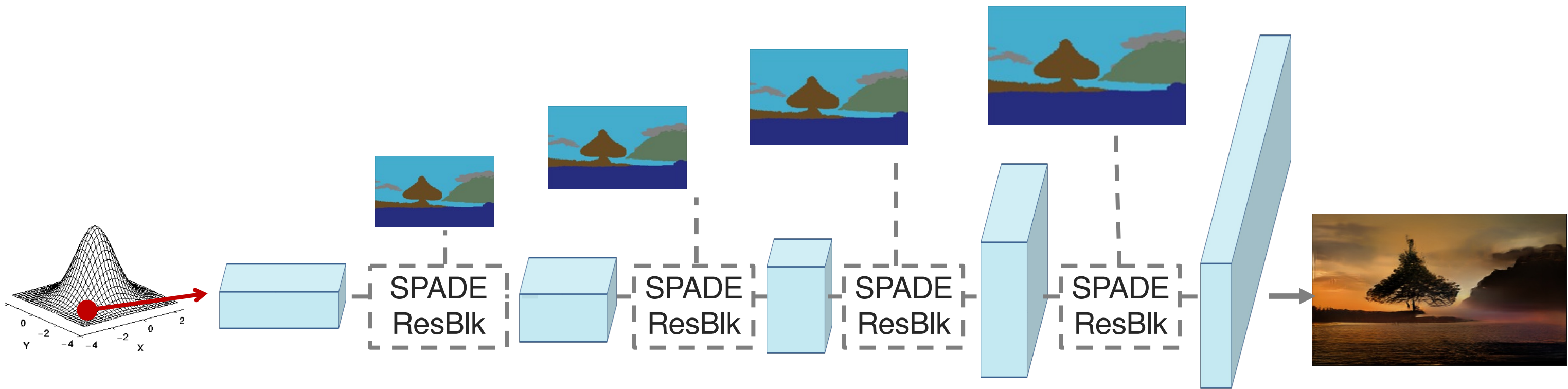
Semantic Control



Style Control

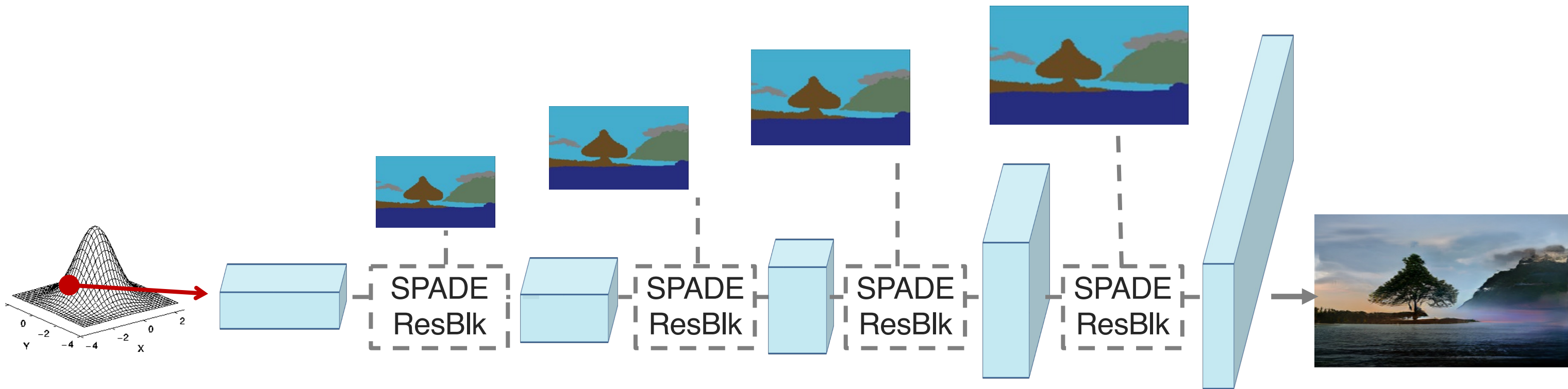


Style Control

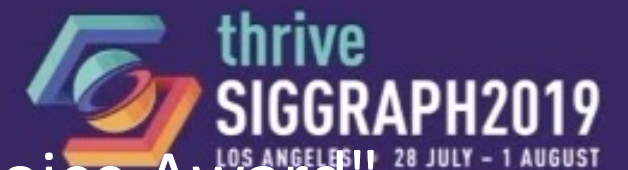
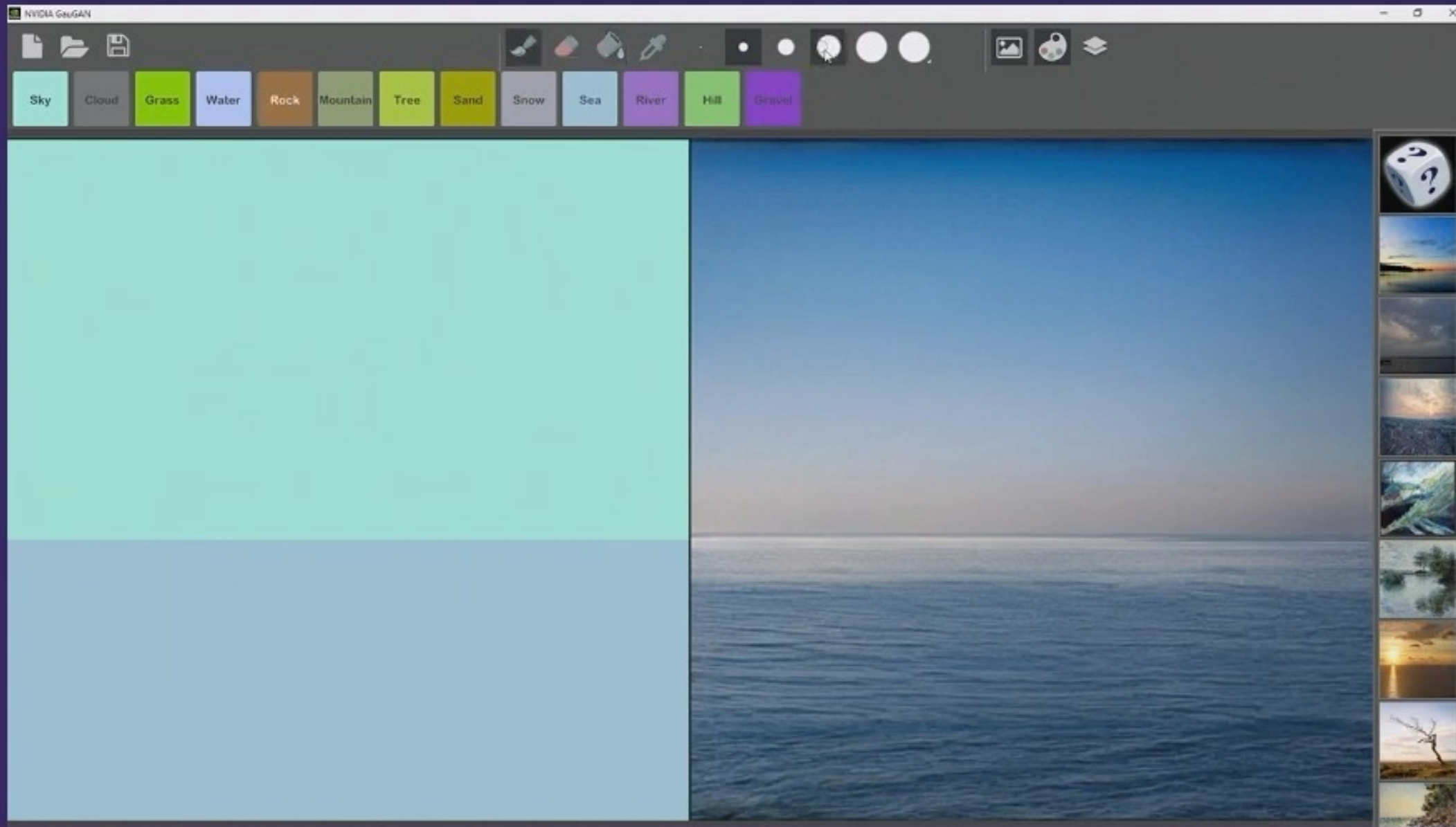


Style Manipulation

Style Control

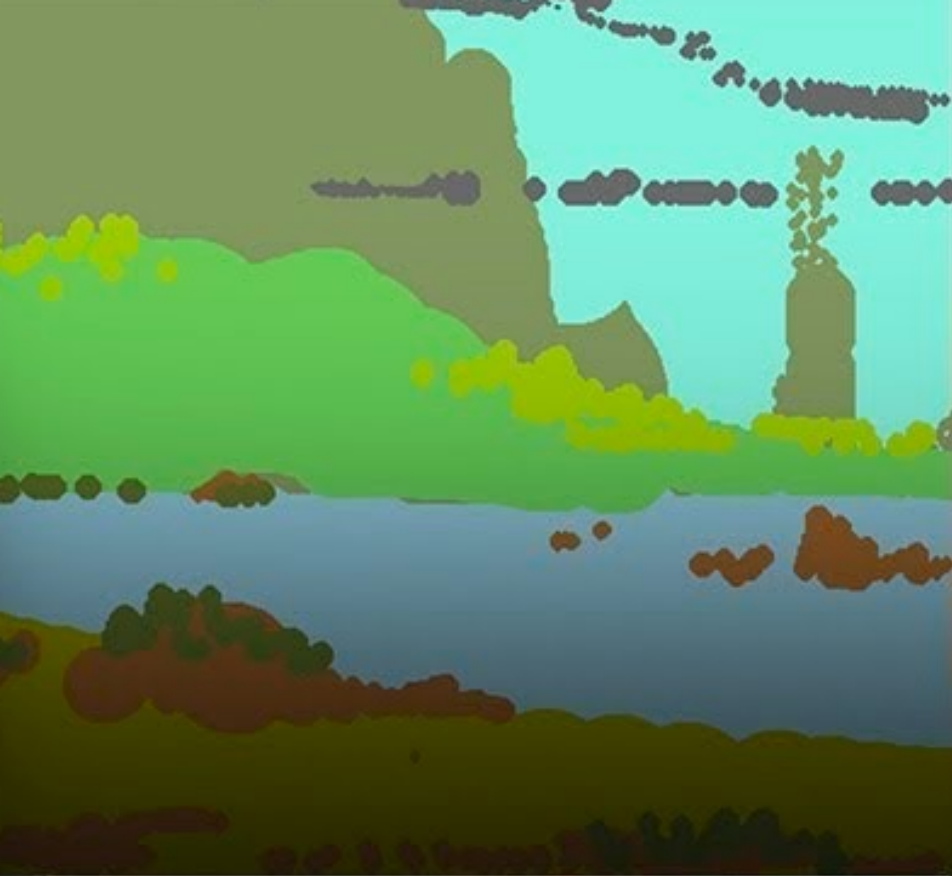


Style Manipulation



SIGGRAPH 2019 Real-time Live! "Best of Show Award" and "Audience Choice Award"

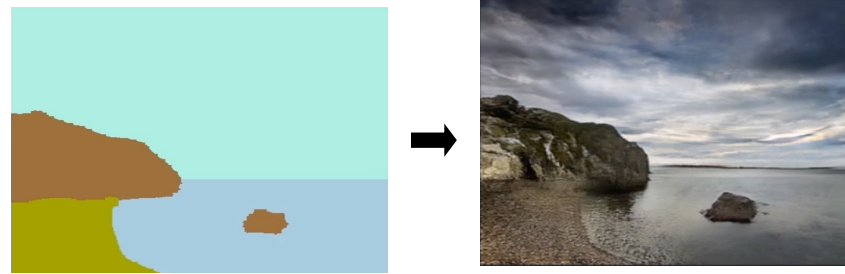
© 2019 SIGGRAPH. ALL RIGHTS RESERVED.



By Darek Zabrocki, Concept Designer and Illustrator

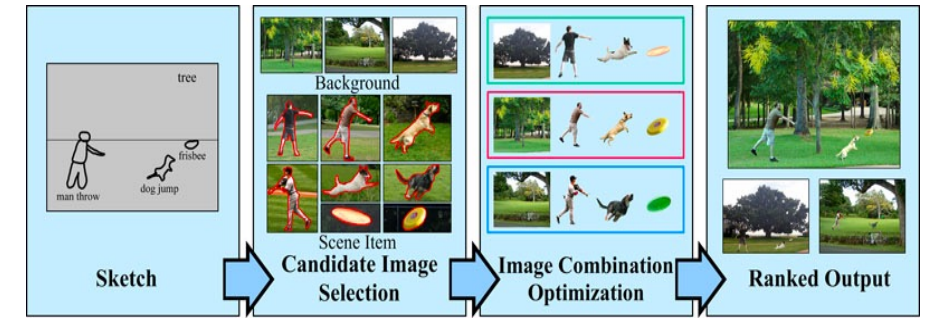
Learning vs. Exemplar-based

Learning-based



[Isola et al], [Wang et al]
[Park et al], SEAN [Zhu et al]

Exemplar-based



[Johnson et al], [Lalonde et al]
[Tao et al], [Bansal et al]

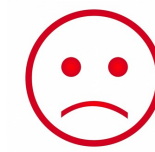
Speed



Local realism



Global realism



Match Input



Thank You!



16-726, Spring 2025