

Part 3 – Image Editing with GANs

Levent Karacan

Computer Vision Lab, Hacettepe University



HACETTEPE
UNIVERSITY
COMPUTER
VISION LAB



HACETTEPE
UNIVERSITY

Works will be presented

- Deep Convolutional Generative Adversarial Networks(DCGAN)
- Image Editing on Learned Manifold(iGAN)
- Conditional Generative Adversarial Networks(cGAN)
 - Image Generation from Text (Text2Im)
 - Stacked Generative Adversarial Networks(StackGAN)
 - Location and Description Conditioned Image Generation(GAWWN)
 - Image to Image Translation(pix2pix)
 - Image Generation from Semantic Segments and Attributes(AL-CGAN)**(Our work)**
 - Unpaired Image to Image Translation(CycleGAN)
- Neural Face Editing

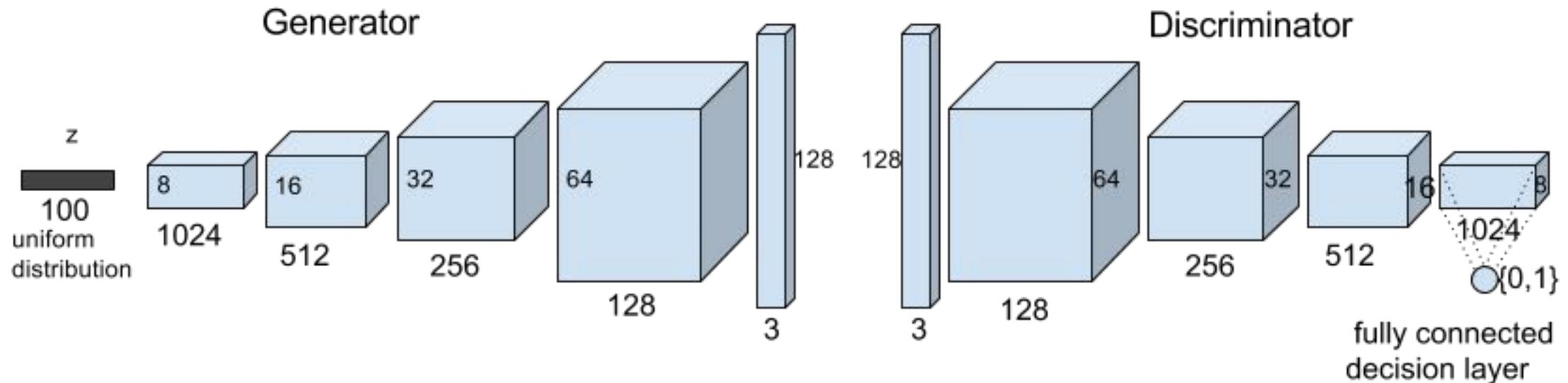
Generative Adversarial Networks(GAN)

Goodfellow vd. 2014(GAN); Radford vd. 2015(DCGAN)

- G tries to generate fake images that fool D.
- D tries to identify fake images.

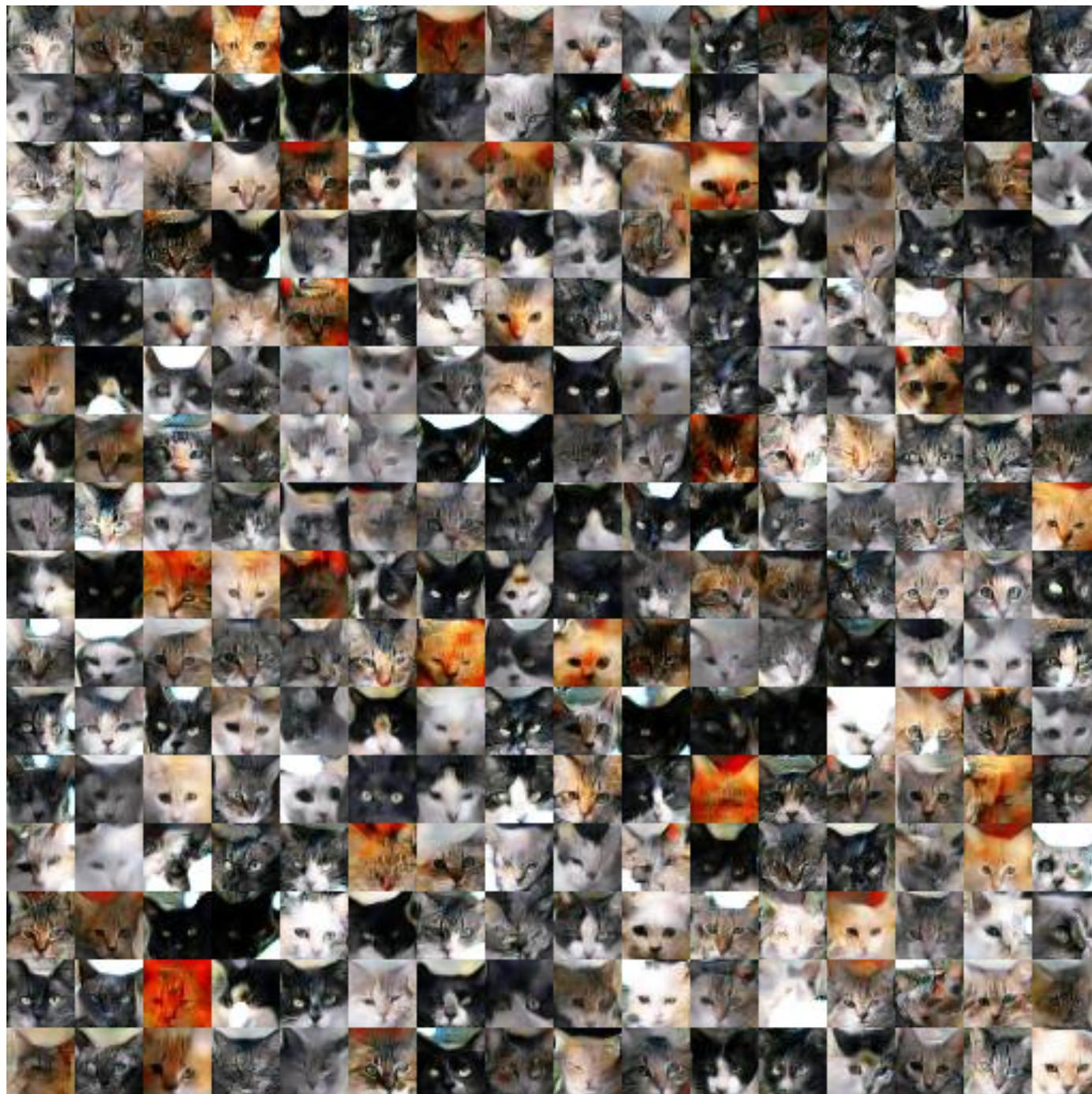
$$\mathcal{L}_{GAN}(G, D) = \mathbb{E}_{x \sim p_{data}(x)} [\log D(x)] + \mathbb{E}_{x \sim p_{data}(x), z \sim p_z(z)} [\log(1 - D(x, G(z)))]$$

$$G^* = \min_G \max_D \mathcal{L}_{GAN}(G, D)$$



DCGAN

- Cats



Source:
<https://github.com/aleju/cat-generator>

DCGAN

- Animes



Source:
<https://github.com/jaylei-cn/animeGAN>

DCGAN

- Album covers



Source:
<https://github.com/jaylei-cn/animeGAN>

DCGAN

- Flowers



DCGAN

- Faces



Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

- An image editing method that aims to find projection \mathbf{z} of input image \mathbf{x} .

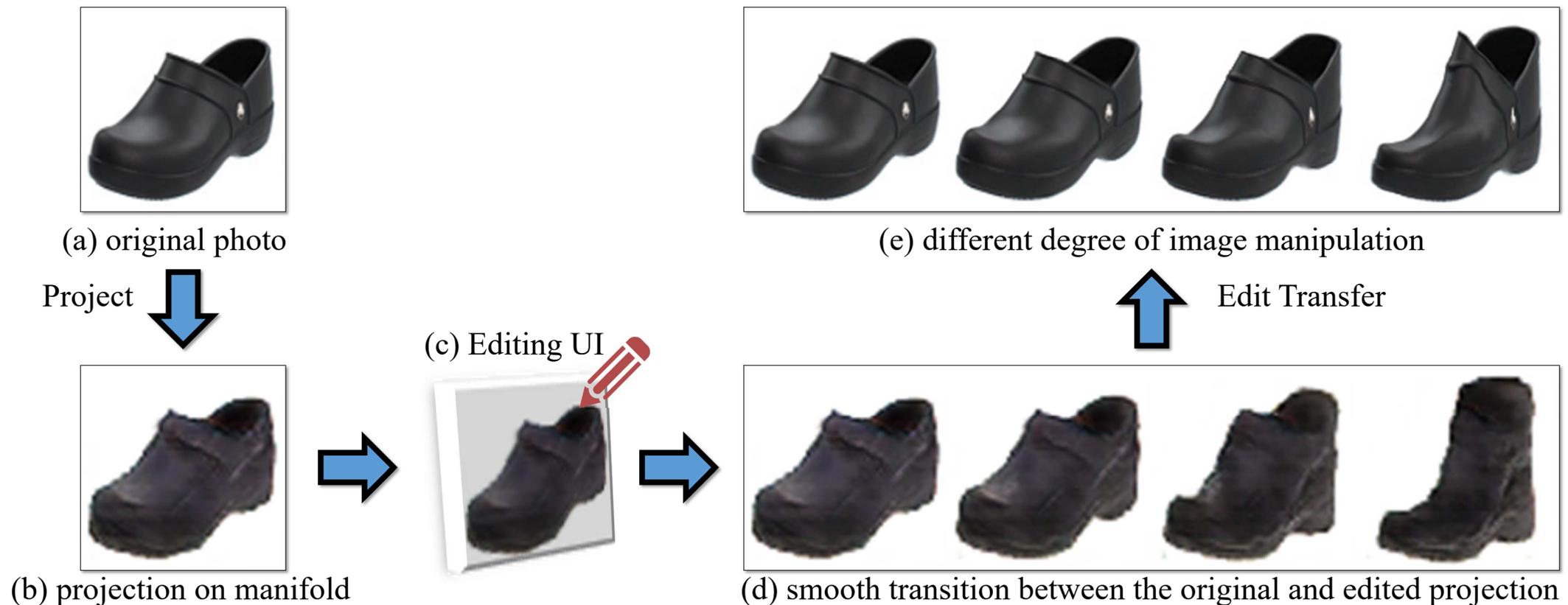


Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

- Find \mathbf{z} that generates the input image \mathbf{x} using generator network .

$$S(G(z_1), G(z_2)) \approx \|z_1 - z_2\|^2$$

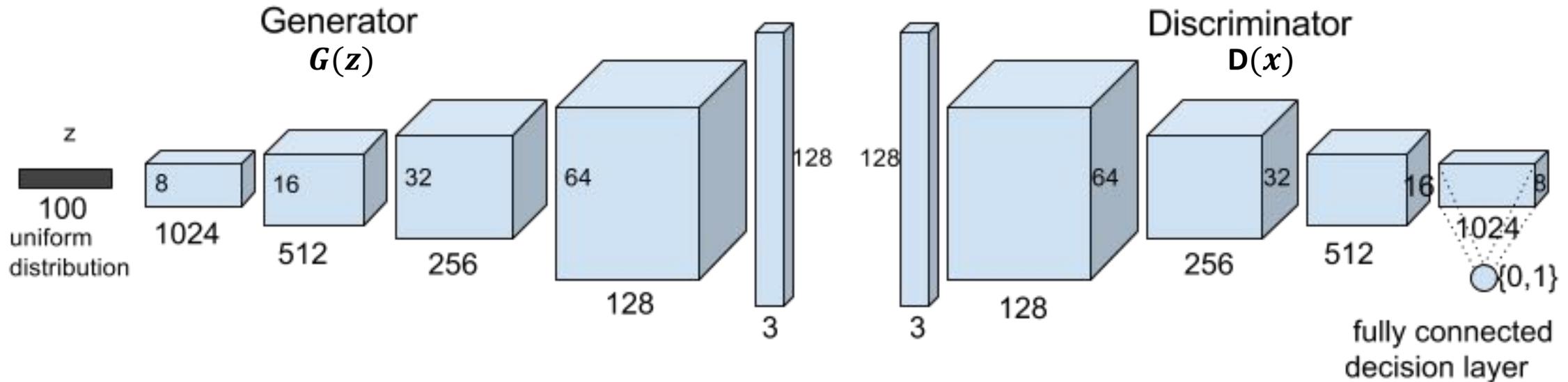


Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

- Images generated from DCGAN trained on shirt image dataset.

$$S(G(z_1), G(z_2)) \approx \|z_1 - z_2\|^2$$



random jittering

linear interpolation

Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

- Projection via optimization. L-BFGS-B method.

$$S(G(z_1), G(z_2)) \approx \|z_1 - z_2\|^2$$

$$z^* = \underset{z \in \tilde{\mathcal{Z}}}{\operatorname{argmin}} \mathcal{L}(G(z), x^R)$$

- Projection via feedforward network.

$$\mathcal{L}(x_1, x_2) = \|C(x_1) - C(x_2)\|^2$$

$$\theta_P^* = \underset{\theta_P}{\operatorname{argmin}} \sum_n \mathcal{L}(G(P(x_n^R; \theta_P), x_n^R))$$

- Hybrid method.

Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

Original photos										
Reconstruction via Optimization										
	0.165	0.164	0.370	0.279	0.350	0.249	0.437	0.255	0.178	0.227
Reconstruction via Network										
	0.198	0.190	0.382	0.302	0.251	0.339	0.482	0.270	0.248	0.263
Reconstruction via Hybrid Method										
	0.133	0.141	0.298	0.218	0.160	0.204	0.318	0.185	0.183	0.190

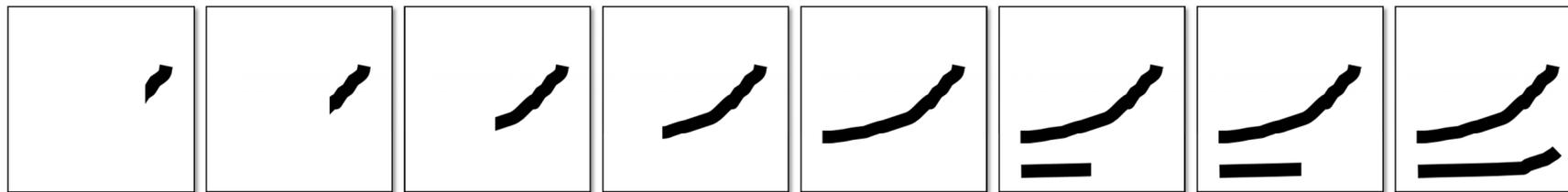
Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

$$S(G(z_1), G(z_2)) \approx \|z_1 - z_2\|^2$$

g : Color, shape and warping constraints for image editing.

$$z^* = \min_{z \in \mathbb{Z}} \left\{ \sum_g \|f_g(G(z)) - v_g\|^2 + \lambda_s \|z - z_0\|^2 \right\}$$



(a) User constraints v_g at different update steps



$G(z_0)$

(b) Updated images according to user edits

$G(z_1)$



(c) Linear interpolation between $G(z_0)$ and $G(z_1)$

Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

Edit Transfer

- A dense correspondence algorithm to estimate both the geometric and color changes induced by the editing process.



Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016

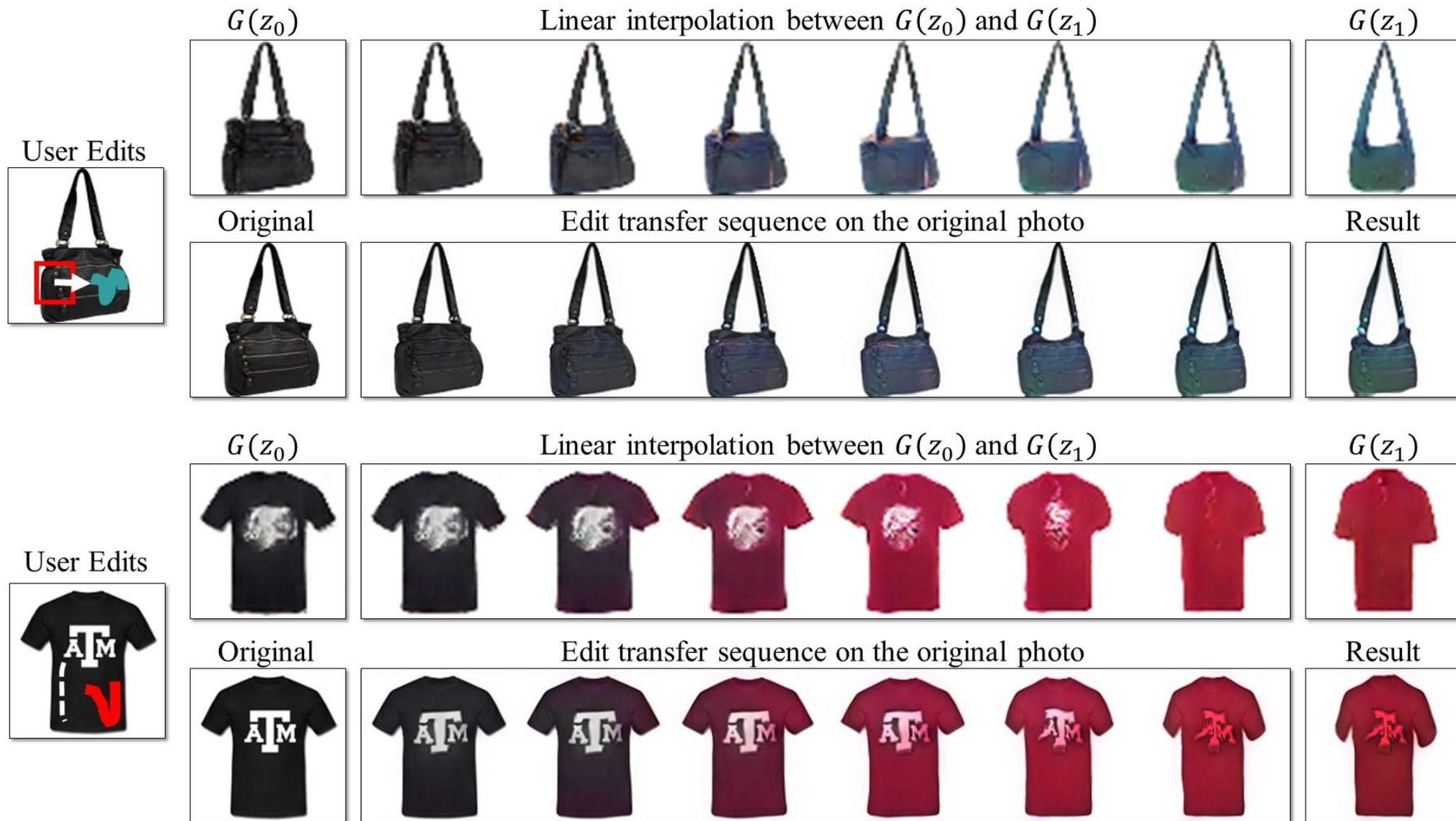
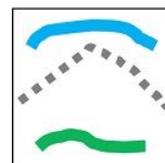
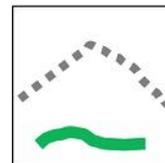
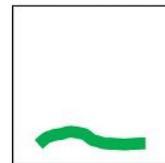


Image Editing on Learned Manifold(iGAN)

Zhu vd. 2016



User edits



Generated images



— Color

■ ■ ■ Sketch

Conditional Generative Adversarial Networks(cGAN) Mirza vd. 2014

- Concatenate condition information \mathbf{x} to noise vector \mathbf{z} and introduce to discriminator.

$$\begin{aligned}\mathcal{L}_{cGAN}(G, D) &= \mathbb{E}_{x, y \sim p_{data}(x, y)} [\log D(x, y)] + \mathbb{E}_{x \sim p_{data}(x), z \sim p_z(z)} [\log(1 - D(x, G(x, z)))] \\ G^* &= \min_G \max_D \mathcal{L}_{cGAN}(G, D)\end{aligned}$$

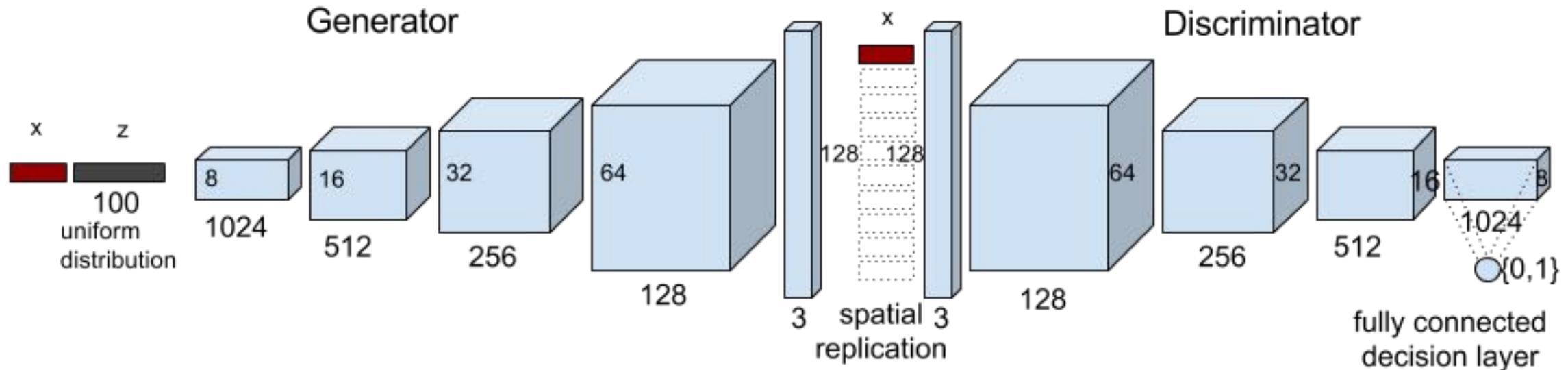


Image Generation from Text(Text2Im) Reed vd. 2016

- Discriminator network tries to classify real image and wrong text as well as real/fake image with right text.
- Condition: Text description embedding.
- CUB bird dataset(11788 images from 200 categories), Oxford-102 flower dataset(8189 images from 102 categories).

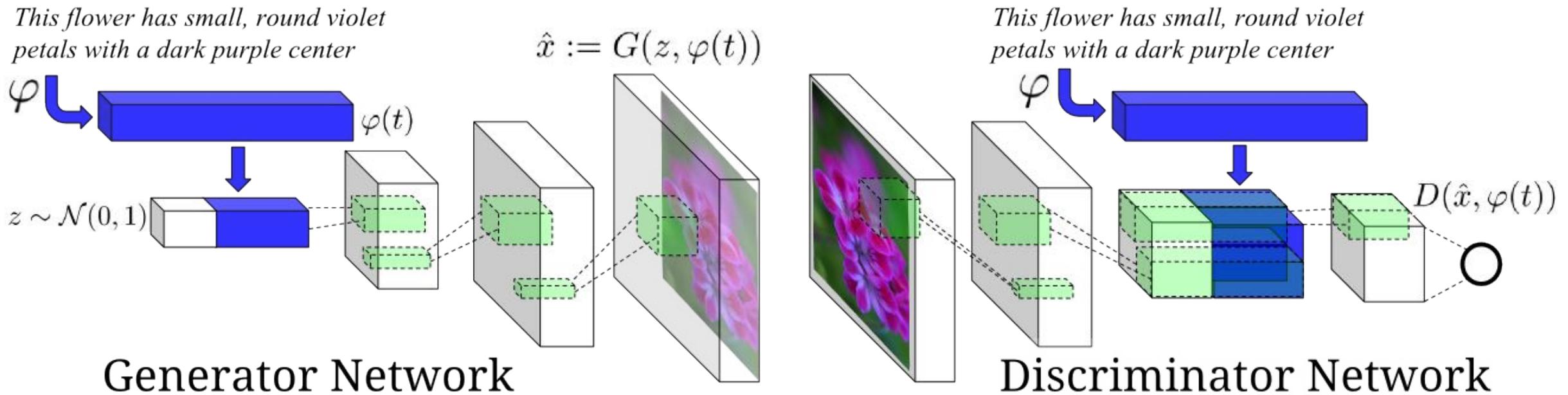


Image Generation from Text(Text2Im) Reed vd. 2016

this small bird has a pink breast and crown, and black primaries and secondaries.



the flower has petals that are bright pinkish purple with white stigma



this magnificent fellow is almost all black with a red crest, and white cheek patch.



this white and yellow flower have thin white petals and a round yellow stamen



Text descriptions (content)

Images (style)



The bird has a **yellow breast** with **grey** features and a small beak.

This is a large **white** bird with **black wings** and a **red head**.

A small bird with a **black head and wings** and features grey wings.

This bird has a **white breast**, brown and white coloring on its head and wings, and a thin pointy beak.

A small bird with **white base** and **black stripes** throughout its belly, head, and feathers.

A small sized bird that has a cream belly and a short pointed bill.

This bird is **completely red**.

This bird is **completely white**.

This is a **yellow** bird. The **wings are bright blue**.



Image Generation from Text(Text2Im) Reed vd. 2016

“Blue bird with black beak”



“This bird is completely red with black wings”



Image Generation from Text(Text2Im) Reed vd. 2016

“Small blue bird with black wings.”



“Small yellow bird with black wings”

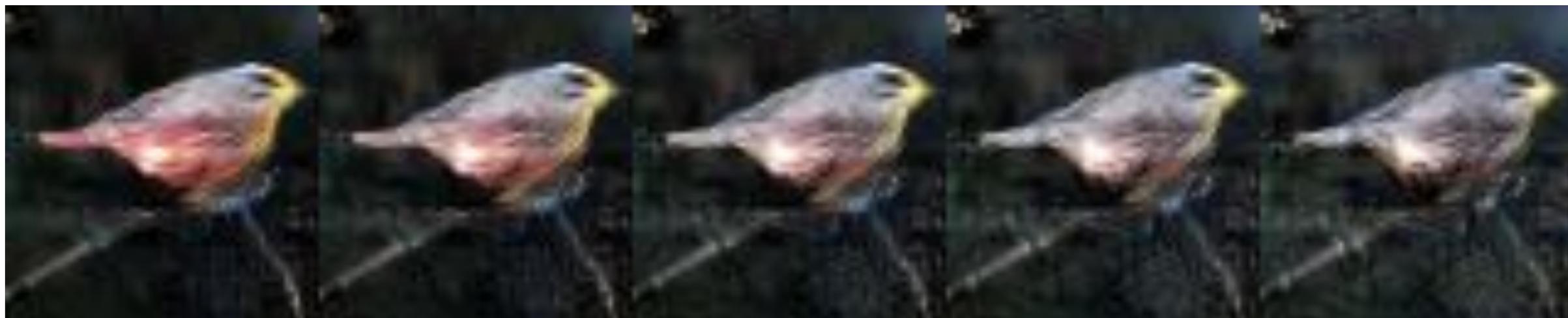


Image Generation from Text(Text2Im) Reed vd. 2016

"This bird is bright!"

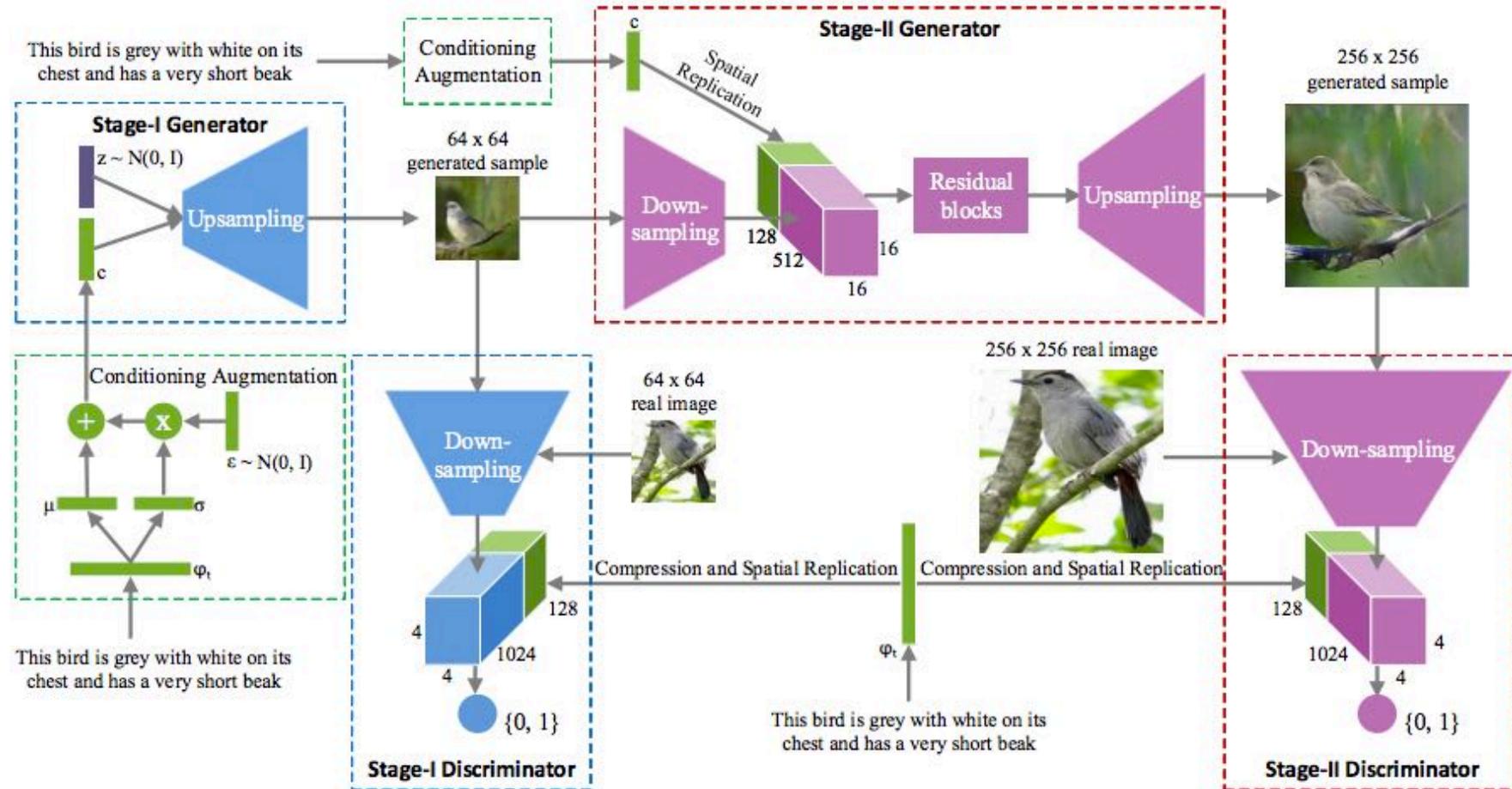


"This bird is dark"



Stacked Generative Adversarial Networks(StackGAN) Han vd. 2016

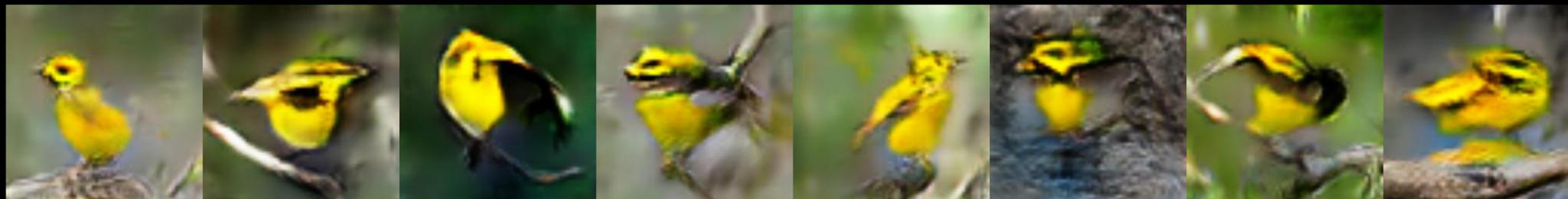
- There are 2 stages.
- Stage-I GAN: Generates low resolution images.
 - Conditioning Augmentation
 - Regularization term is added to generator.
$$D_{KL}(\mathcal{N}(\mu(\varphi_t) \parallel \mathcal{N}(0, I))$$
- Stage-II GAN: Generates high resolution detailed images.
 - Noise vector is not used.



Stacked Generative Adversarial Networks(StackGAN) Han vd. 2016

A small yellow bird with a black crown and a short black pointed beak

Stage-I



Stage-II



Stacked Generative Adversarial Networks(StackGAN) Han vd. 2016

This small bird has a white breast, light grey head, and black wings and tail

Stage-I



Stage-II



Stacked Generative Adversarial Networks(StackGAN) Han vd. 2016

This flower has long thin yellow petals and a lot of yellow anthers in the center

Stage-I



Stage-II



Stacked Generative Adversarial Networks(StackGAN) Han vd. 2016

This flower is white, pink, and yellow in color, and has petals that are multi colored

Stage-I



Stage-II



Location and Description Conditioned Image Generation(GAWWN)

Reed vd. 2016



This bird is completely black.



This bird is bright blue.

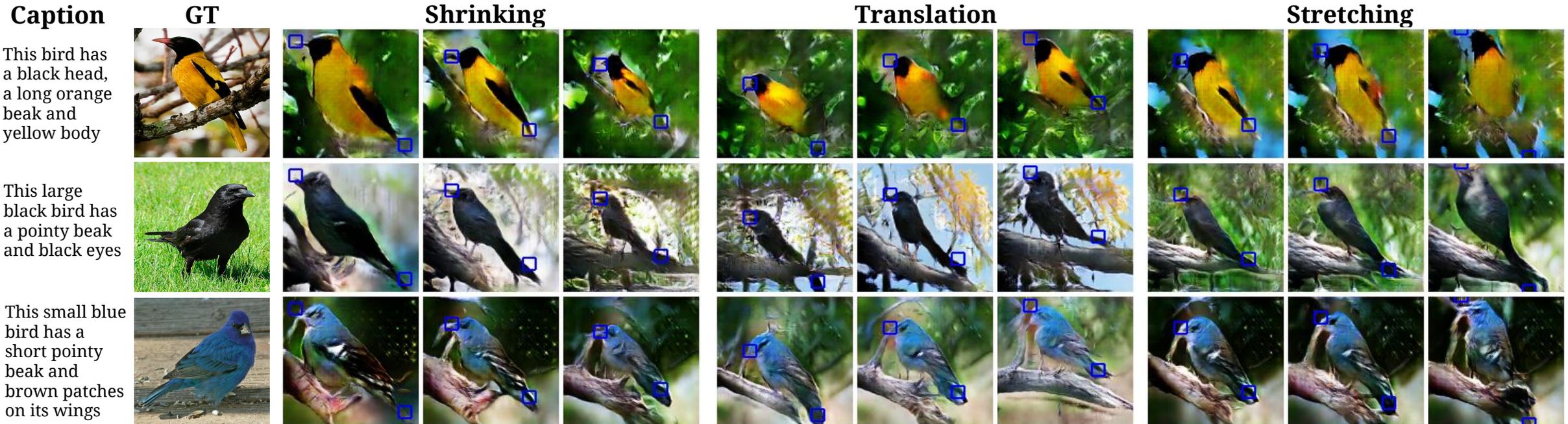


a man in an orange jacket, black pants and a black cap wearing sunglasses skiing

Location and Description Conditioned Image Generation(GAWWN)

Reed vd. 2016

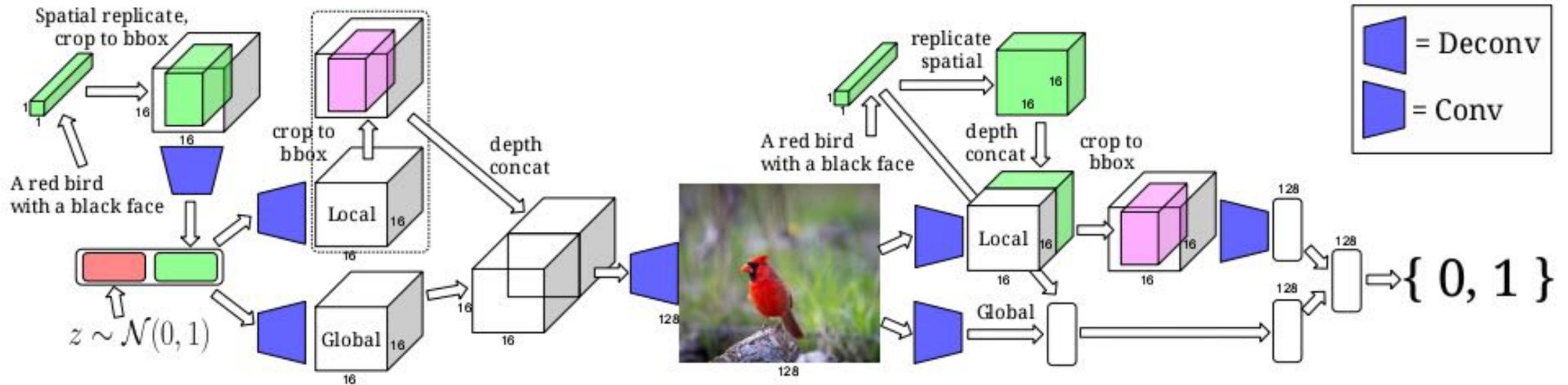
- Keypoint conditioned architecture.



Location and Description Conditioned Image Generation(GAWWN)

Reed vd. 2016

- Bounding box conditioned architecture.



Location and Description Conditioned Image Generation(GAWWN)

Reed vd. 2016

- Bounding box conditioned architecture.

Caption

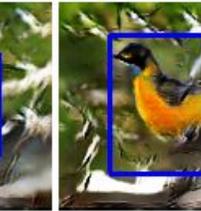
GT

Shrinking

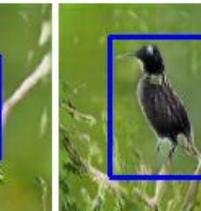
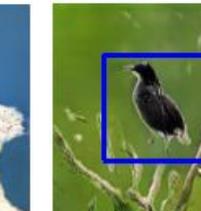
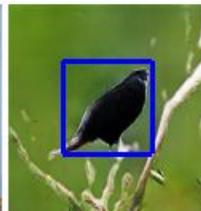
Translation

Stretching

This bird has a black head, a long orange beak and yellow body



This large black bird has a pointy beak and black eyes



This small blue bird has a short pointy beak and brown patches on its wings

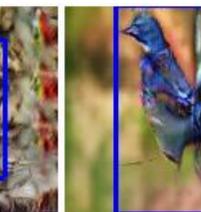
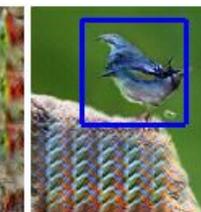
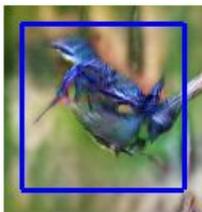
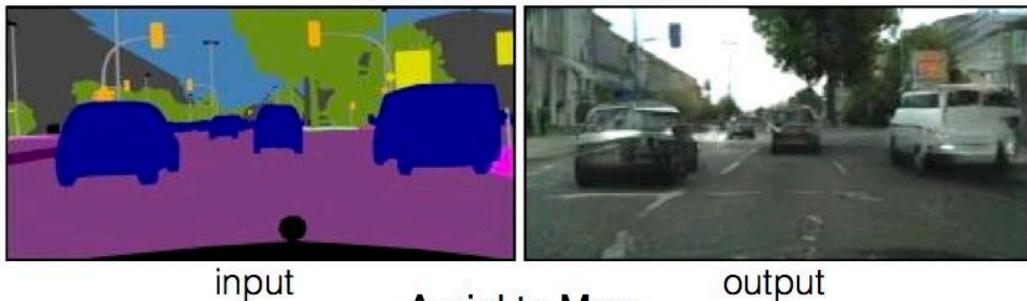
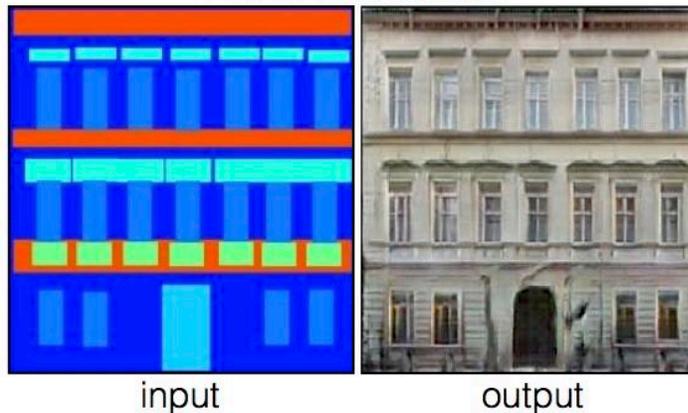


Image to Image Translation (pix2pix) Isola et al. 2017

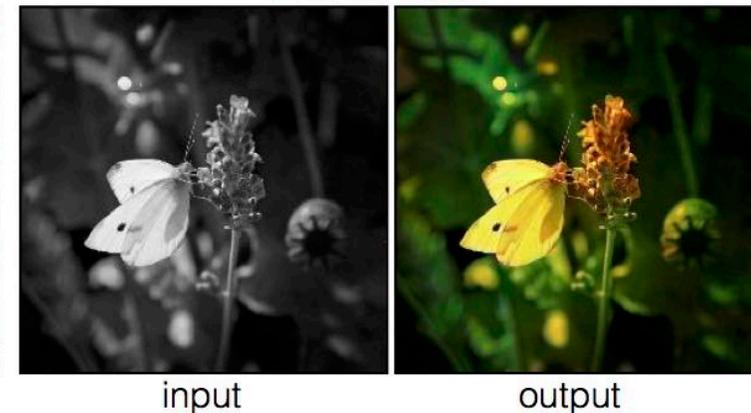
Labels to Street Scene



Labels to Facade



BW to Color



Aerial to Map



Day to Night



Edges to Photo



Image to Image Translation(pix2pix) Isola vd. 2017

$$G^* = \underset{G}{\operatorname{argmin}} \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G)$$

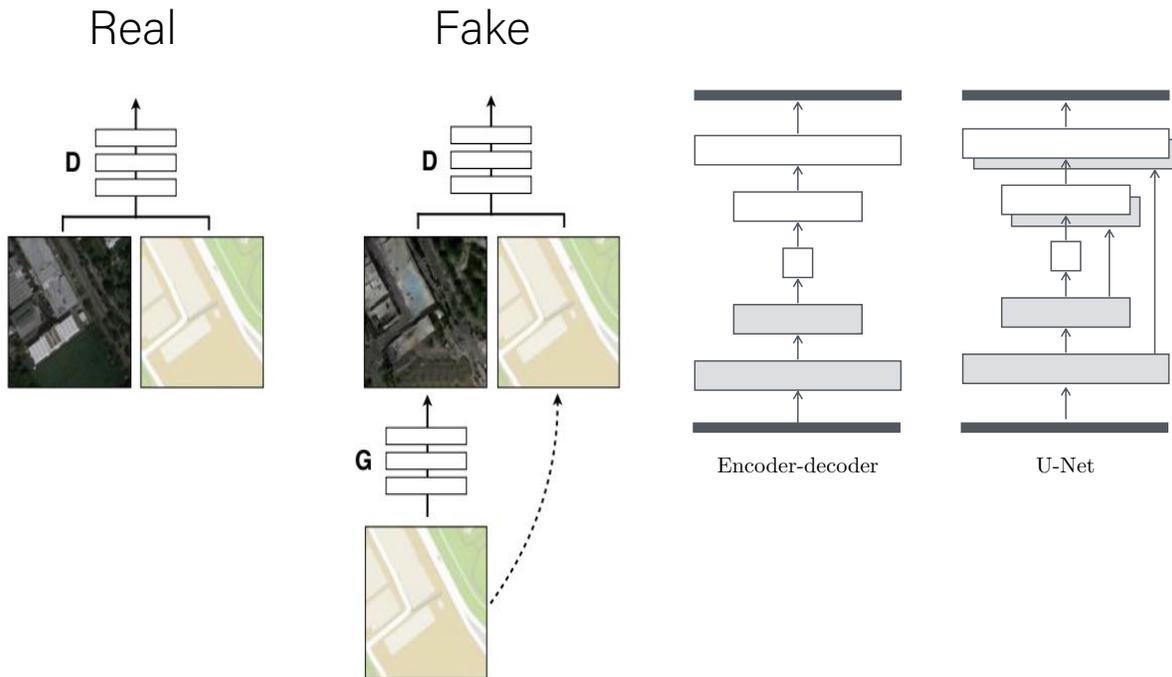
Adversarial Loss

$$\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{x,y \sim p_{data}(x,y)} [\log D(x, y)] +$$

$$\mathbb{E}_{x \sim p_{data}(x)} [\log(1 - D(x, G(x)))]$$

L1 Loss

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y \sim p_{data}(x,y), z \sim p_z(z)} [\|y - G(x, z)\|_1]$$



- G tries to generate fake images that fool D.
- D tries to identify fake images.

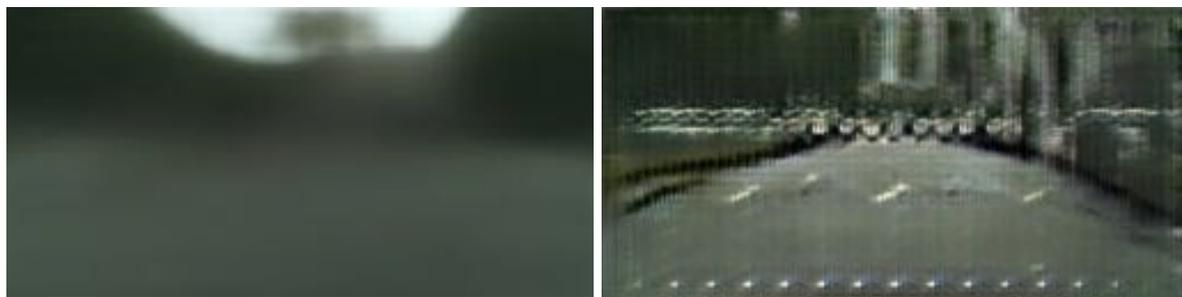
- Noise vector is removed, Instead dropout is used to provide stochasticity.
- Skip connections on Generative model
- PatchGAN is proposed for discriminator instead of pixel GAN.

Image to Image Translation (pix2pix) Isola vd. 2017

L1

L1+cGAN

Encoder-decoder



U-Net



- U-Net provides to include low-level features to be used to generate more realistic images.
- PatchGAN provides to generate sharper images.

Input

Ground truth

L1

cGAN

L1 + cGAN



Image to Image Translation(pix2pix) Isola vd. 2017

Input

Real

L1

cGAN

cGAN + L1



Isola, P., Zhu, J.Y., Zhou, T. and Efros, A.A. "Image-to-image translation with conditional adversarial networks.". In *CVPR 2017*.

Image to Image Translation(pix2pix) Isola vd. 2017

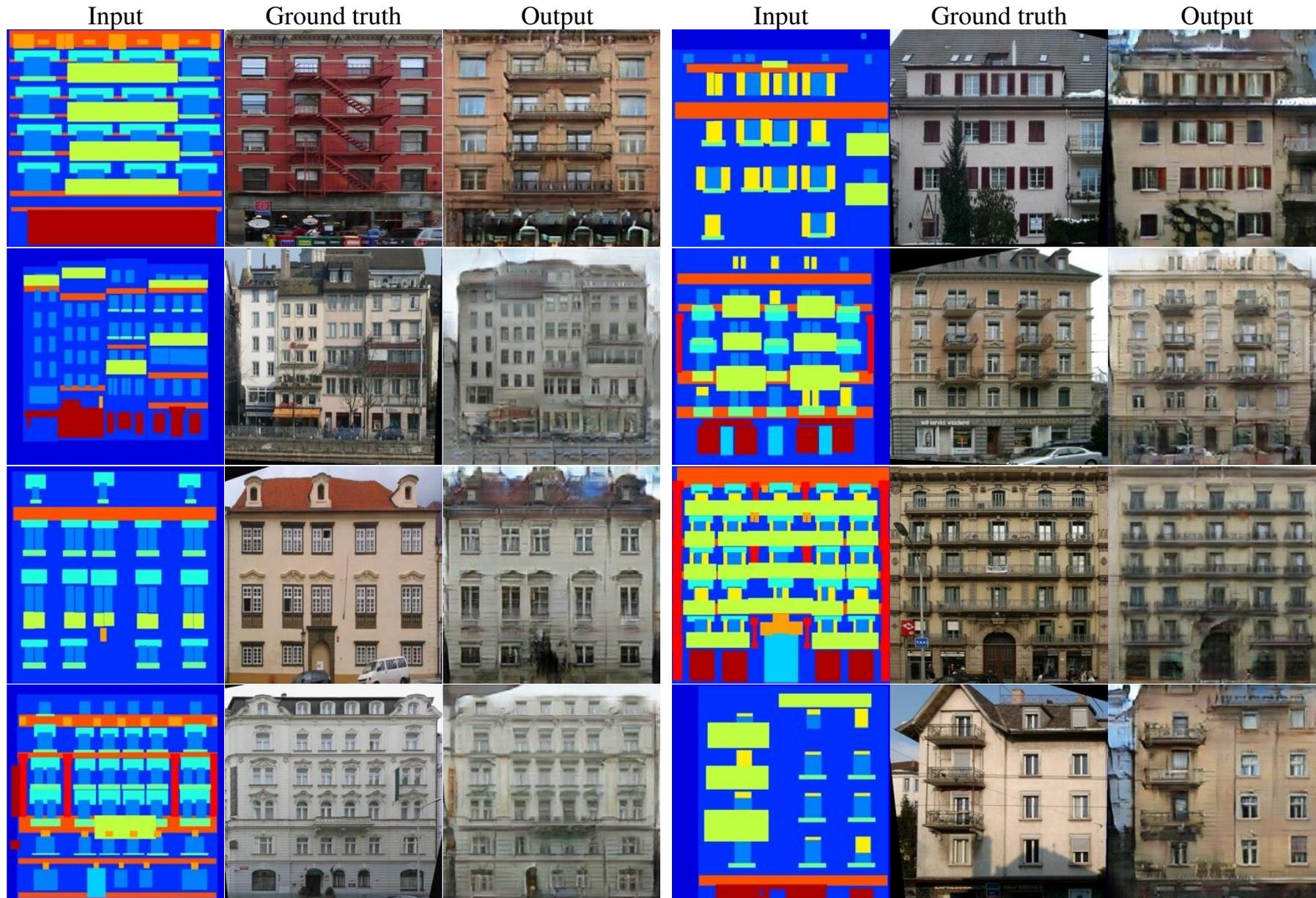


Image to Image Translation(pix2pix) Isola vd. 2017

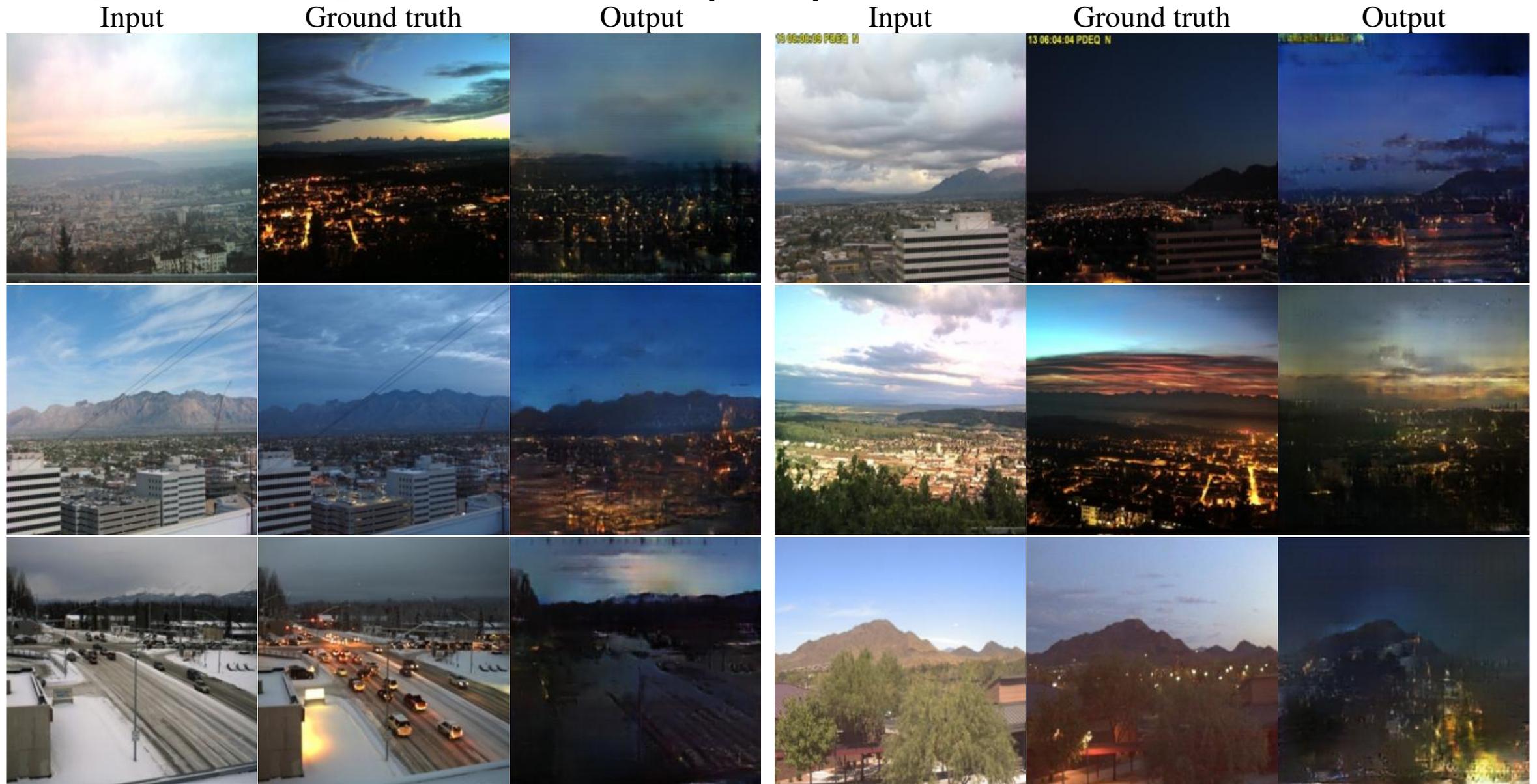


Image to Image Translation(pix2pix) Isola vd. 2017



Image to Image Translation(pix2pix) Isola vd. 2017



Image to Image Translation(pix2pix) Isola vd. 2017



Attribute and Layout Conditioned Image Generation(AL-CGAN) Our work

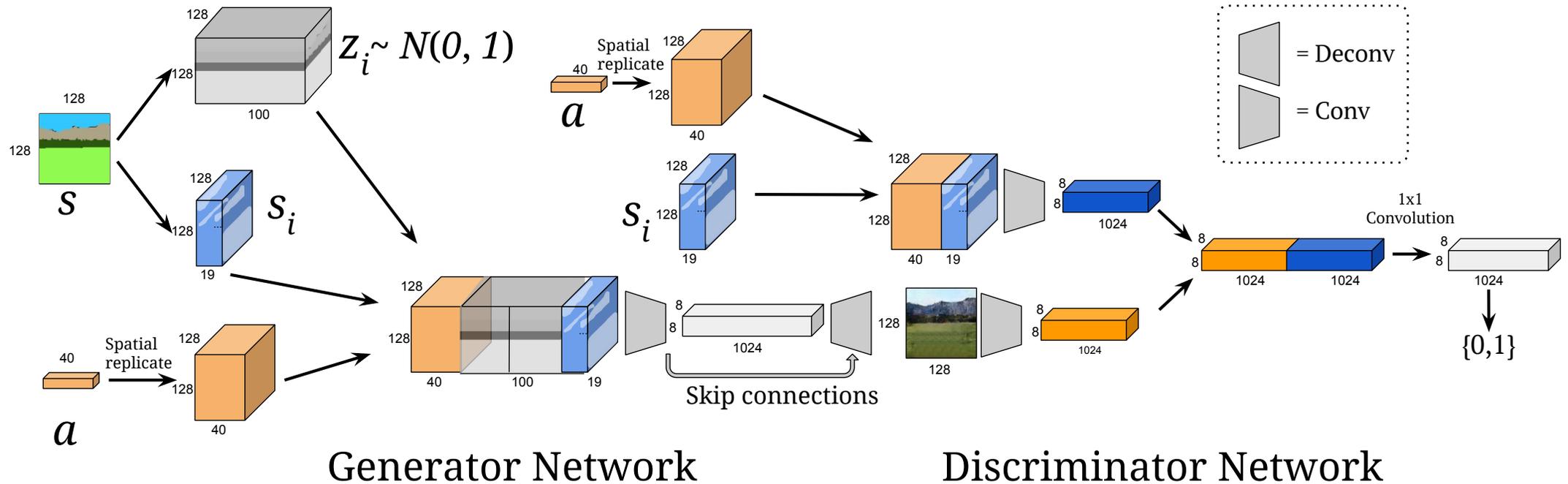


Attribute and Layout Conditioned Image Generation (ALCGAN) Our work

$$\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{x,s,a \sim p_{data}(x,s,a)} [\log D(x, s, a)] + \mathbb{E}_{s,a \sim p_{data}(s,a), z \sim p_z(z)} [\log(1 - D(x, G(z, s, a)))]$$

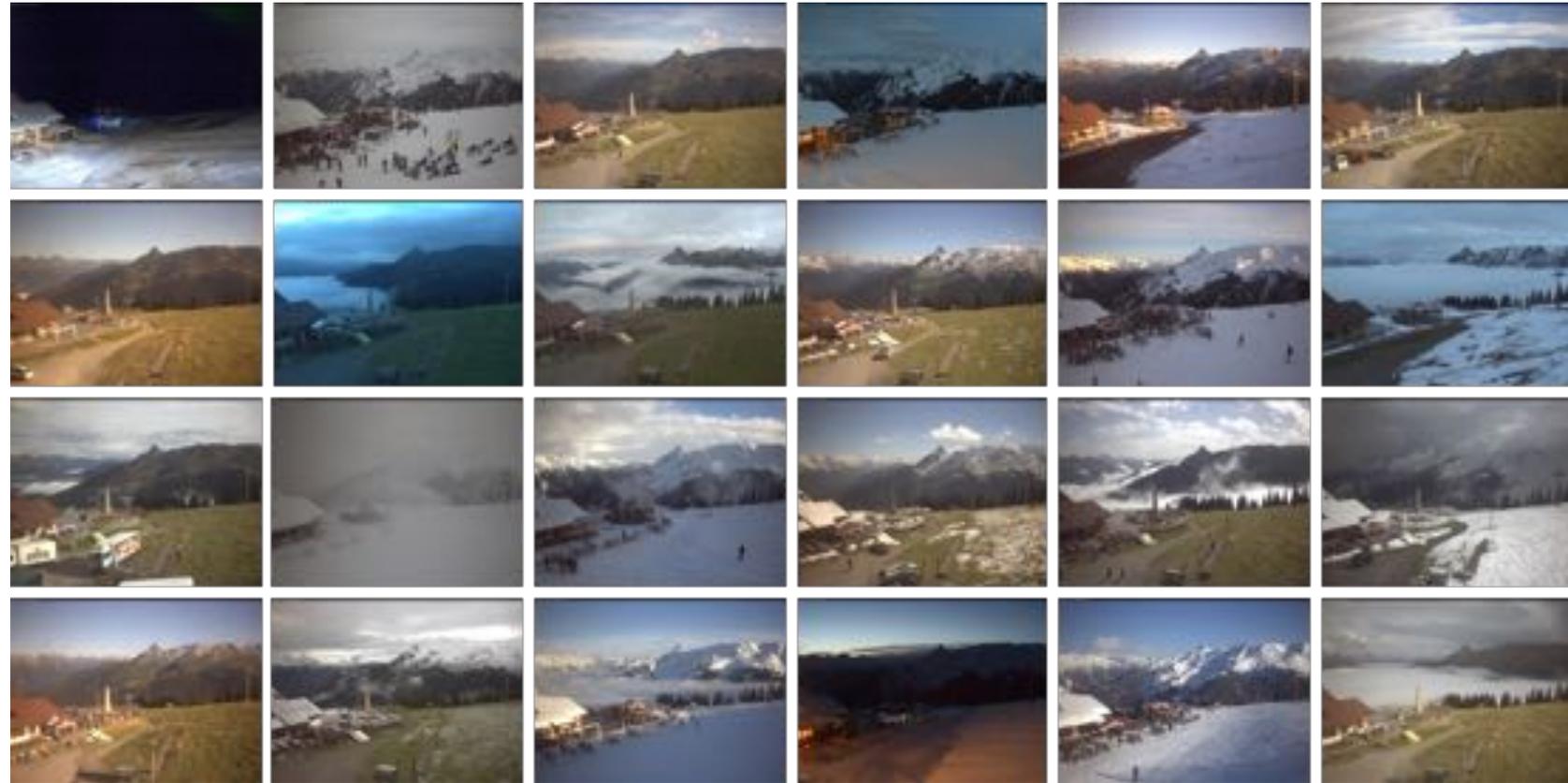
$$\min_G \max_D \mathcal{L}_{cGAN}(G, D)$$

- The noise vectors z are specific to the semantic layout.
- This provides the diversity in generated samples.



Dataset Laffont vd. 2014

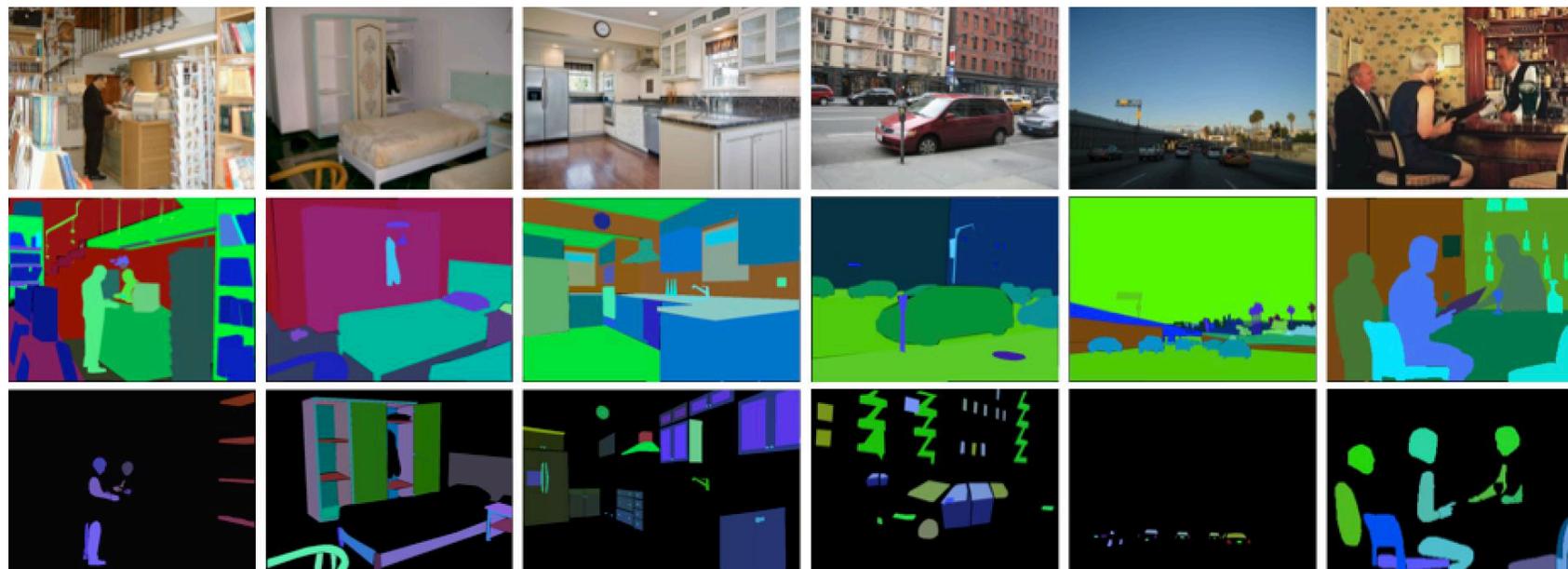
- Transient Attribute Dataset
 - 8571 outdoor images from 101 web cams located in different places.
 - 40 dimensional transient attributes for each image.
 - We annotate semantic layouts of 101 scenes with predefined 18 categories e.g. sky, tree, building, mountain, etc.



P.-Y. Laffont, Z. Ren, X. Tao, C. Qian, and J. Hays, “Transient attributes for high-level understanding and editing of outdoor scenes,” *ACM Transactions on Graphics*, vol. 33, no. 4, 2014 .

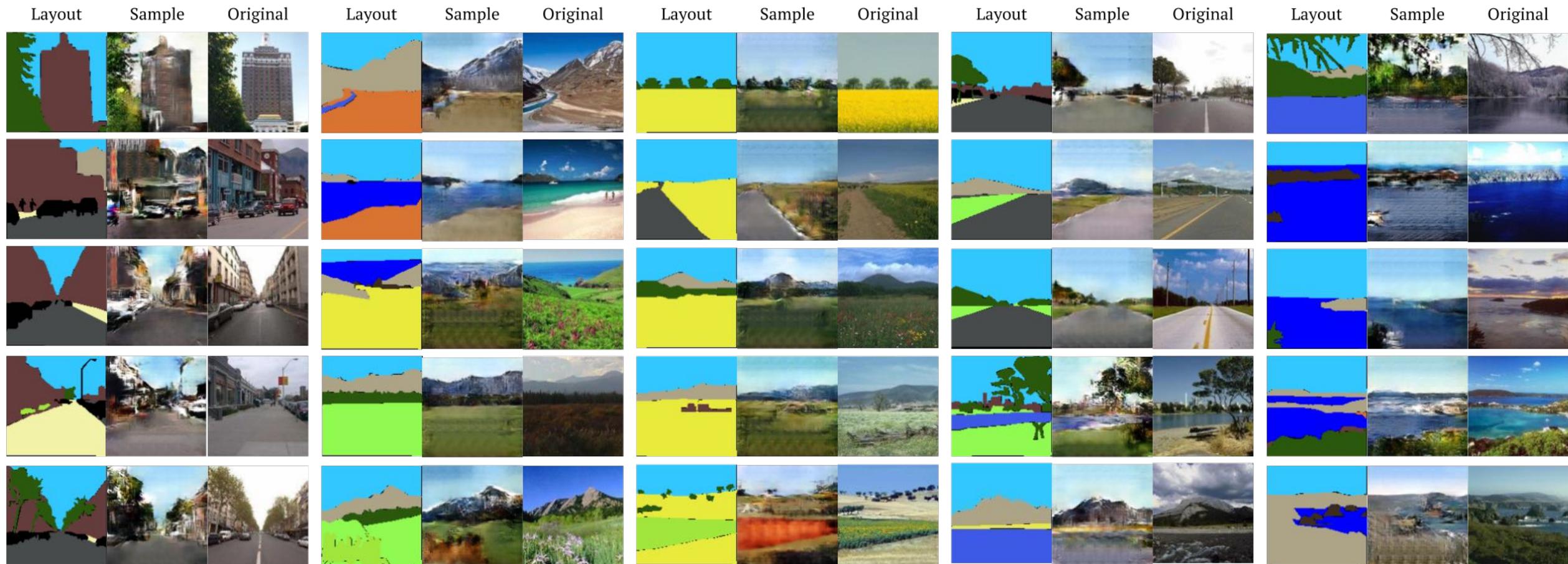
Dataset Zhou vd. 2017

- ADE20K
 - 22210 indoor and outdoor scenes with semantically labeled layouts.
 - We selected 9201 outdoor scenes according to predefined 18 categories.
 - We predicted transient attributes for each image using a deep transient model.



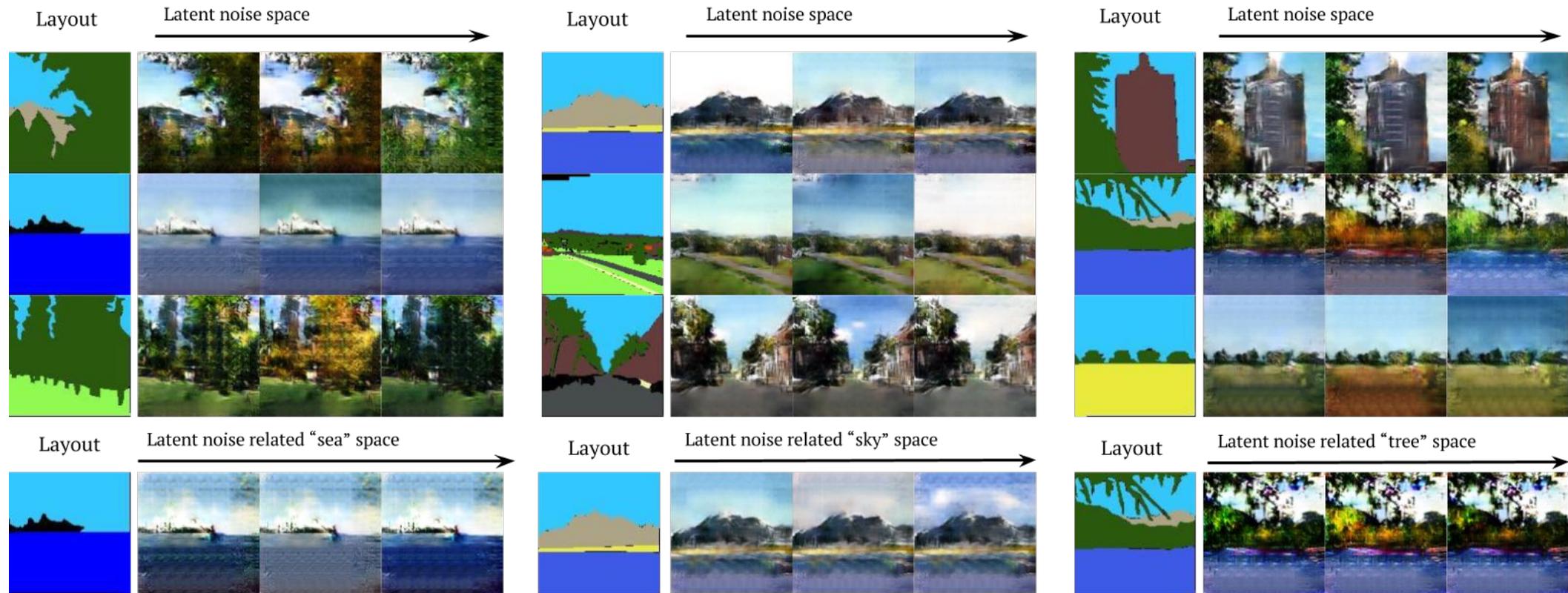
R. Baltenberger, M. Zhai, C. Greenwell, S. Workman, and N. Jacobs. A Fast Method for Estimating Transient Scene Attributes. In *WACV 2016*.

Attribute and Layout Conditioned Image Generation(AL-CGAN) Our work



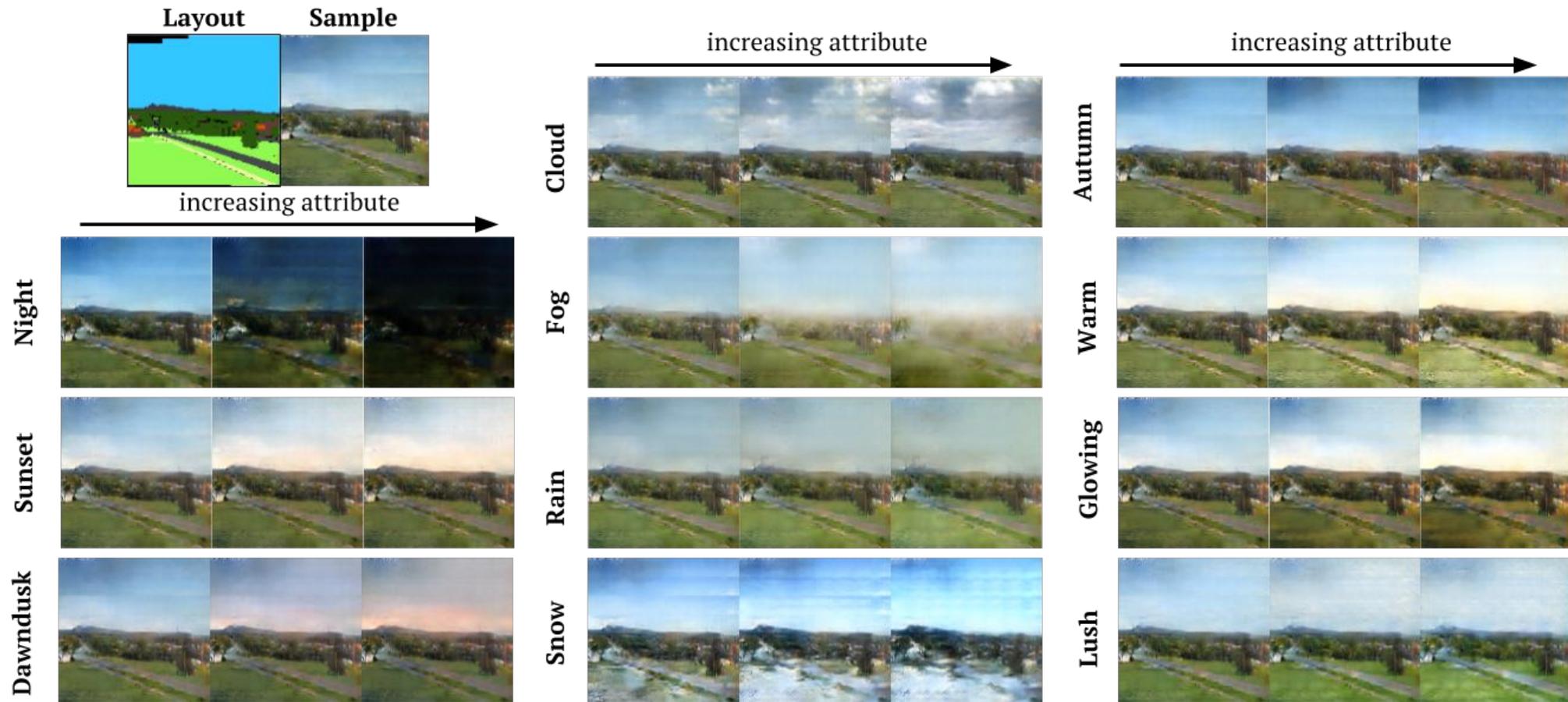
Attribute and Layout Conditioned Image Generation(AL-CGAN) Our work

- Diversity

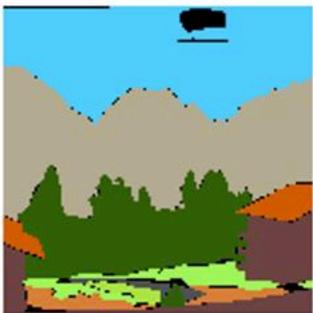
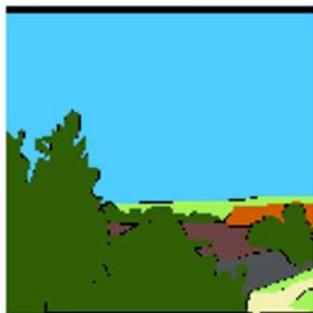


Attribute and Layout Conditioned Image Generation(AL-CGAN) Our work

- Diversity by transient attributes.

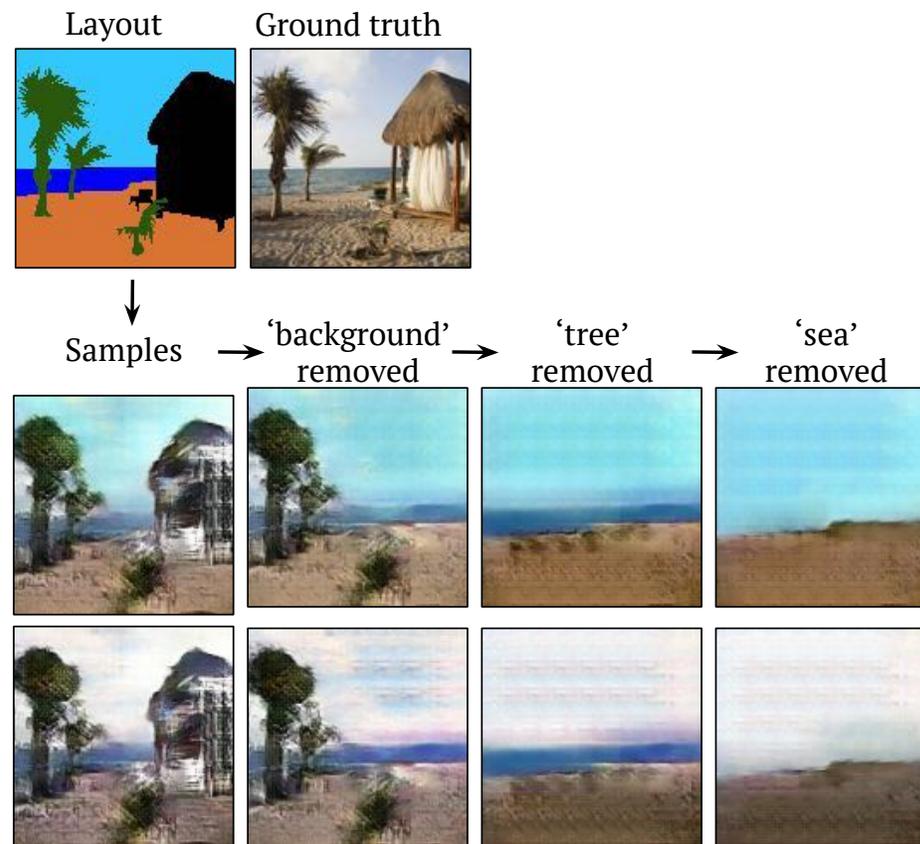
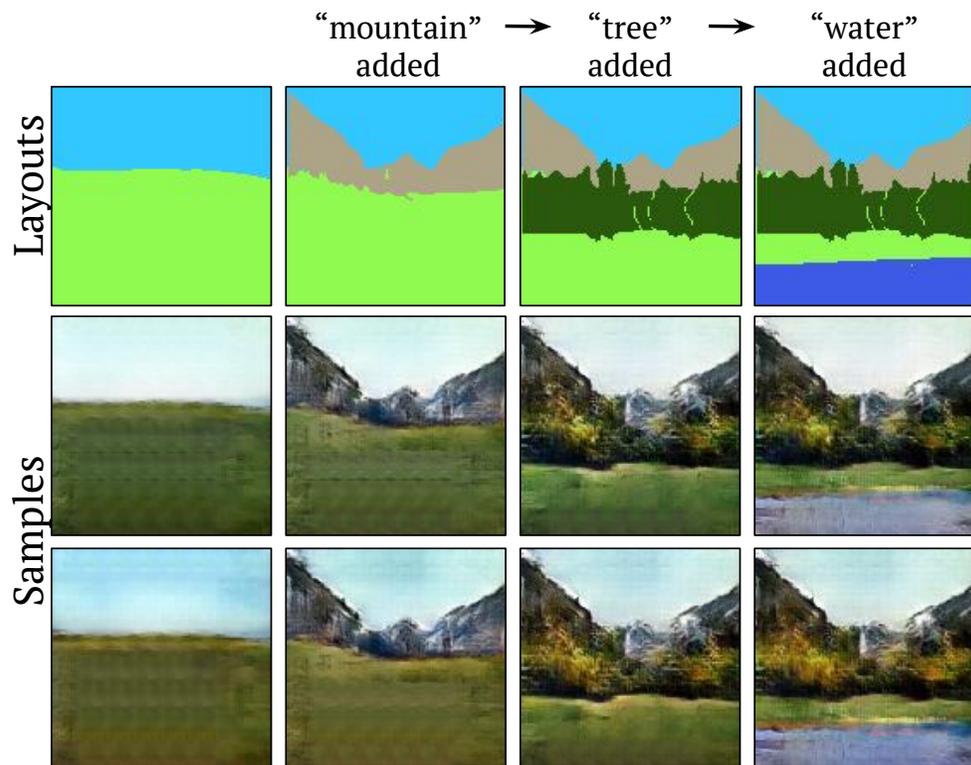


Layouts



Attribute and Layout Conditioned Image Generation(AL-CGAN) Our work

- Object adding / subtracting.

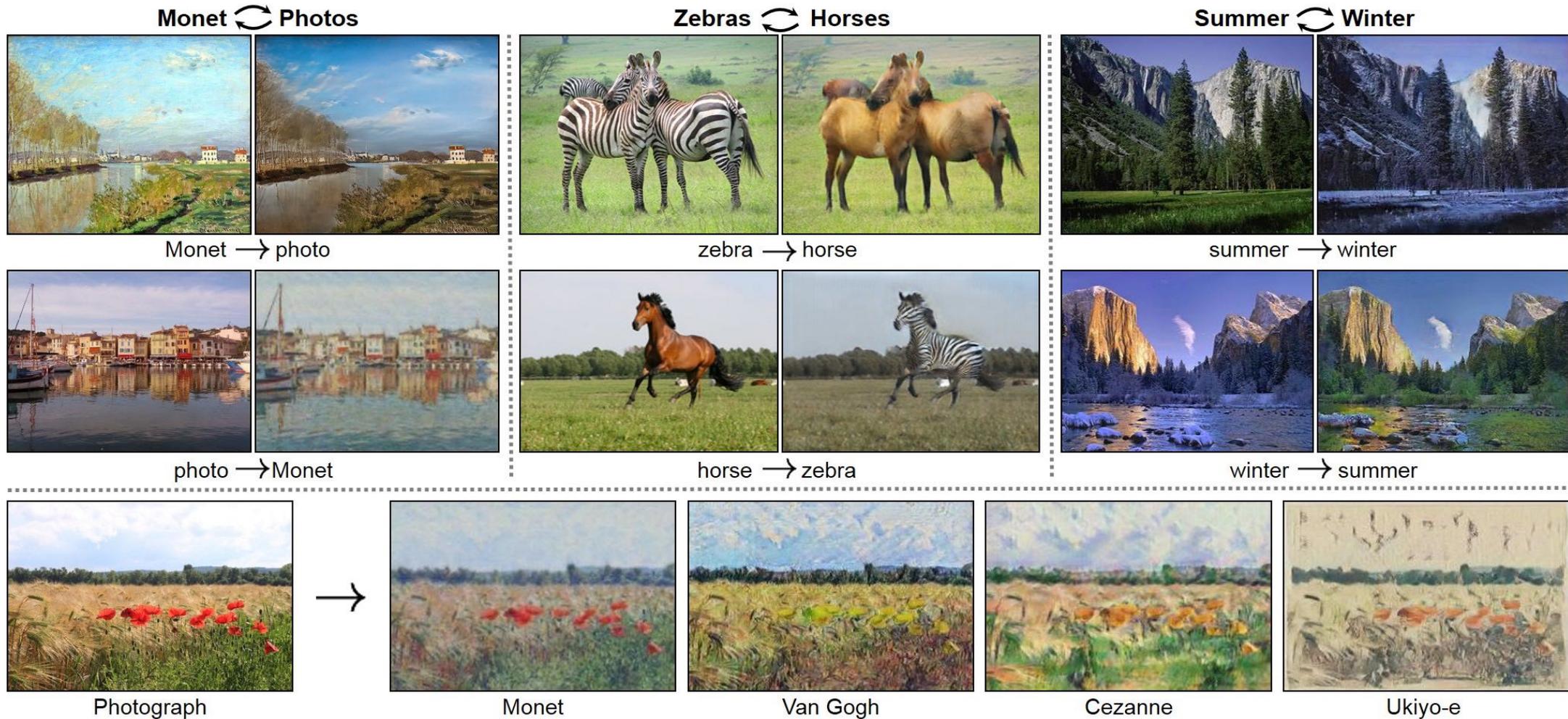


Object Adding

AL-CGAN vs pix2pix

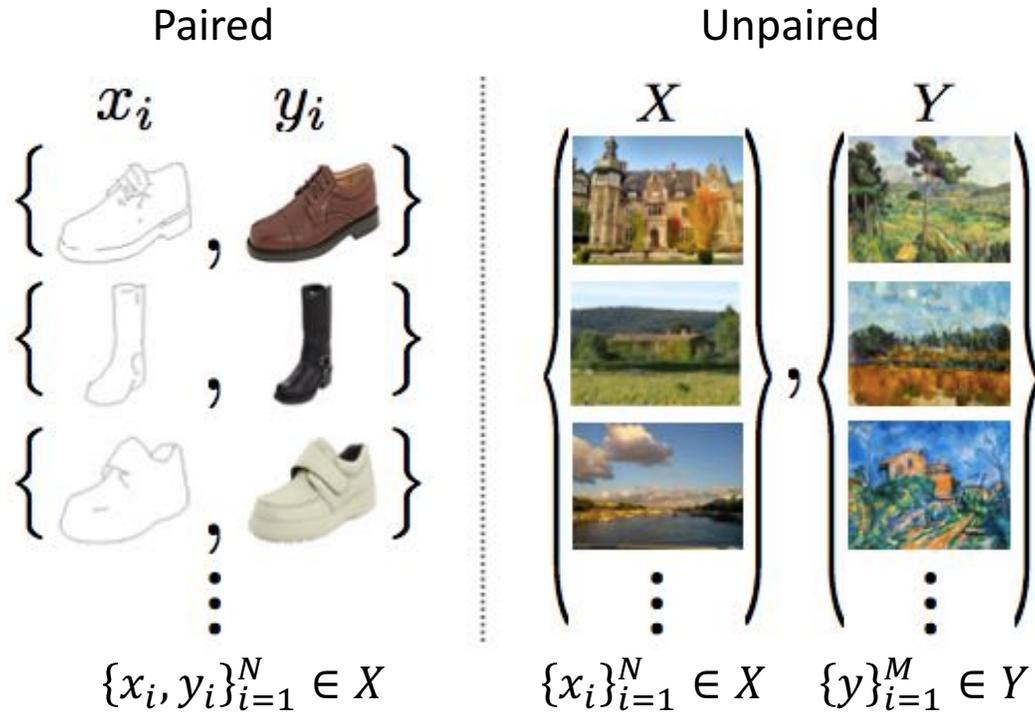


Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017



Unpaired Image to Image Translation(CycleGAN)

Zhu vd. 2017



Cycle Consistency

"if we translate, e.g., a sentence from English to French, and then translate it back from French to English, we should arrive back at the original sentence."

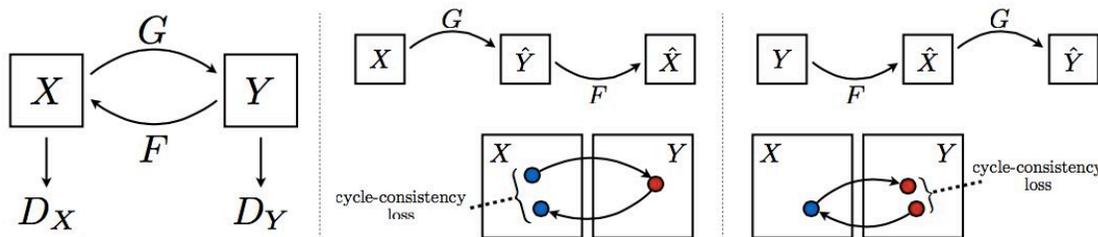
$$G: X \rightarrow Y$$

$$F: Y \rightarrow X$$

$$F(G(x)) \approx x$$

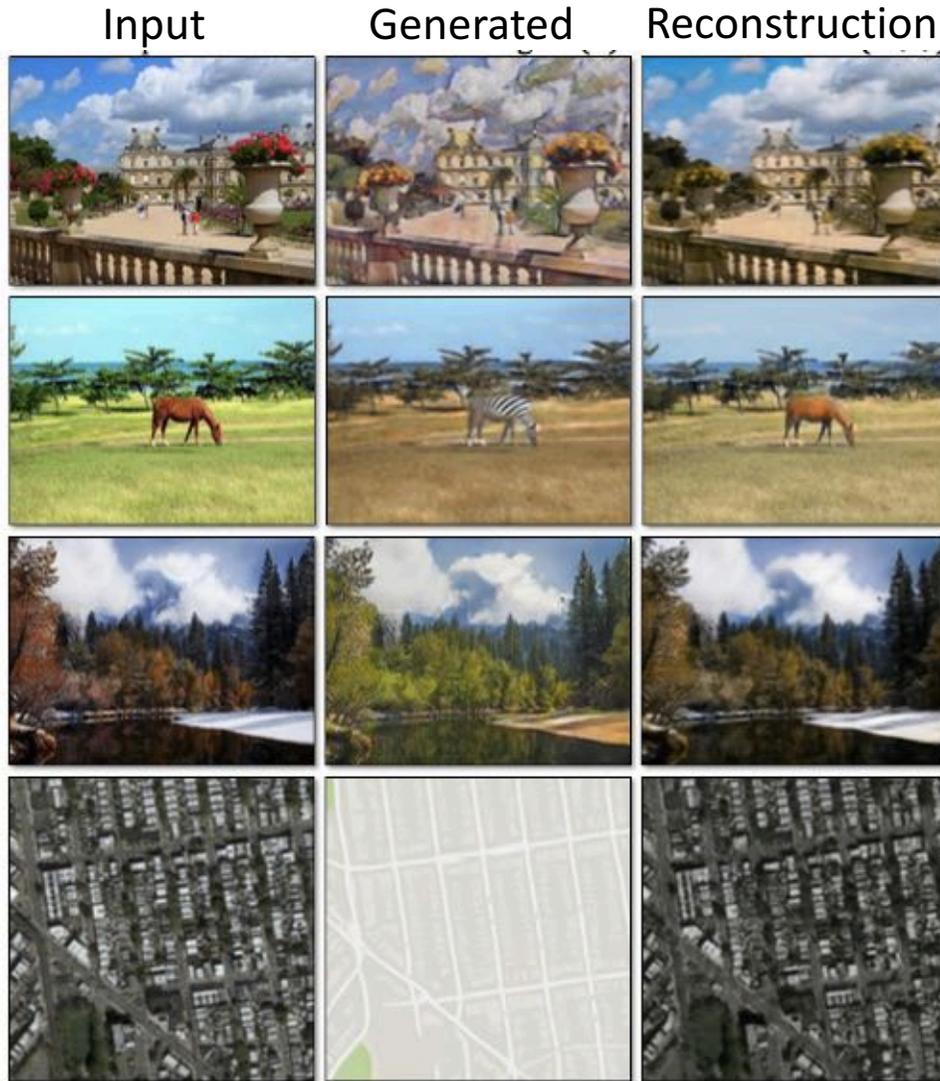
$$G(G(y)) \approx y$$

$$\mathcal{L}_{cyc}(G, F) = \mathbb{E}_{x \sim p_{data}(x)} [\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_{data}(y)} [\|G(F(y)) - y\|_1]$$



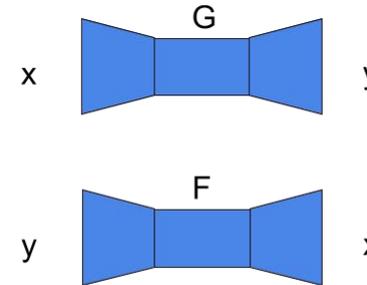
$$\mathcal{L}(G, F, D_X, D_Y) = \mathcal{L}_{GAN}(G, D_Y, X, Y) + \mathcal{L}_{GAN}(F, D_X, Y, X) + \lambda \mathcal{L}_{cyc}(G, F)$$

Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017



- Two encoder-decoder networks are jointly trained.

$$F \circ G: X \rightarrow X \text{ ve } G \circ F: Y \rightarrow Y$$

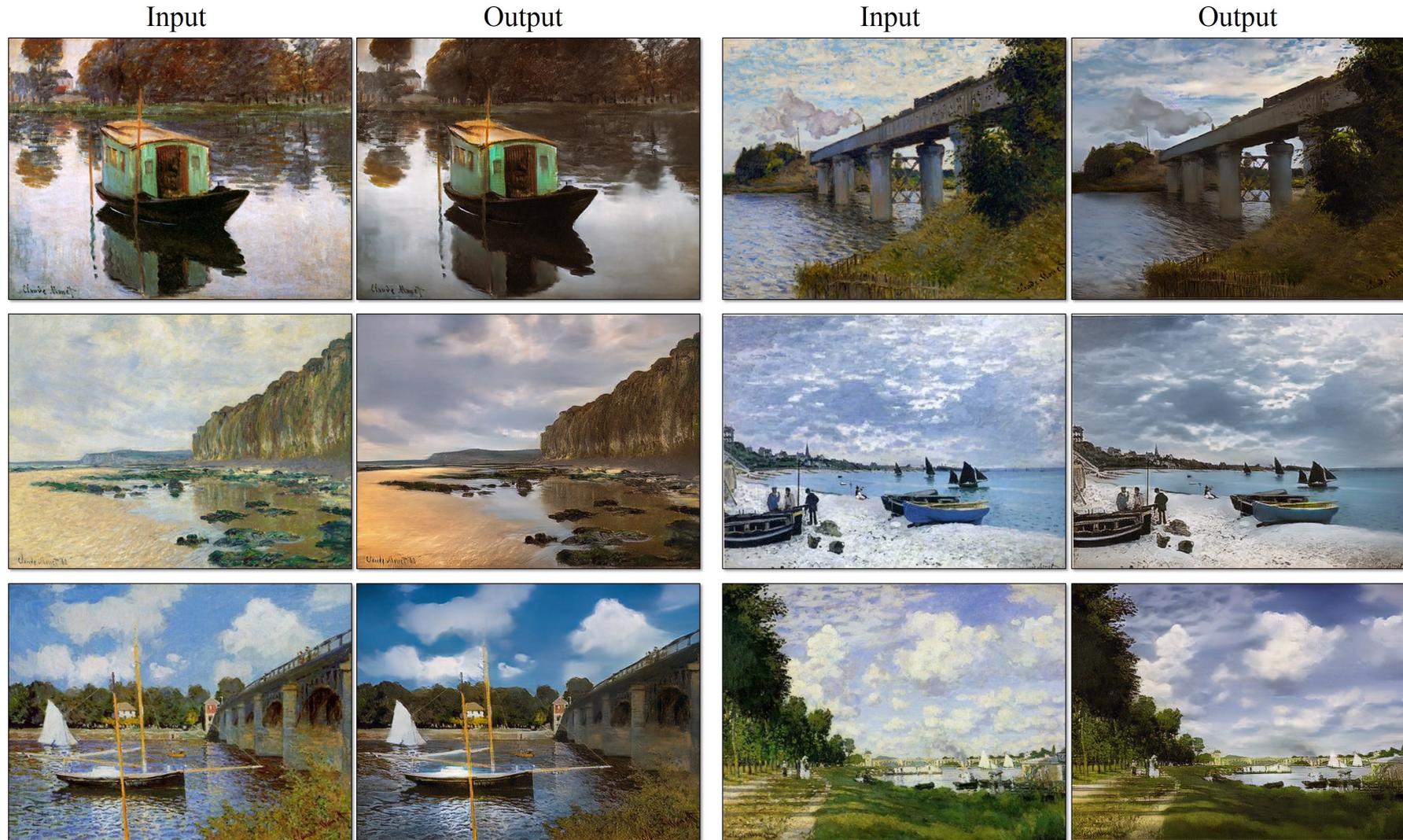


- 70×70 PatchGAN, which try to classify whether 70×70 overlapping image patches are real or fake is used.
- Adversarial training.

Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017

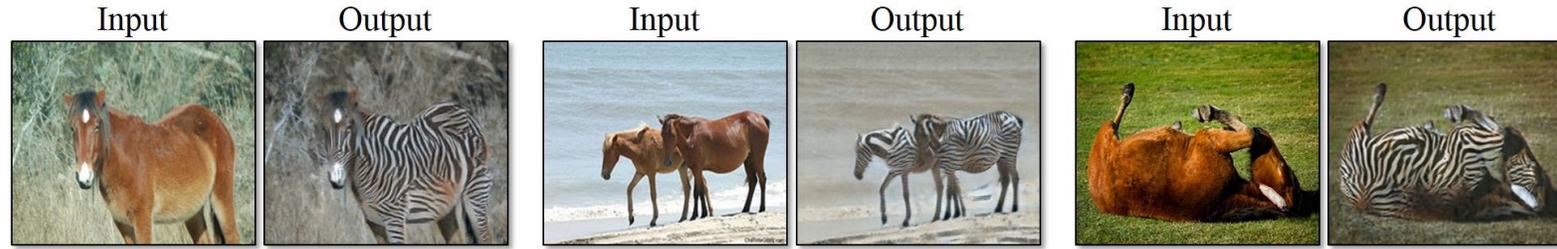


Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017



Unpaired Image to Image Translation(CycleGAN)

Zhu vd. 2017



horse → zebra



zebra → horse



apple → orange

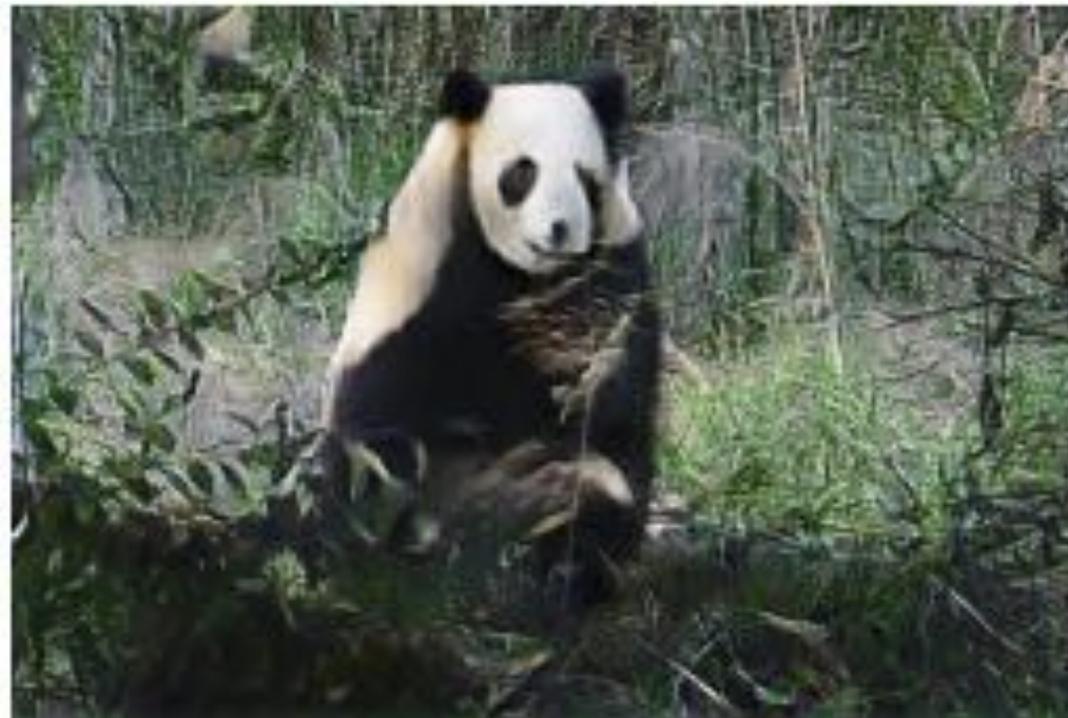


orange → apple

Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017



Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017



Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017



Source: <https://github.com/tatsuyah/CycleGAN-Models>

Unpaired Image to Image Translation(CycleGAN) Zhu vd. 2017



Source: <https://github.com/tatsuyah/CycleGAN-Models>

Unpaired Image to Image Translation(CycleGAN)

Zhu vd. 2017

A failure case



Neural Face Editing with Intrinsic Image Disentangling Shu et al. 2017

- An end-to-end GAN that infers a face-specific disentangled representation of intrinsic face properties.

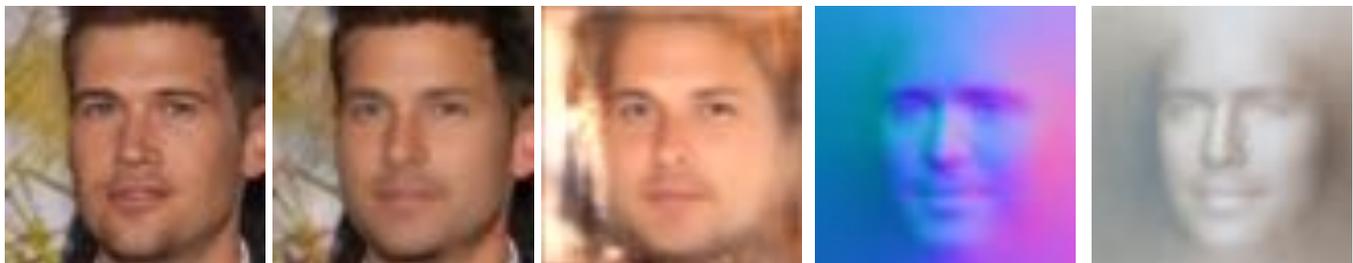
- Shape
- Albedo
- Lighting
- Alpha matte

- A given face image I_{fg} is the result of a rendering process: $f_{rendering}$

$$I_{fg} = f_{rendering}(A_e, N_e, L)$$

$$I_{fg} = f_{image-formation}(A_e, S_e) = A_e \odot S_e$$

$$S_e = f_{shading}(N_e, L)$$



(a) input (b) recon (c) albedo (d) normal (e) shading



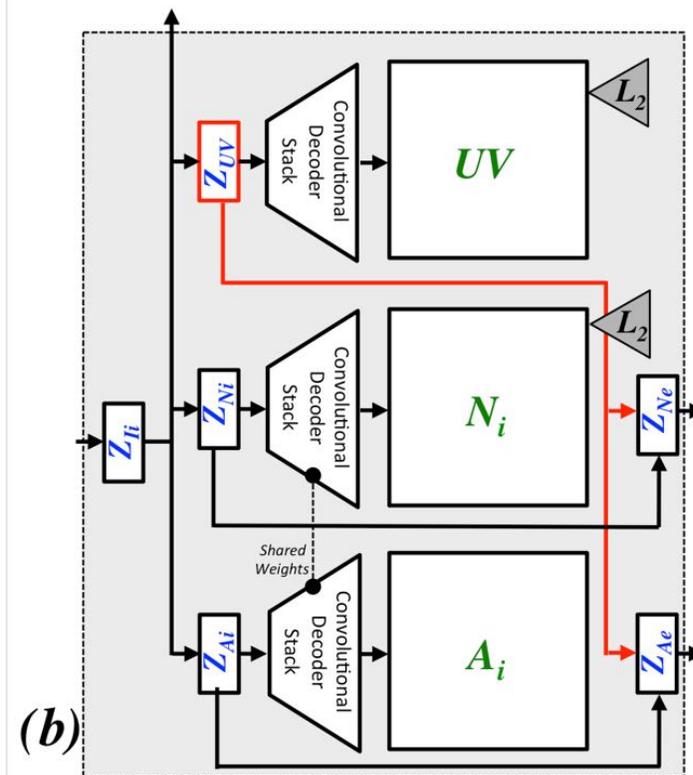
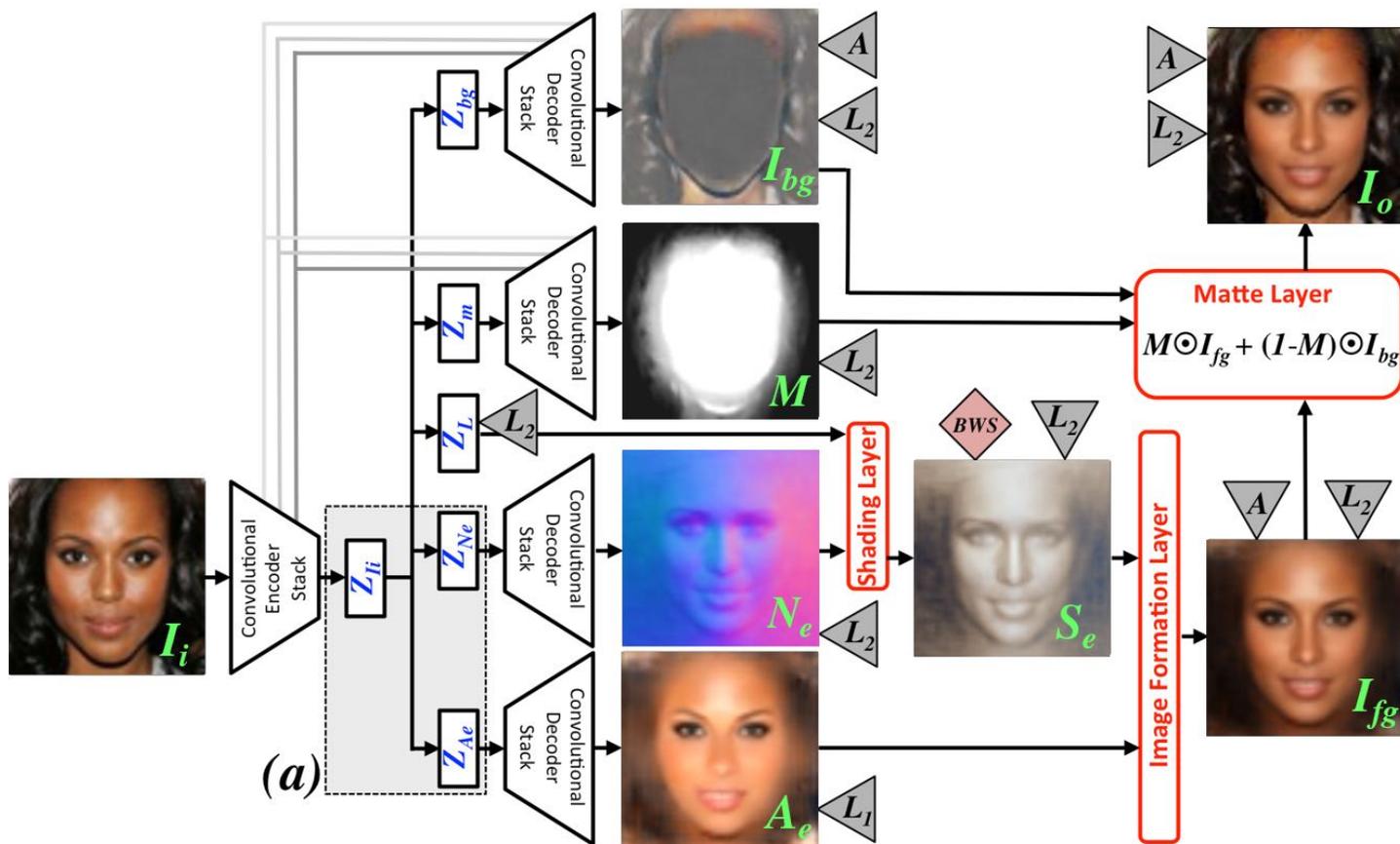
(f) relit (g) smile (h) beard (i) eyewear (j) older

Neural Face Editing with Intrinsic Image Disentangling Shu et al. 2017

$$I_{fg} = f_{rendering}(A_e, N_e, L)$$

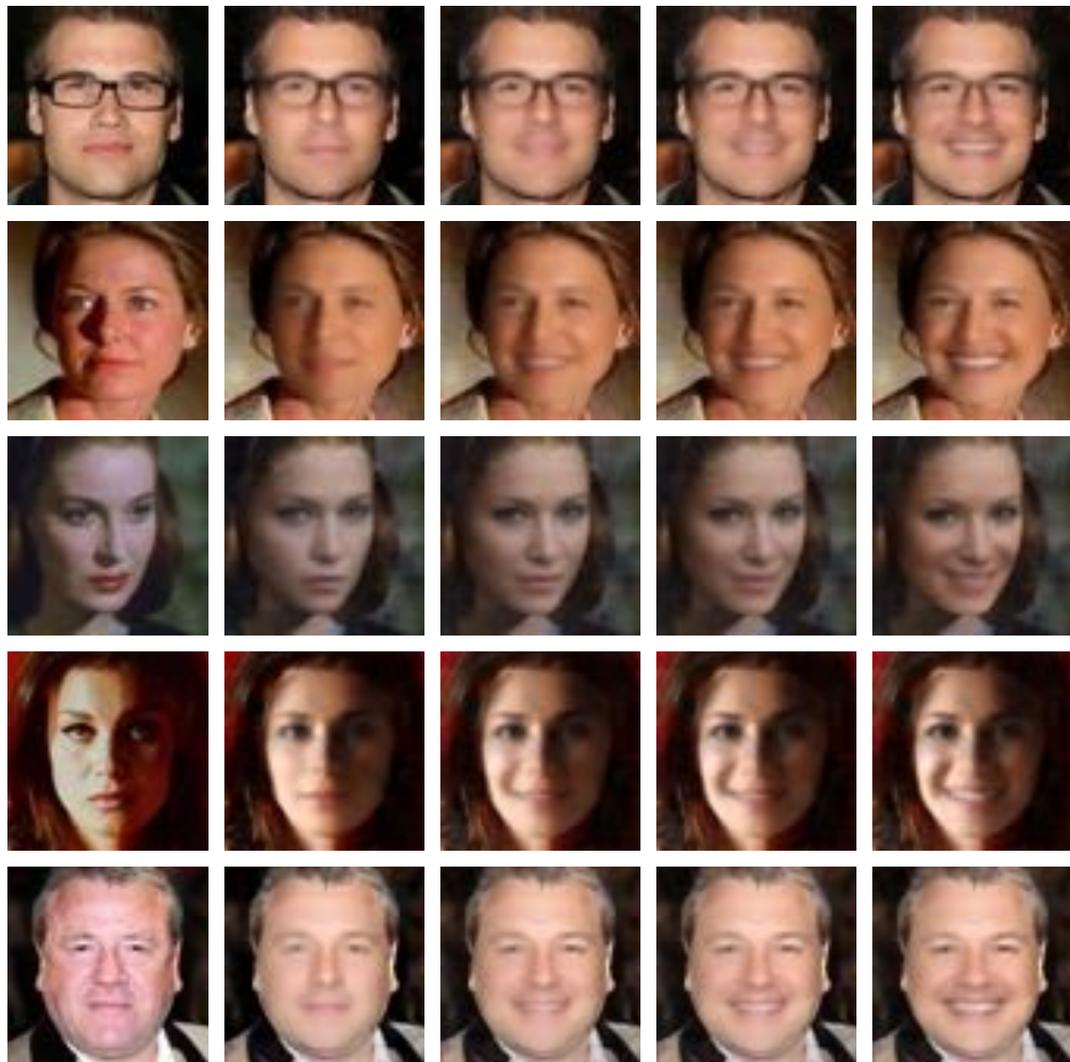
$$I_{fg} = f_{image-formation}(A_e, S_e) = A_e \odot S_e$$

$$S_e = f_{shading}(N_e, L)$$



Neural Face Editing with Intrinsic Image Disentangling Shu et al. 2017

Smiling



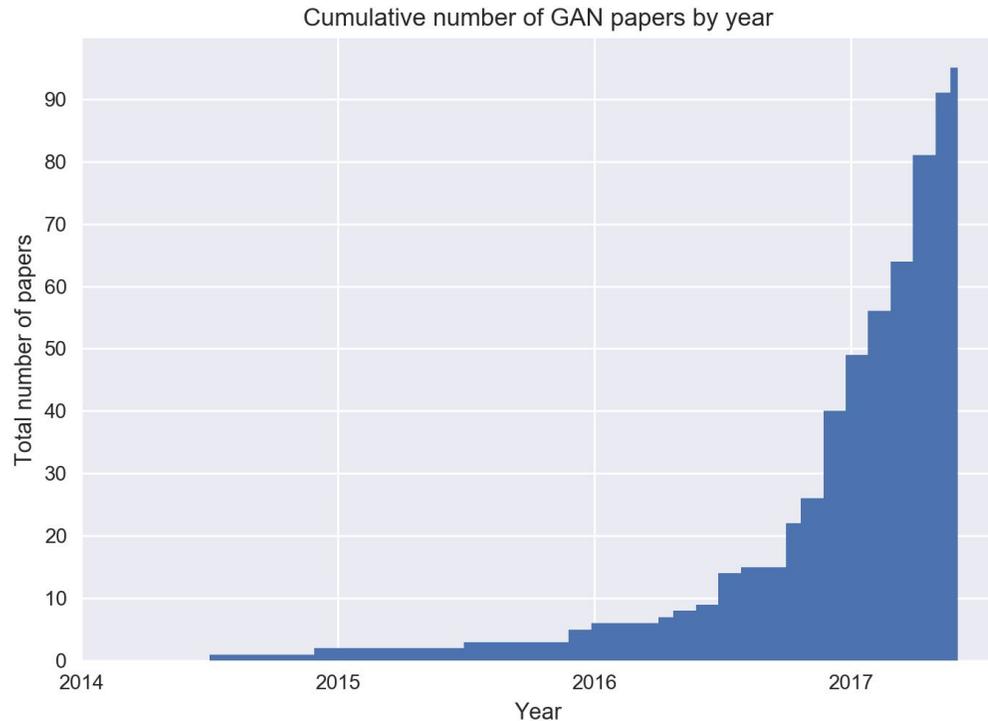
(a) input (b) recon (c) (d) (e)

Aging



(a) input (b) recon (c) (d) (e)

Conclusion



- Every week, new GAN papers are coming out.
- Very active topic in Machine Learning and Computer Vision.
- Adversarial loss started to be used for different problems in new papers in premier conferences.
- It has big potential for other areas.