

Generative models

Outline

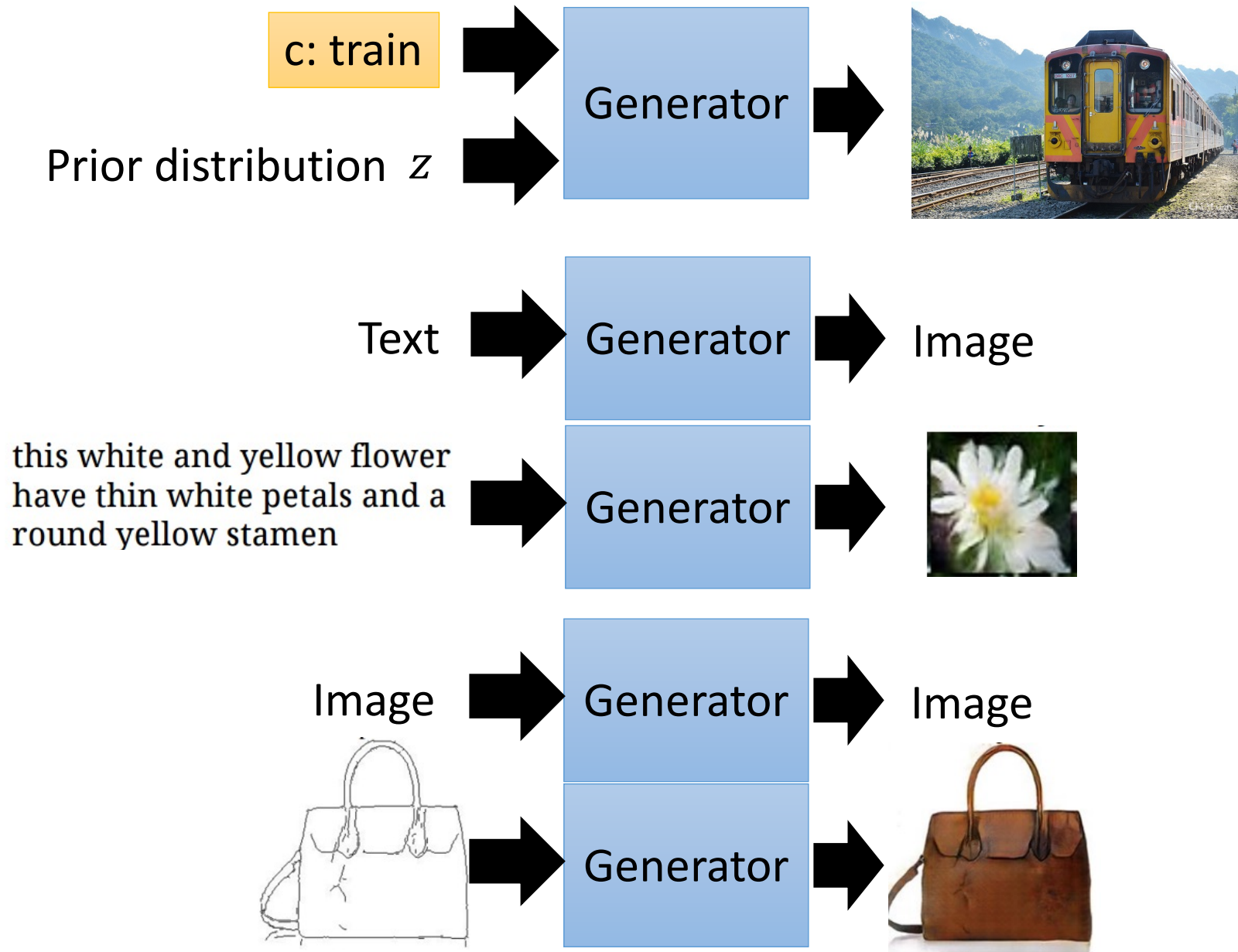
1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
- 5. Conditional GANs**

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
- 5. Conditional GANs**
 - 1. Principle**

Motivation



Conditional GAN

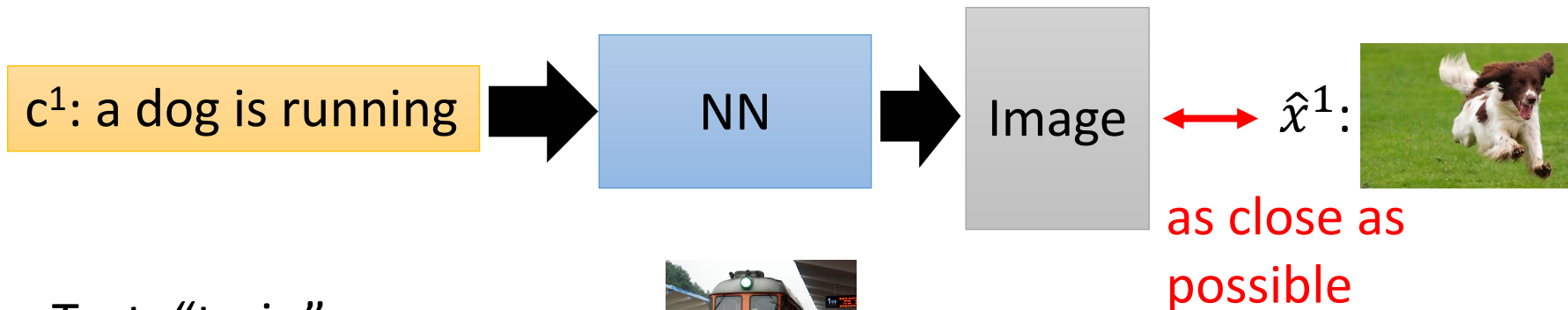
c^1 : a dog is running \hat{x}^1 :



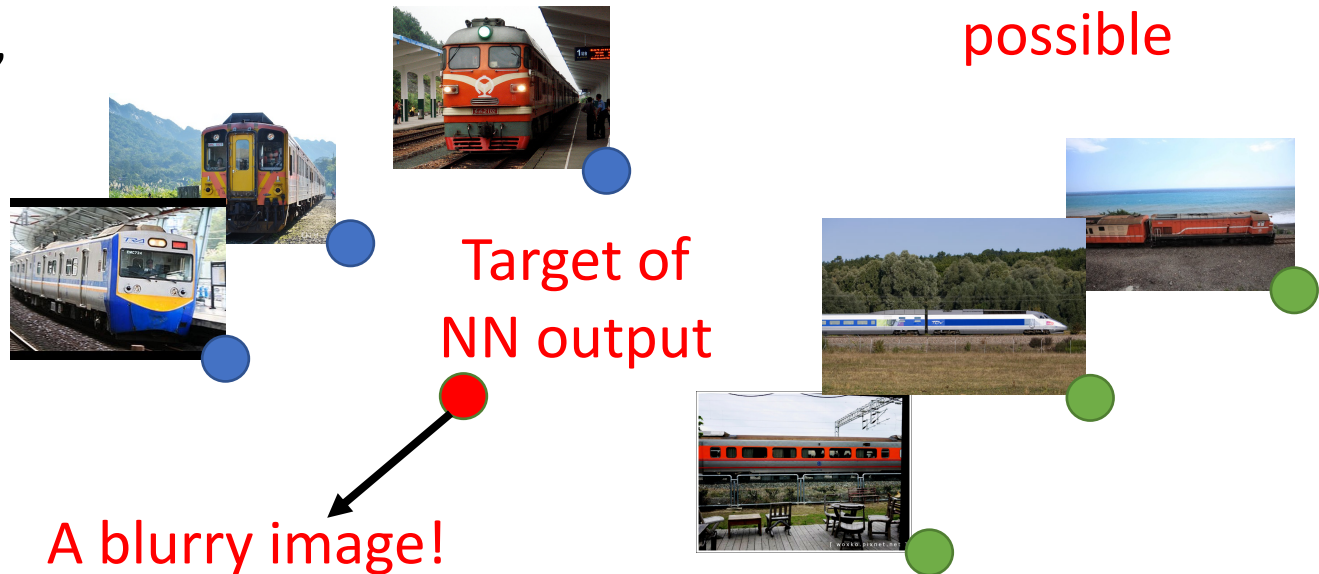
c^2 : a bird is flying \hat{x}^2 :



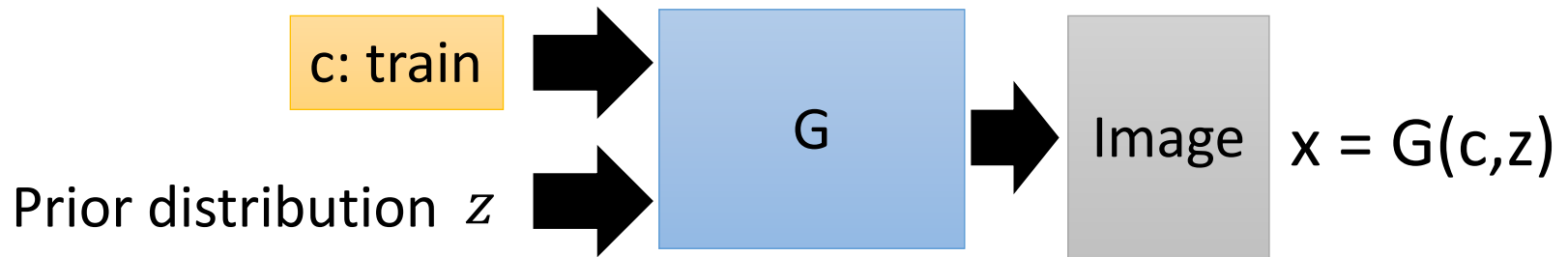
- **Text to image** by traditional supervised learning



Text: "train"



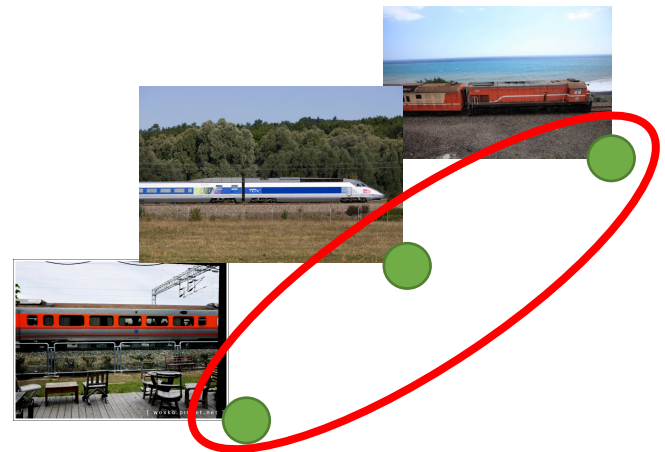
Conditional GAN



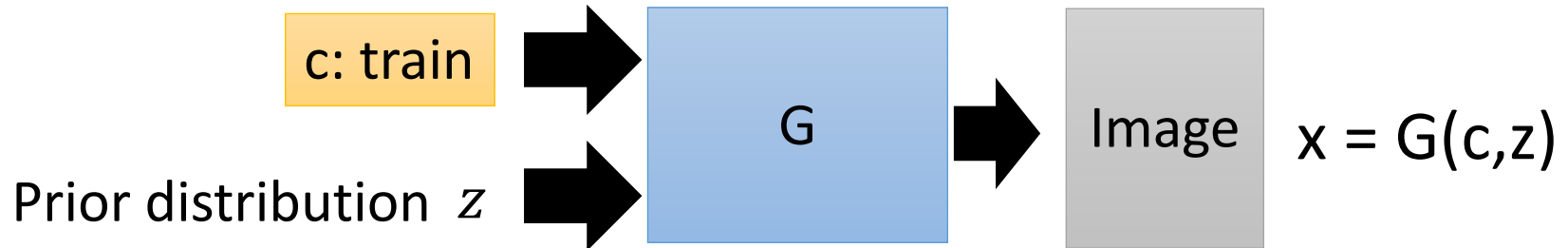
It is a distribution

Approximate the
distribution of real data

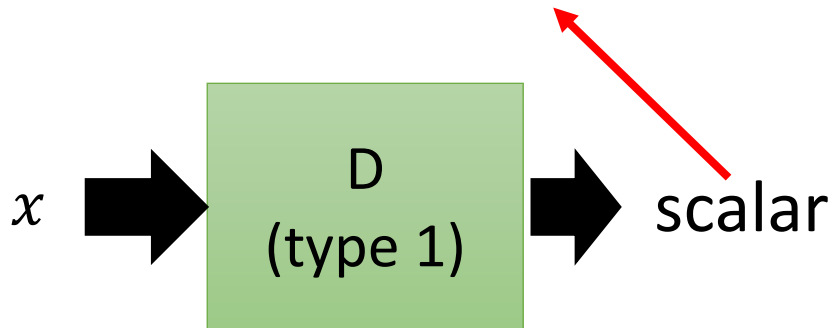
Text: "train"





Conditional GAN



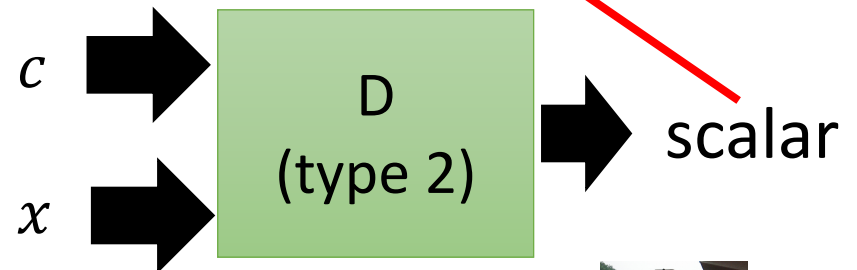
x is realistic or not



Positive example: 

Negative example: 

x is realistic or not +
 c and x are matched or not



Positive example: (train, )

Negative example: (train, )

Extra neg

(cat, )

Conditional GAN (cGAN model)

GAN

$$V(\textcolor{red}{G}, \textcolor{green}{D}) = \mathbb{E}_{x \sim P_{data}} [\log \textcolor{green}{D}(x)] + \mathbb{E}_{\textcolor{red}{x} \sim P_{\textcolor{red}{G}}} [\log(1 - \textcolor{green}{D}(x))]$$

$$\textcolor{red}{G}^* = \arg \min_{\textcolor{red}{G}} \max_{\textcolor{green}{D}} V(\textcolor{red}{G}, \textcolor{green}{D})$$

cGAN

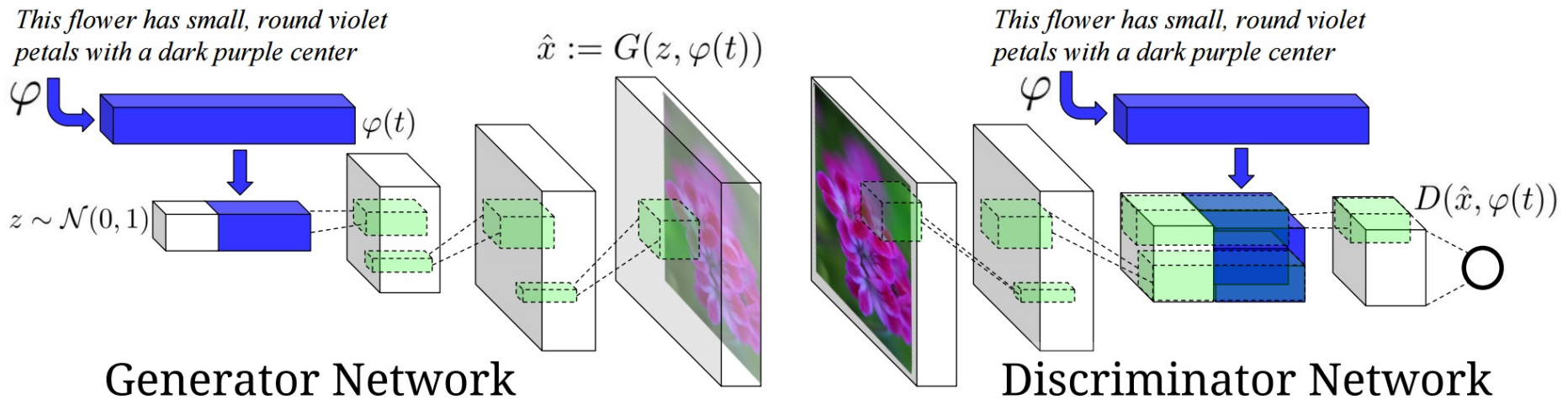
$$\min_G \max_D \left(\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{data}(\mathbf{x}, \mathbf{y})} [\log D(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{y} \sim p_{\mathbf{y}}, \mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z}, \mathbf{y}), \mathbf{y}))] \right)$$

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
 1. Principle
 2. **Text2Image**

Text2Image: architecture example



- Positive samples:
 - real image + right texts
- Negative samples:
 - fake image + right texts
 - Real image + wrong texts

Text2Image results

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



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



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



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

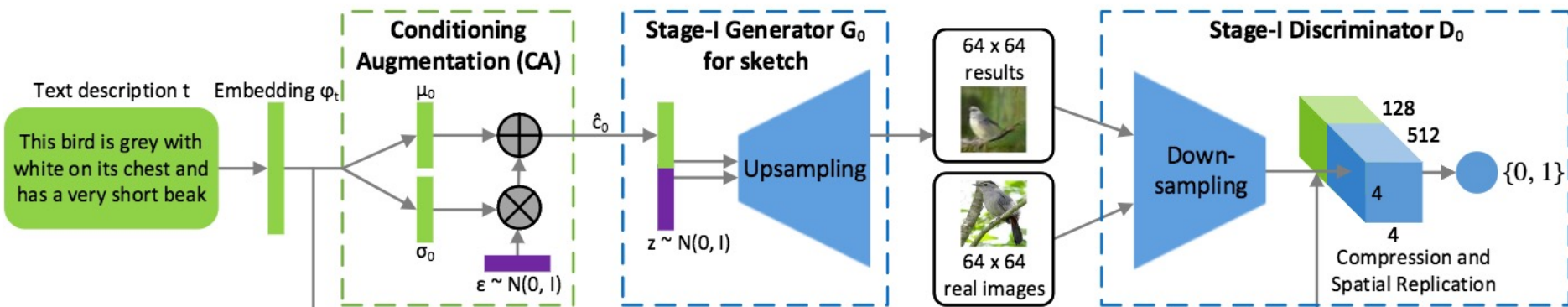


Text2Image results

Caption	Image
this flower has white petals and a yellow stamen	
the center is yellow surrounded by wavy dark purple petals	
this flower has lots of small round pink petals	

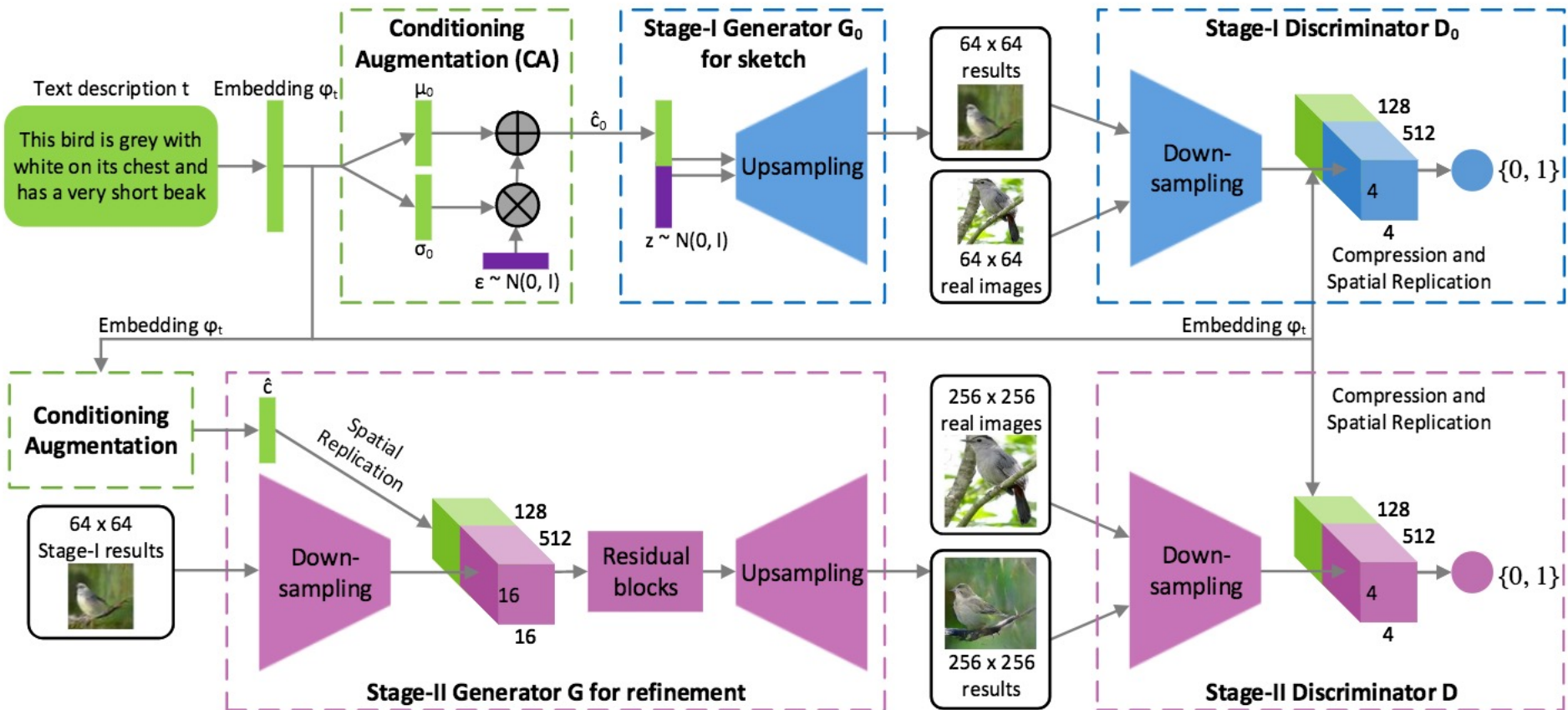
Text2Image: architecture example (2)

Generating higher resolution images (from 64 to 256)



Text2Image: architecture example (2)

Generating higher resolution images (from 64 to 256)



StackGAN results

This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face

This bird is white with some black on its head and wings, and has a long orange beak

This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments

(a) Stage-I images



(b) Stage-II images

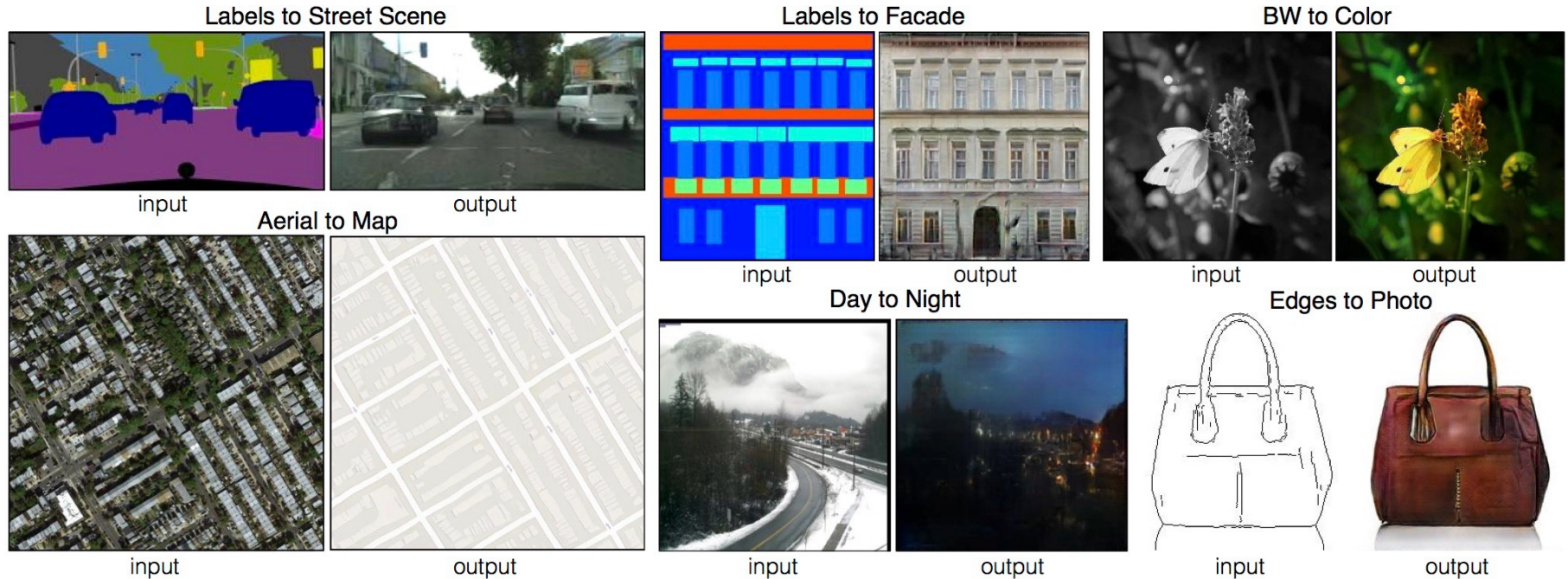


Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
 1. Principle
 2. Text2Image
 3. **Image2Image**

Image-based Conditional GAN



- Conditioned on an image of different modality
- Image-to-Image Translation => **pix2pix**

Image-to-image pix2pix

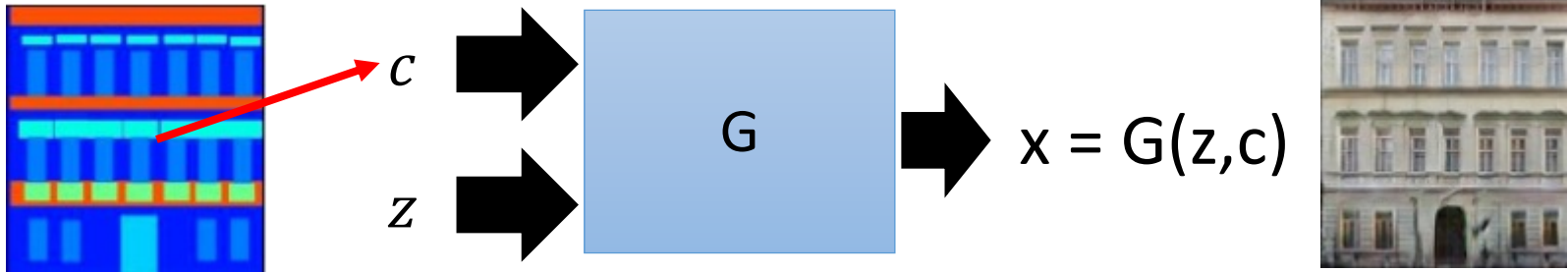
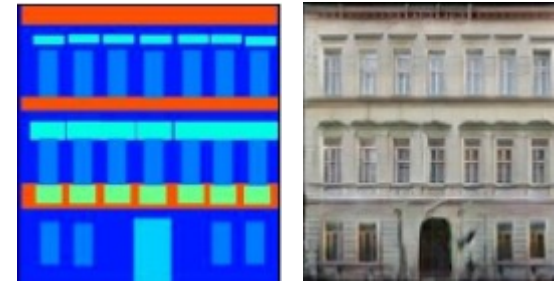
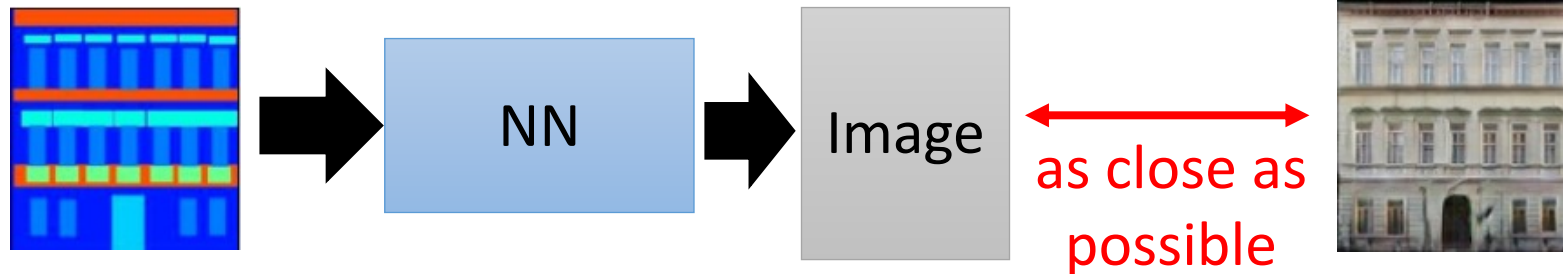


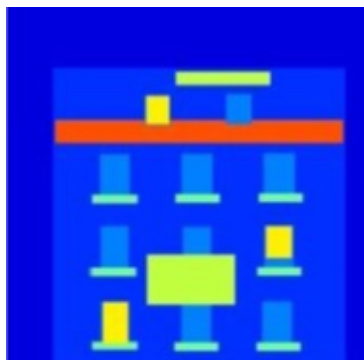
Image-to-image pix2pix



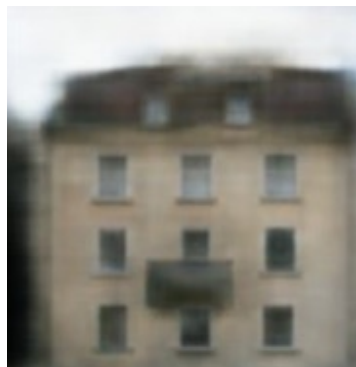
- Traditional supervised approach



Testing:



input



close

It is blurry
because it is
the average of
several images.

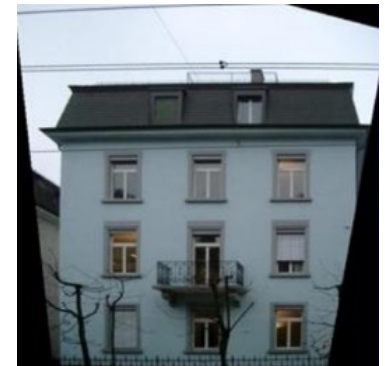
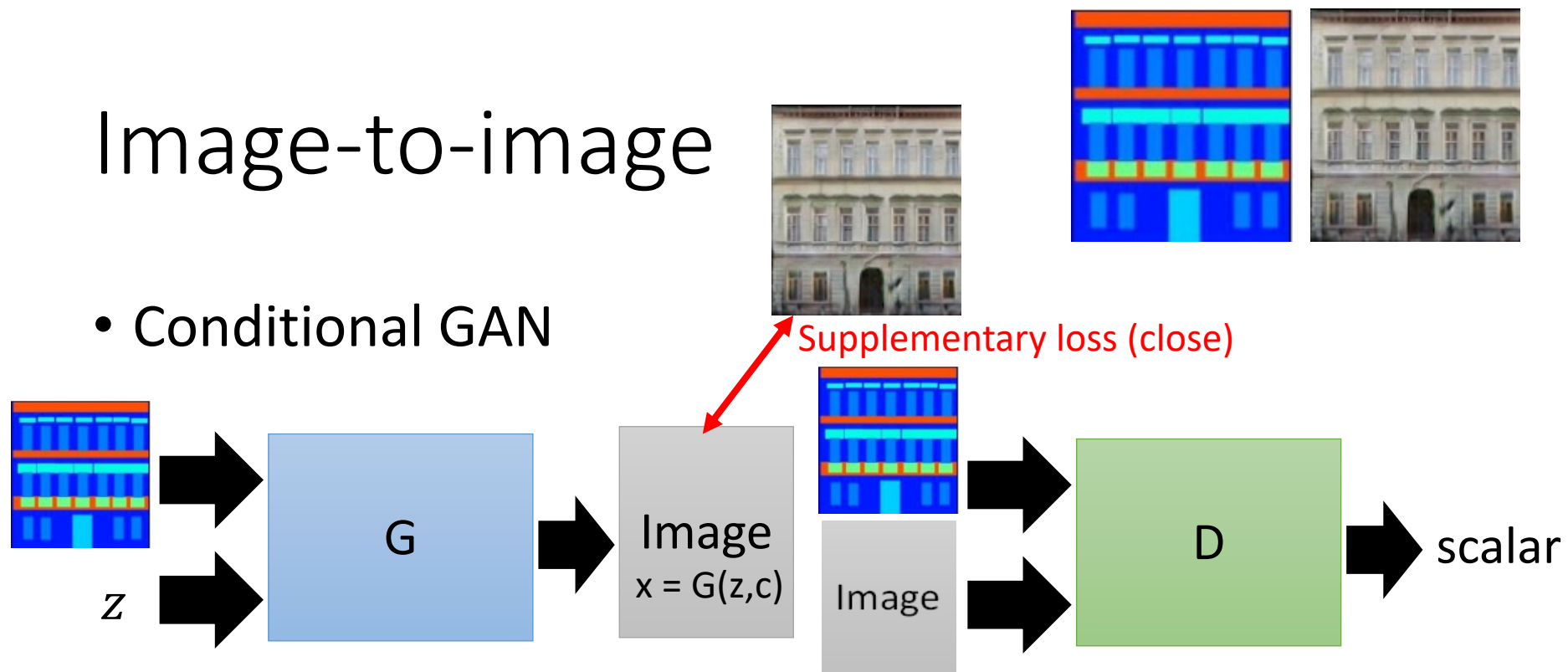


Image-to-image

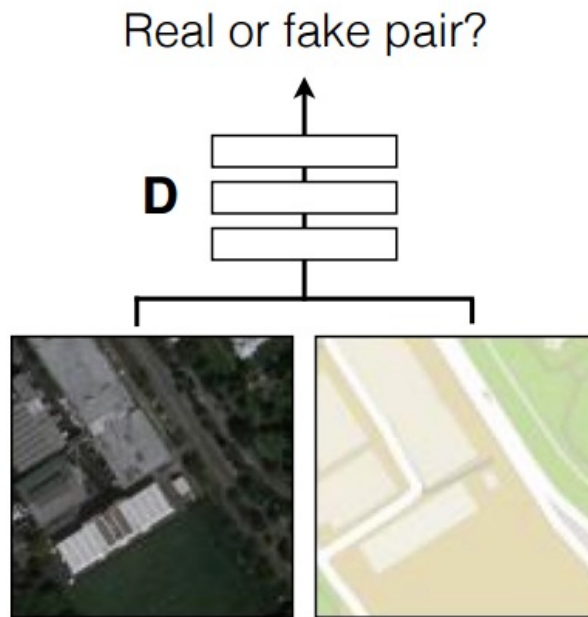
- Conditional GAN



Testing:



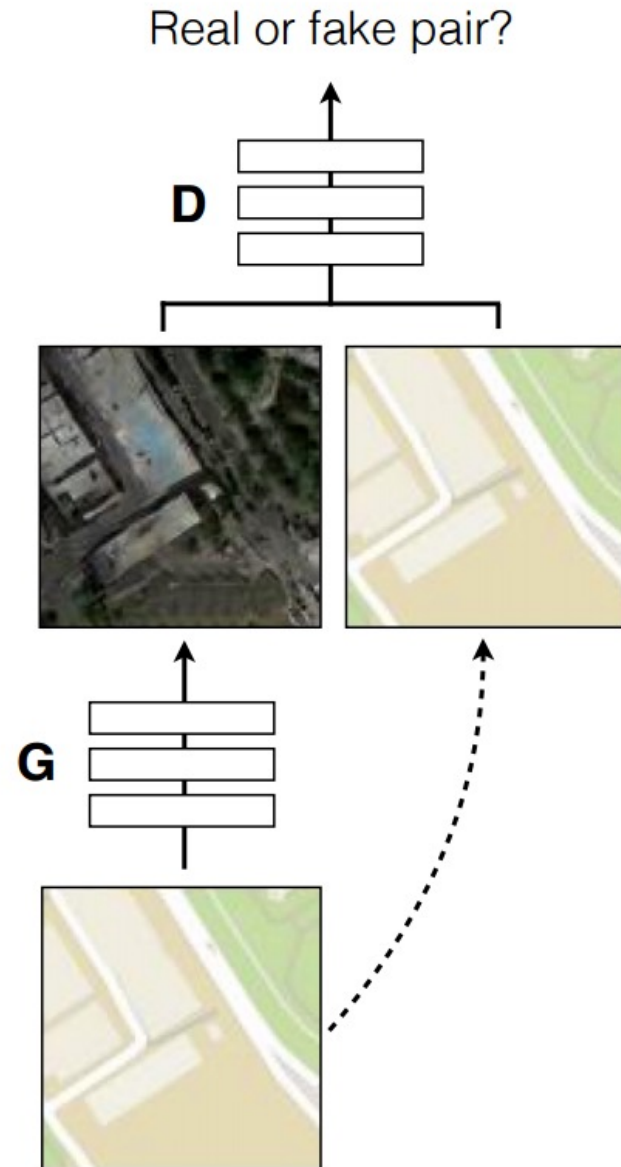
Positive examples



G tries to synthesize fake images that fool **D**

D tries to identify the fakes

Negative examples



Label2Image

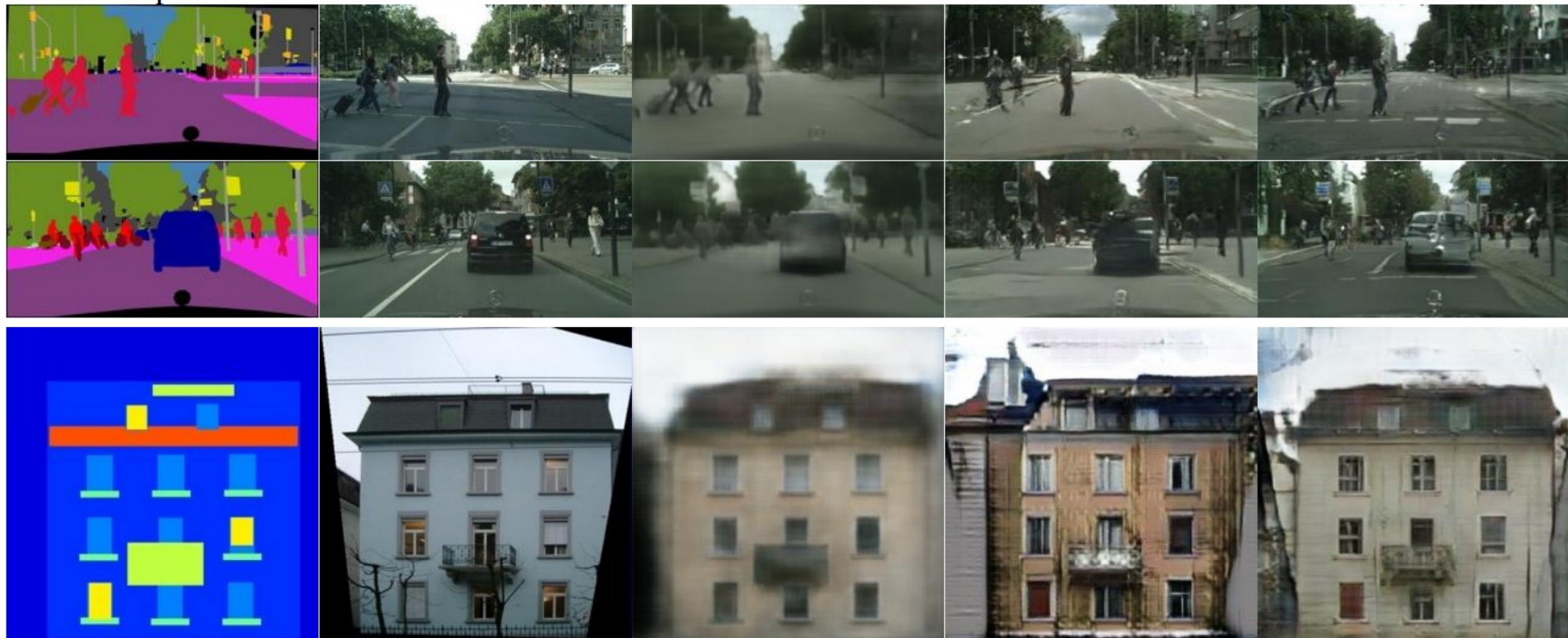
Input

Ground truth

L1

cGAN

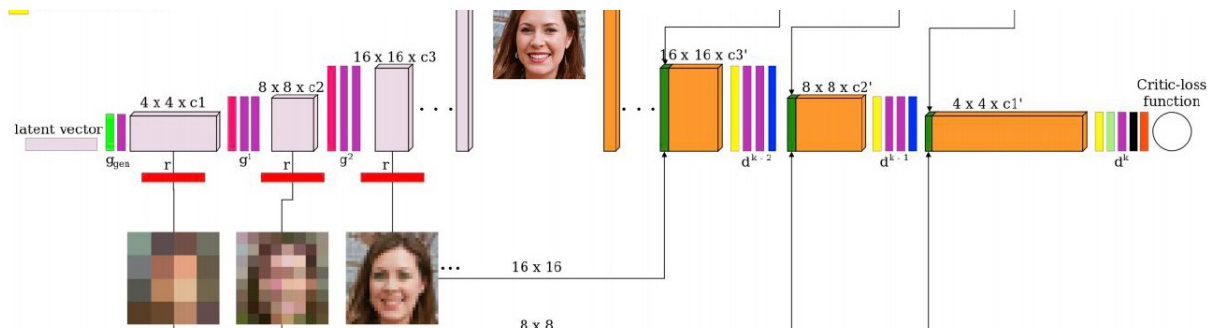
L1 + cGAN



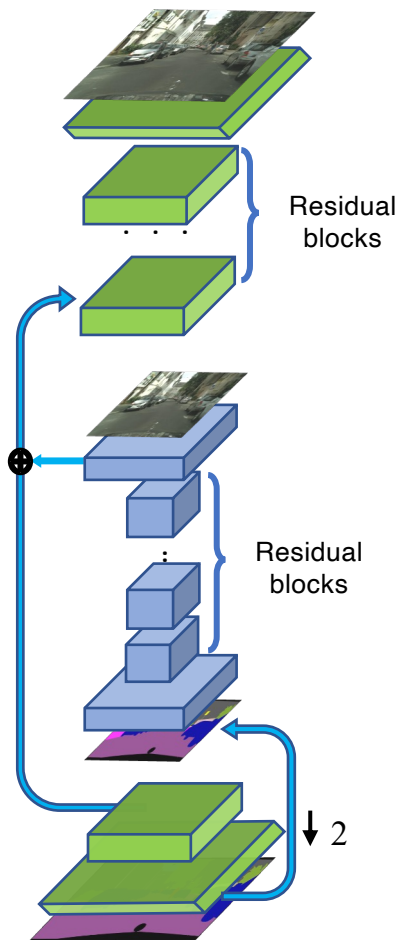
Edges2Image



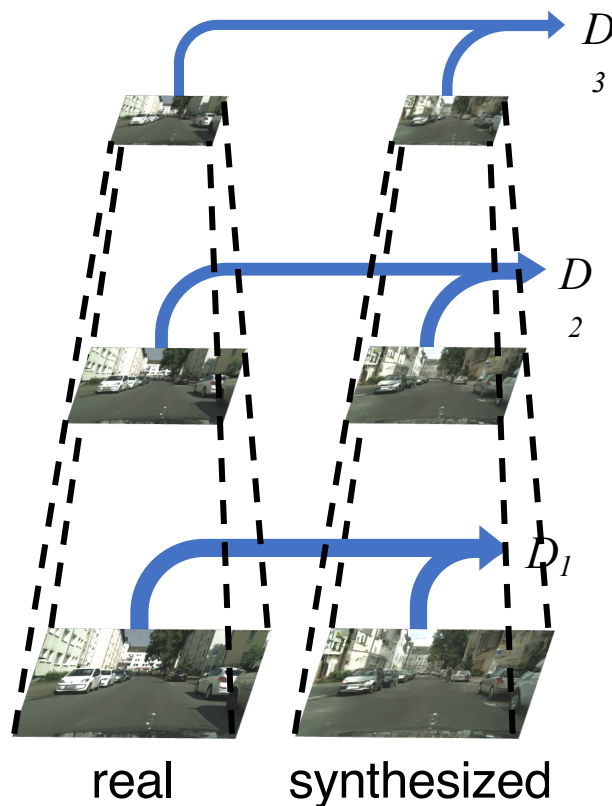
Pix2pixHD



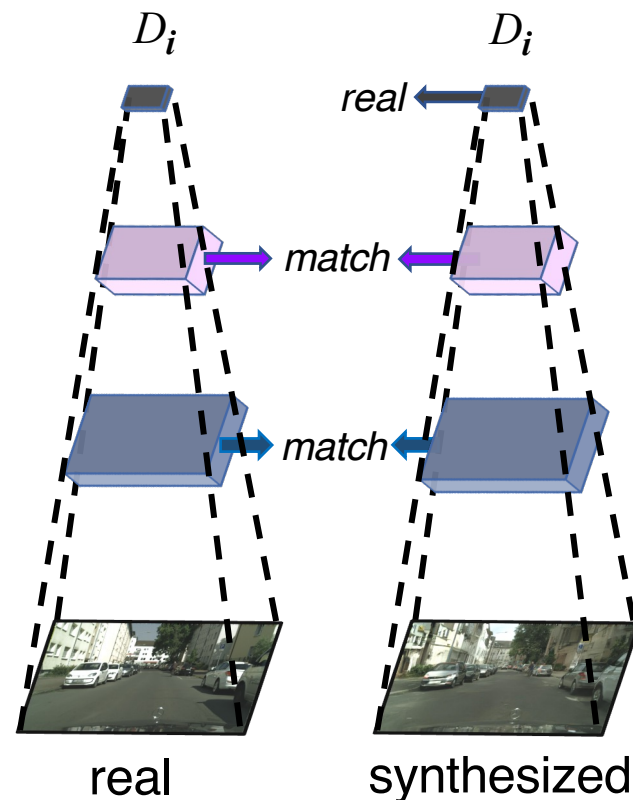
Coarse-to-fine Generator



Multi-scale Discriminators



Robust Objective





Semantic Map



pix2pix



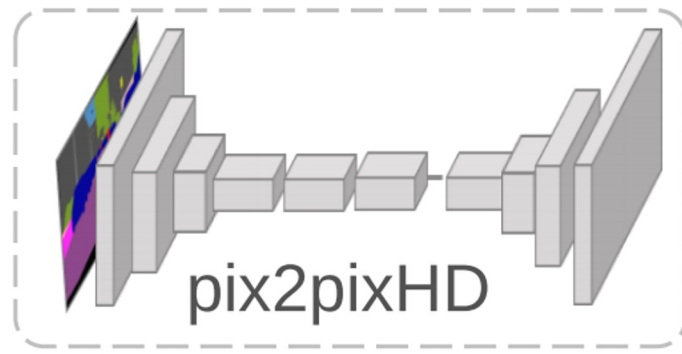
CRN



Ours

Improving Segmentation2Image strategy?

Limitation of the approach:



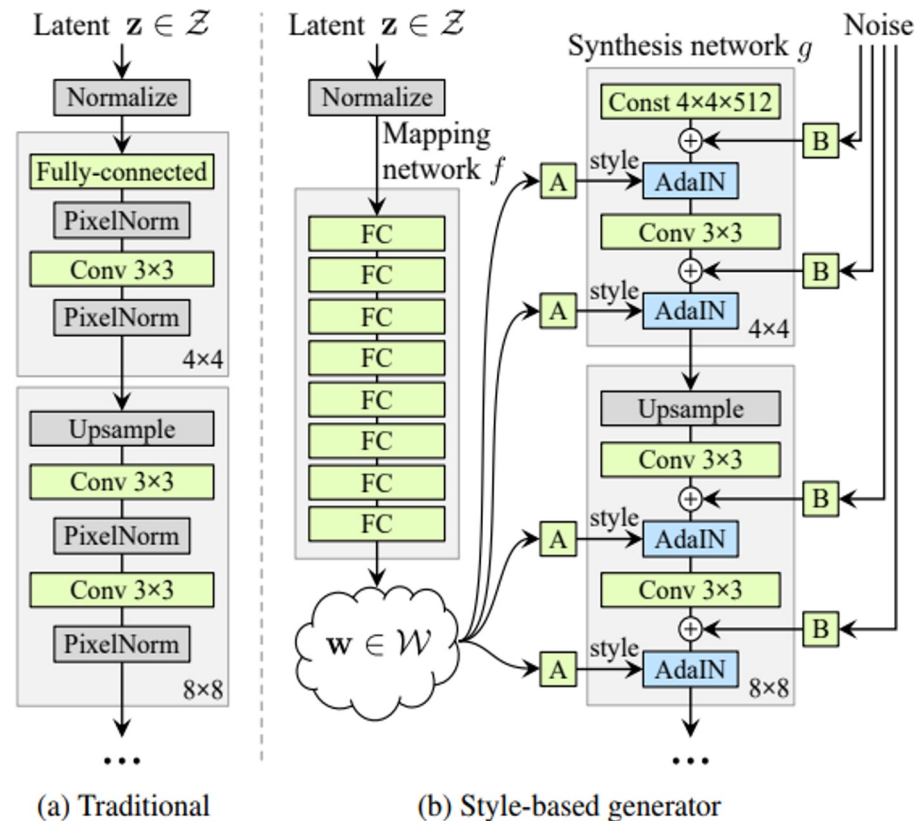
Directly feed the semantic layout as input to the deep network, which is processed through stacks of convolution, normalization, and nonlinearity layers.

However, this is suboptimal as the normalization layers tend to “wash away” semantic information in input semantic segmentation masks.

Improving Segmentation2Image strategy

Proven effective for recent generative adversarial networks such as **StyleGAN**

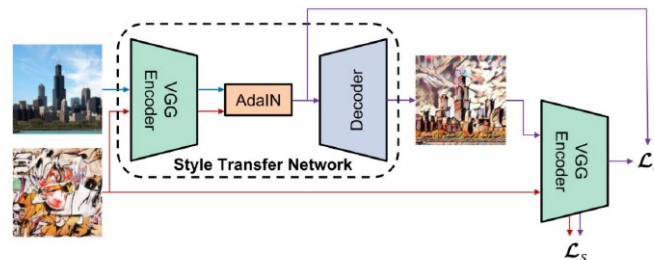
Can we do the same for conditional GAN?
Conditional Normalization Layers?



Improving Segmentation2Image strategy

Recall: Adaptive instance normalization

$$\text{AdaIN}(x, y) = \sigma(y) \left(\frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$

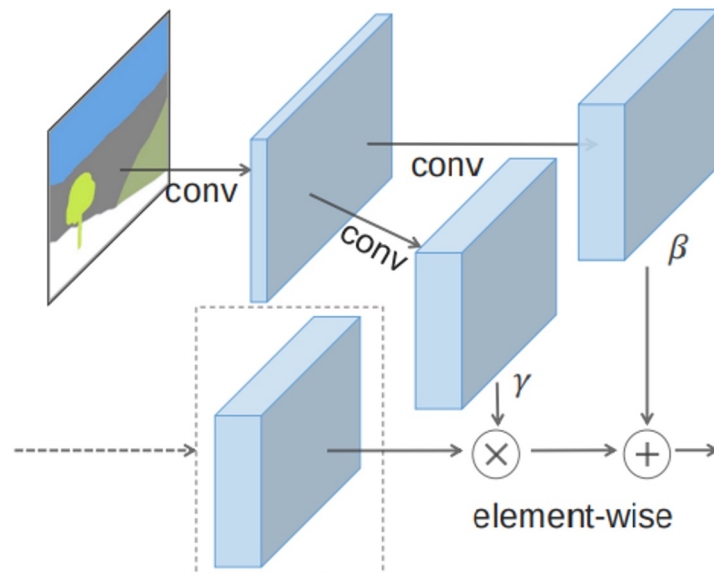


SPADE block= spatially-adaptive denormalization:
Same idea but per class c over each channel i (N =batch size)

$$\gamma_{c,y,x}^i(\mathbf{m}) \frac{h_{n,c,y,x}^i - \mu_c^i}{\sigma_c^i} + \beta_{c,y,x}^i(\mathbf{m})$$

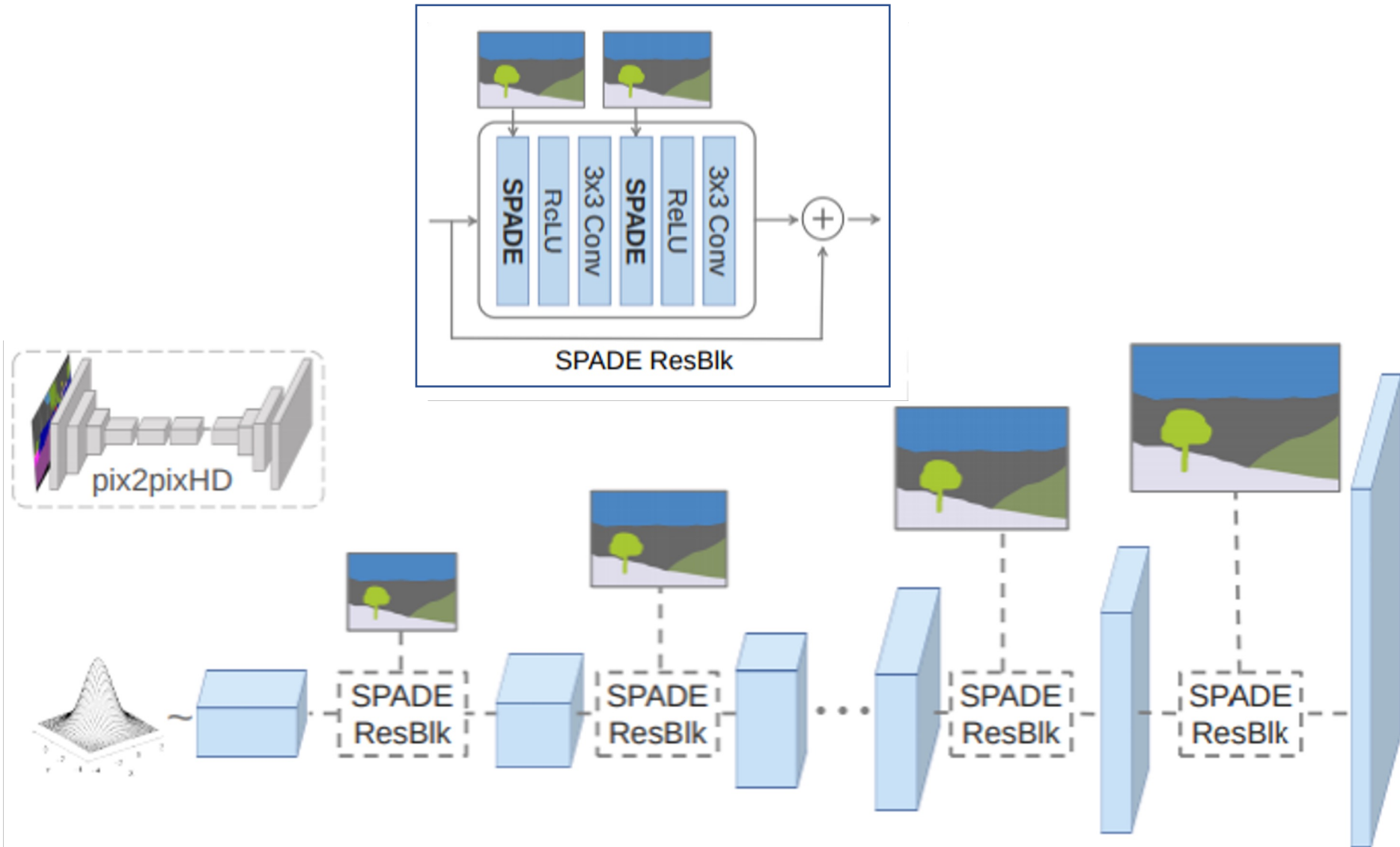
$$\mu_c^i = \frac{1}{NH^iW^i} \sum_{n,y,x} h_{n,c,y,x}^i$$

$$\sigma_c^i = \sqrt{\frac{1}{NH^iW^i} \sum_{n,y,x} (h_{n,c,y,x}^i)^2 - (\mu_c^i)^2}$$

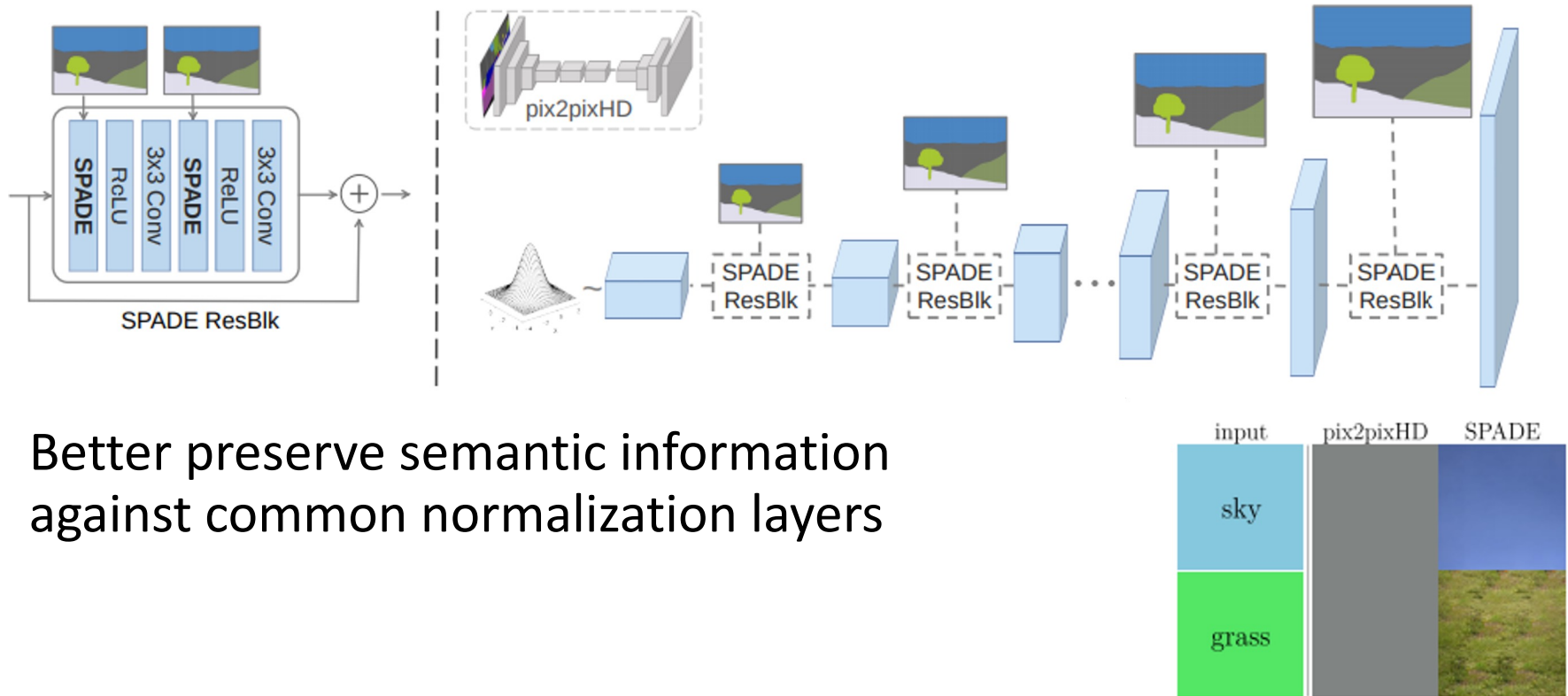


SPADE paper = [Semantic Image Synthesis with Spatially-Adaptive Normalization CVPR 2019]

SPADE Generator

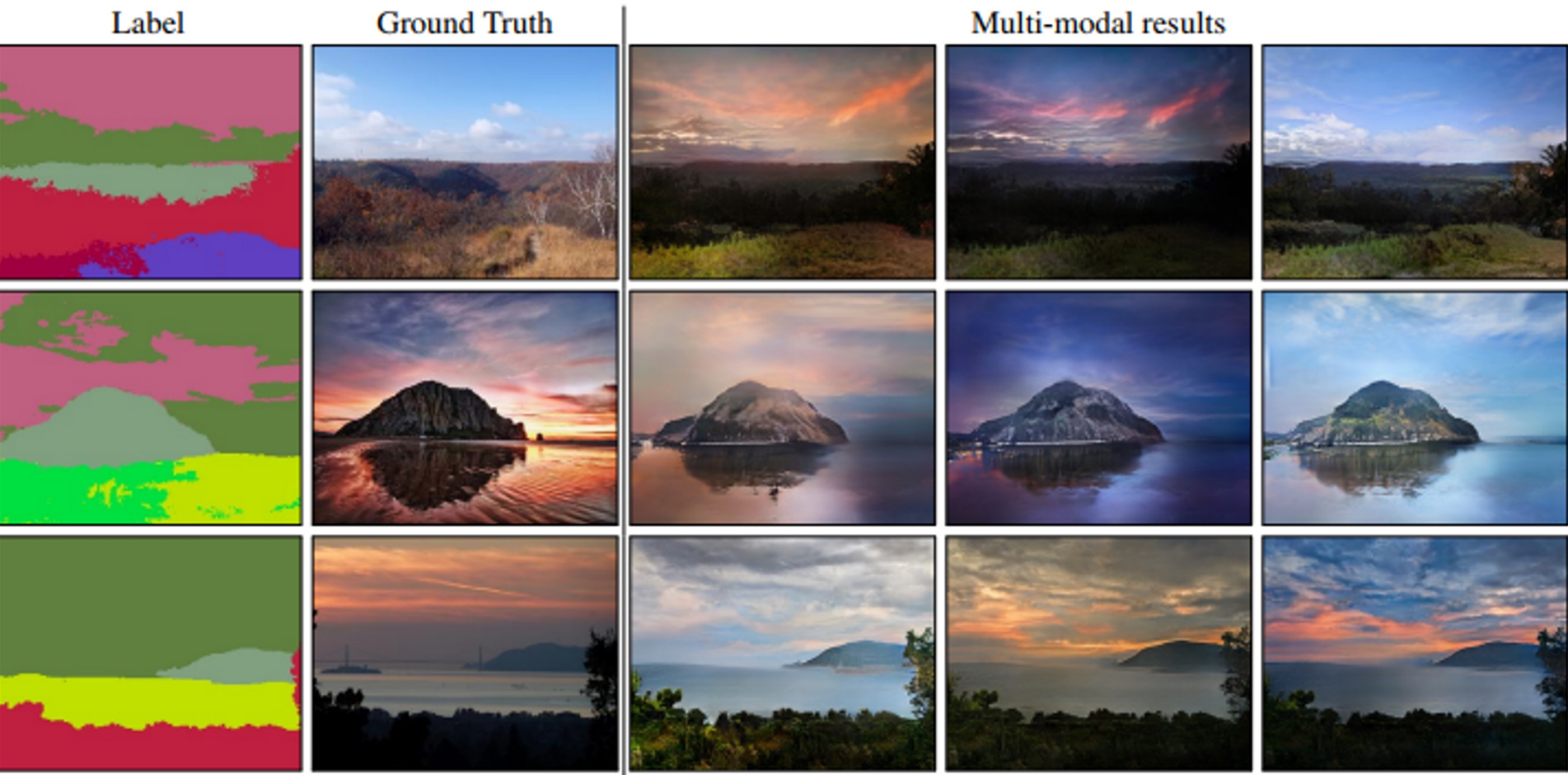


SPADE Generator



Better preserve semantic information
against common normalization layers

SPADE results



Spade and follow-up approaches

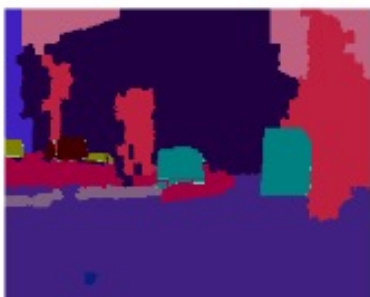
Label map

Ground truth

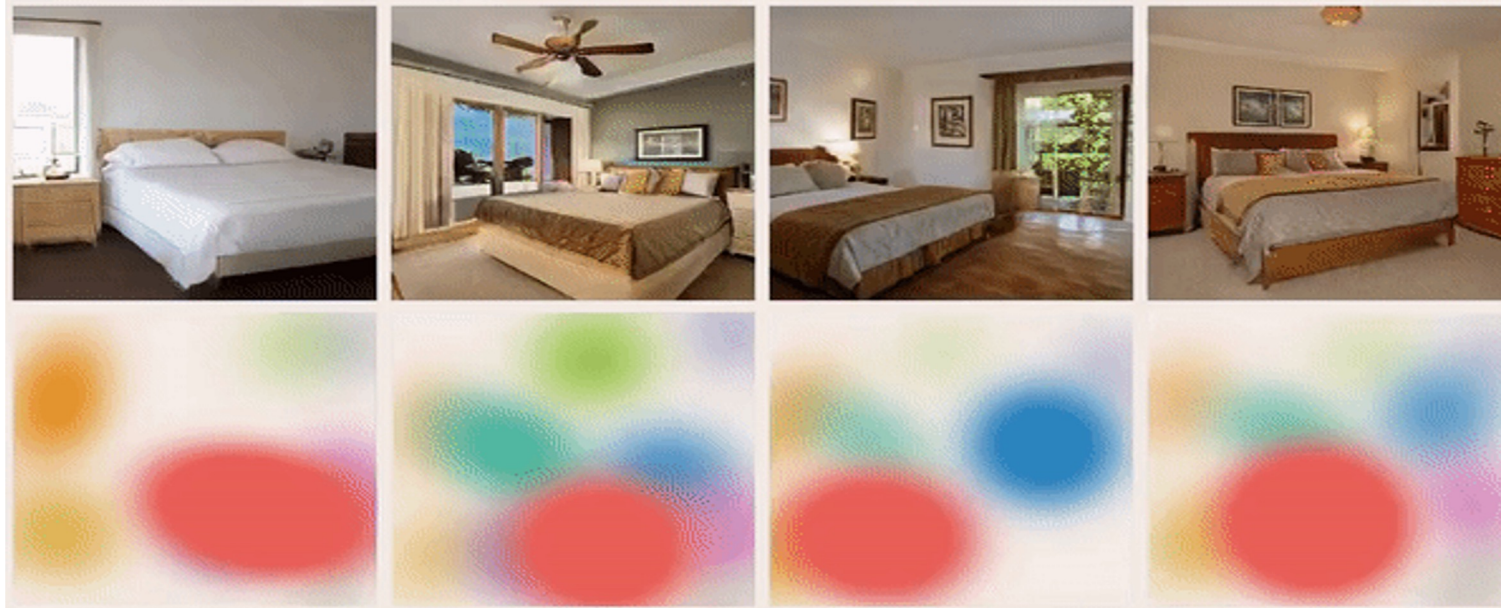
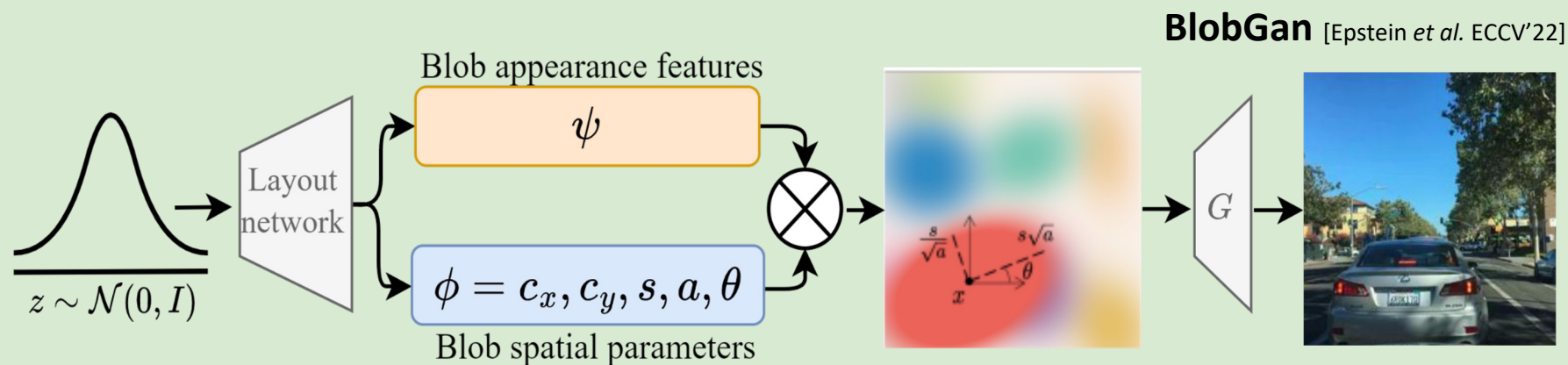
SPADE

CC-FPSE

OASIS

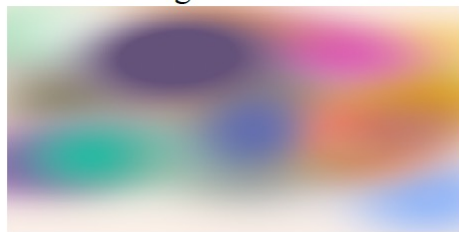


Editing with *conditional-or-structured-latent* GANs



Editing with *conditional-or-structured-latent* GANs

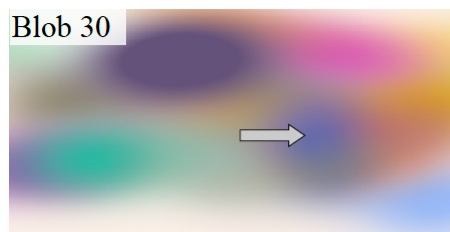
Original blobs



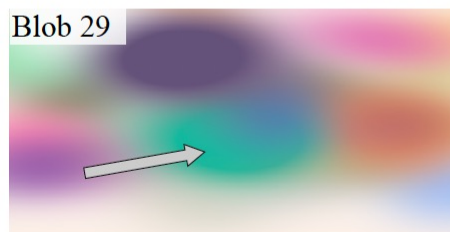
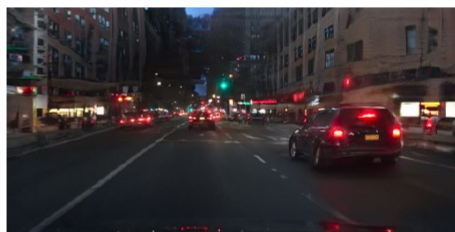
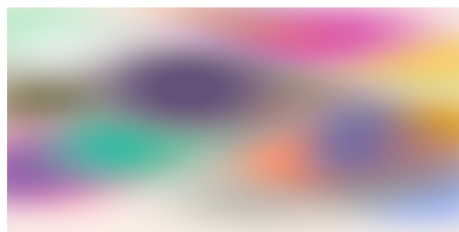
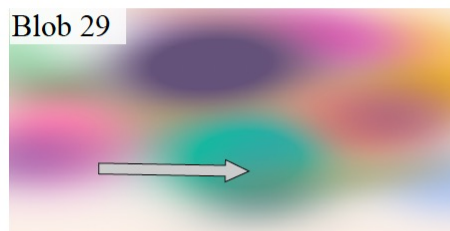
Original image



Edited blobs



Edited image



Editing with *conditional-or-structured-latent* GANs

Example of Counterfactual optimization for editing



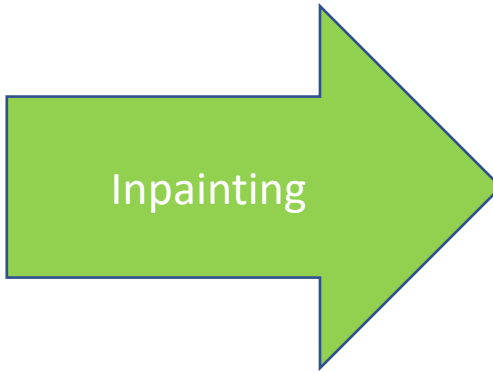
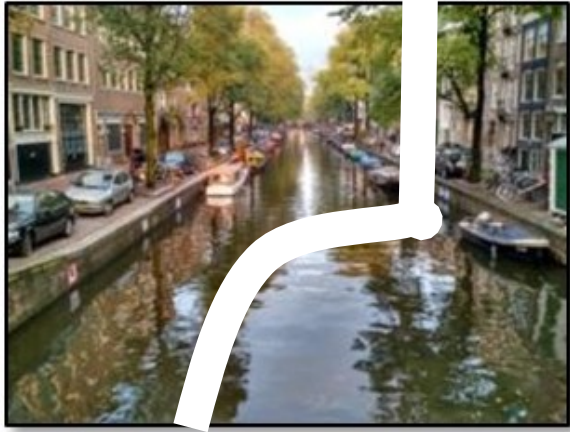
Generative models

Outline

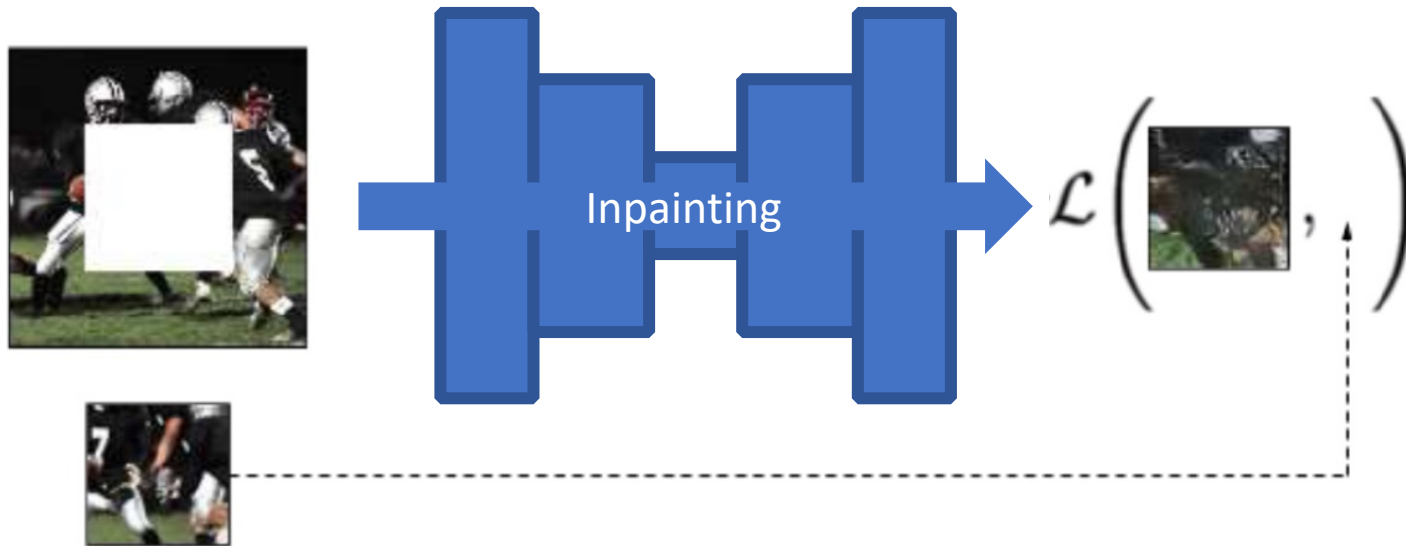
1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
 1. Principle
 2. Text2Image
 3. Image2Image
 4. **Inpainting and general missing data encoder**

Inpainting task

- Complete the missing part



Inpainting as unsupervised learning with GAN loss

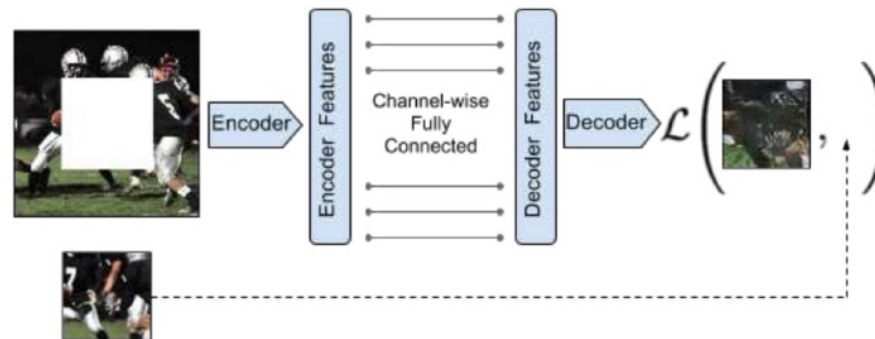


Reconstruct missing pixels by decoding using context

Loss defined on the predicted patch and the real one (known at training time)

First proposition -- Architecture

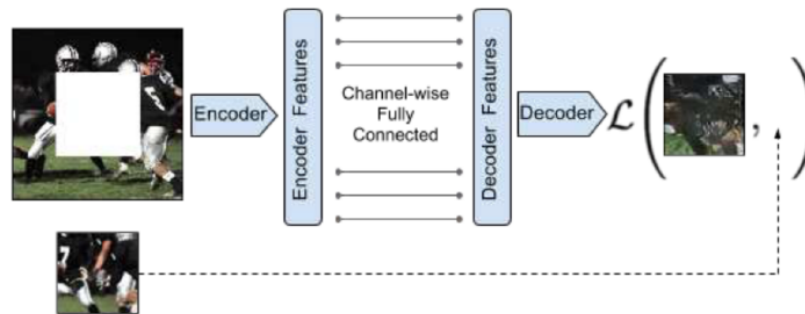
- Architecture: Encoder/Fully connected/Decoder



- DC-GAN for inpainting task
- **Input:** $227 \times 227 \times 3$ image
- **Output:** encoder context features ($6 \times 6 \times 256$)

Channel-wise fully-connected layer

- **Input / output:** $6 \times 6 \times 256$ channels
- **First layer:** Channel-wise fully-connected
(each 6×6 input connected to the corresponding 6×6 output)
- **Second layer:** Stride 1 convolution to mix channels

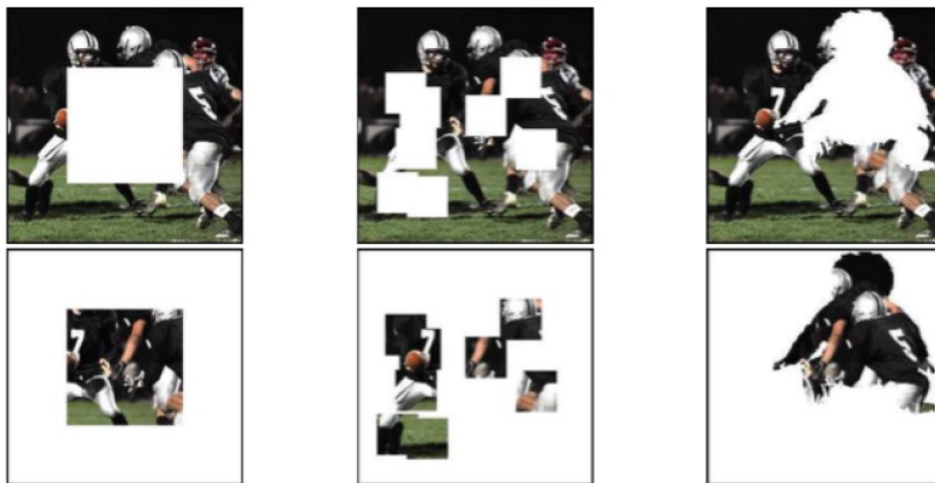


Decoder

- **Architecture:** Same as DC-GAN: 5 up-convolutional layers (“deconv” + ReLU)
- **Input:** decoder context features $6 \times 6 \times 256$
- **Output:** $227 \times 227 \times 3$ image

Training: Masking the images

- **How to define the mask ?**
 - ▶ Center region of the image
 - ▶ Random regions (chosen solution)
 - ▶ Random segmentation mask from VOC (said to be equivalent to random regions)
- **Formal definition:** Defined by a mask $\hat{M} \in \{0, 1\}^{227 \times 227}$ with 1 if the pixel should be masked



Training: Loss - Overview

- Trained completely from scratch to fill-up the masked areas
- **Problem:** multiple plausible solutions
- **Solution:** combining 2 losses:
 - ▶ \mathcal{L}_{rec} **L2 reconstruction loss:** learn the structure of the missing region (average multiple modes in prediction)
 - ▶ \mathcal{L}_{adv} **Adversarial loss:** make it look real (pick a mode from the distribution)

$$\min_F \mathcal{L} = \lambda_{rec} \mathcal{L}_{rec} + \lambda_{adv} \mathcal{L}_{adv}$$

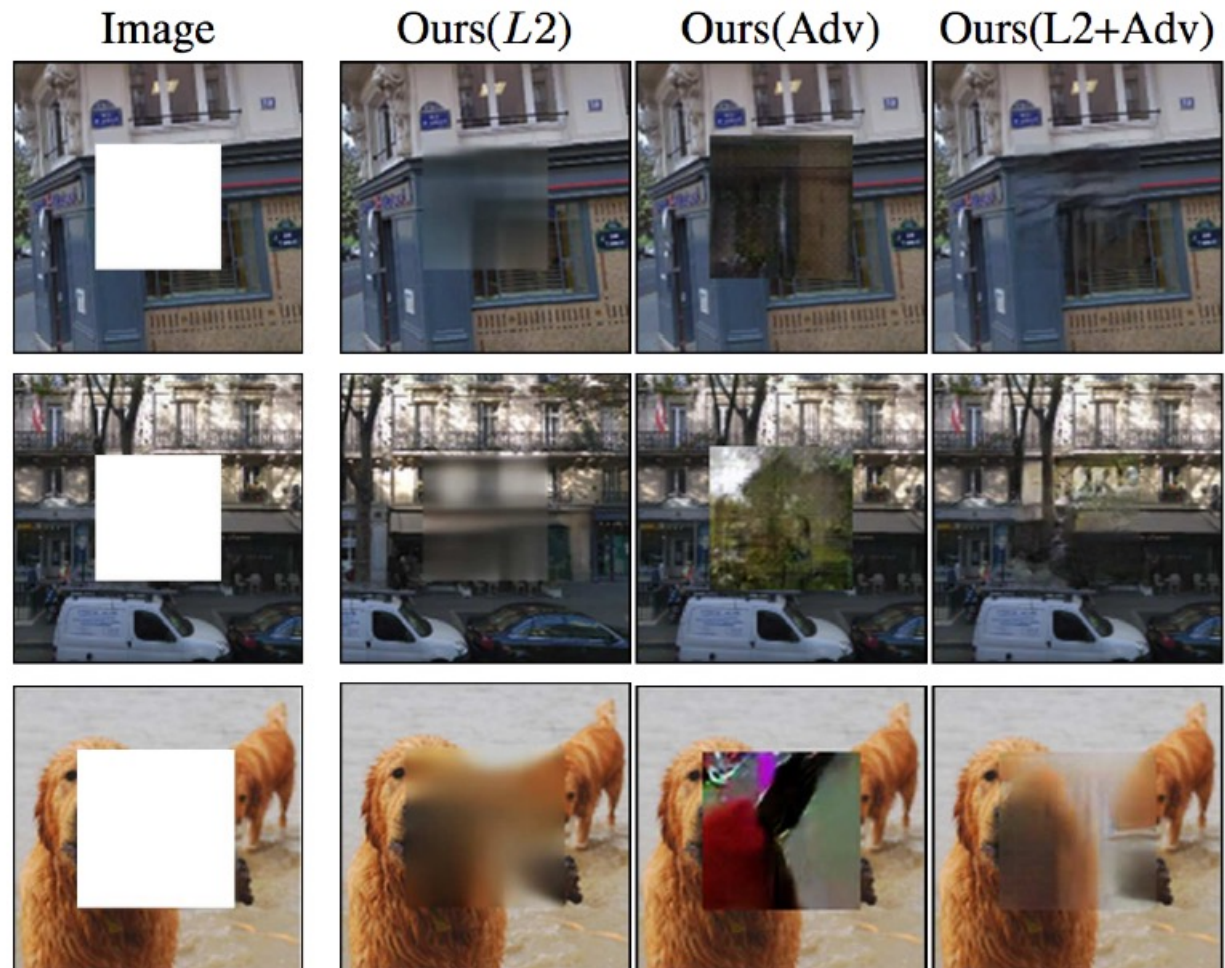
$$\mathcal{L}_{rec}(x) = \left\| \hat{M} \odot \left(x - F((1 - \hat{M}) \odot x) \right) \right\|_2$$

$$\mathcal{L}_{adv} = \max_D \mathbb{E}_{x \in \mathcal{X}} \left[\log(D(x)) + \log \left(1 - D(F((1 - \hat{M}) \odot x)) \right) \right]$$

- Rq: The encoder-decoder is the generator, D is a CNN

Results

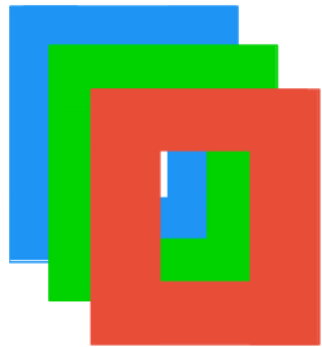
Dataset: StreetView Paris and ImageNet



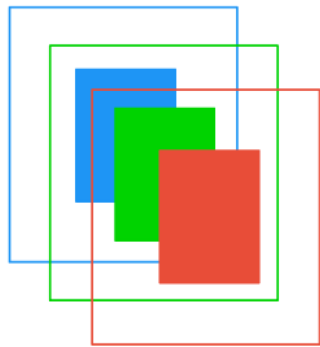
Semantic inpainting - Qualitative results



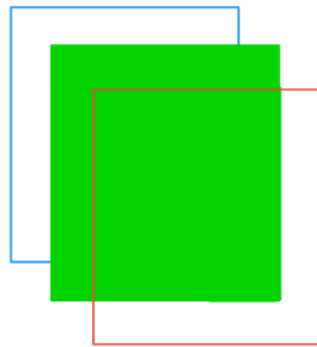
Generalizing inpainting: missing data encoder



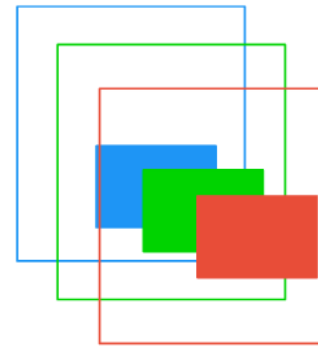
(1) inpainting



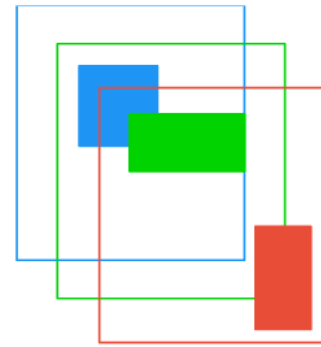
(2) reverse inpainting



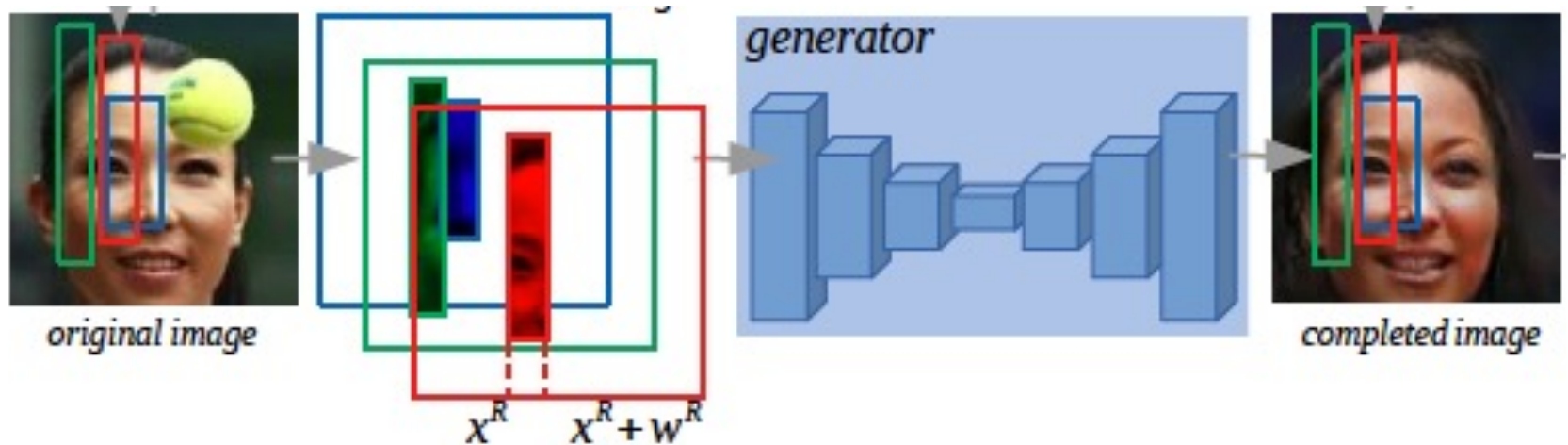
(3) colorization



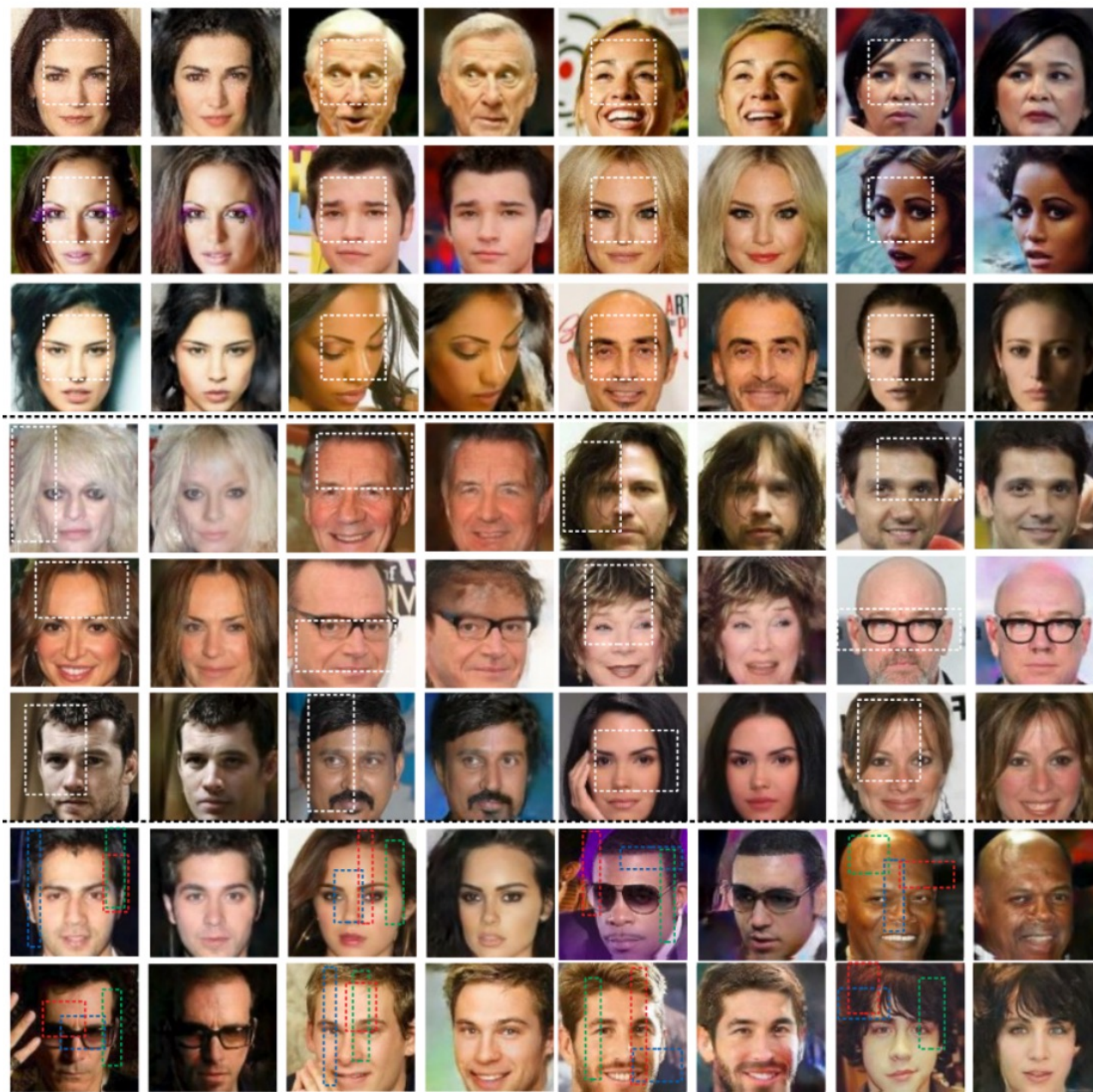
(4) random extrapolation



(5) random extrapolation + colorization



Results



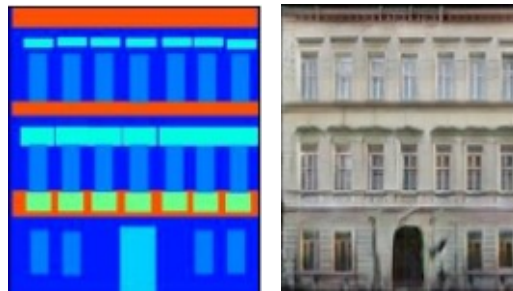
Generative models

Outline

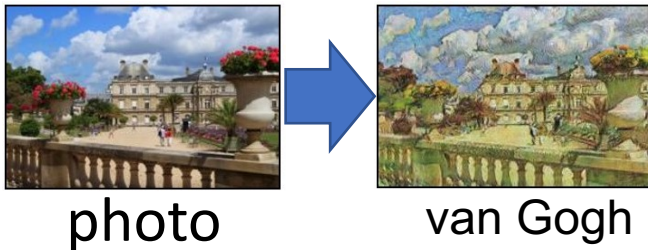
1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
 1. Principle
 2. Text2Image
 3. Image2Image
 4. Inpainting and general missing data encoder
 5. **Learning unpaired Transformation**

Unpaired Transformation

paired data



Transform an object from one domain to another *without paired data*



Domain X



Domain Y



Cycle GAN

<https://arxiv.org/abs/1703.10593>

<https://junyanz.github.io/CycleGAN/>

Domain X



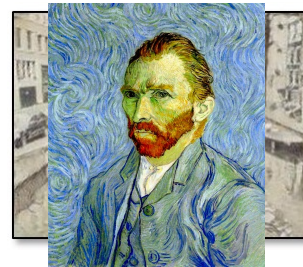
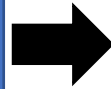
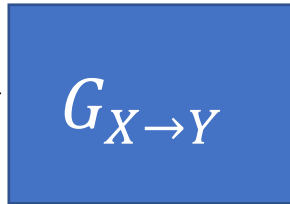
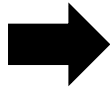
Domain Y



Domain X



ignore input



Become similar
to domain Y

Not what we want



→ scalar



Input image
belongs to
domain Y or not



Domain Y

Cycle GAN

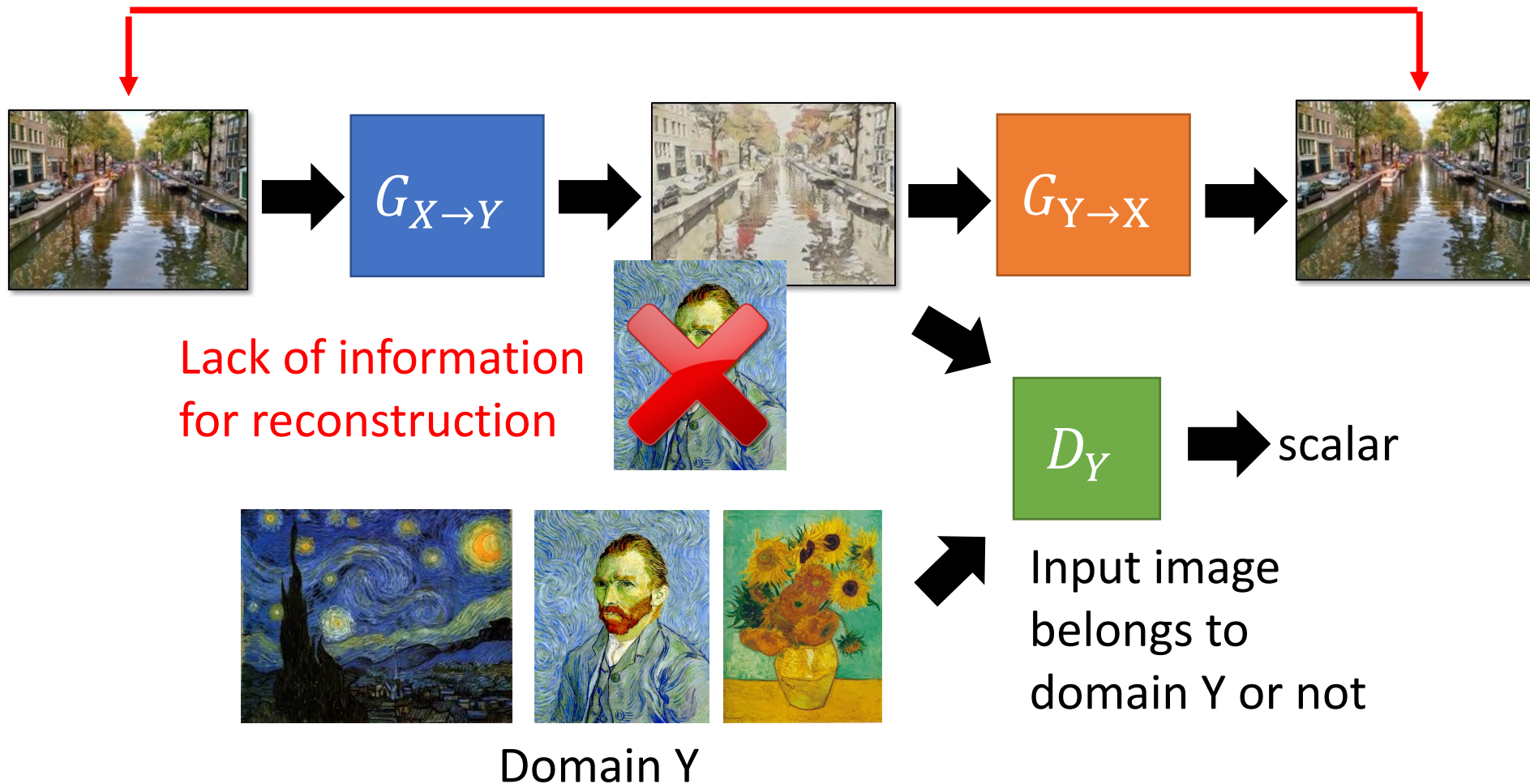
Domain X



Domain Y



as close as possible



Cycle GAN

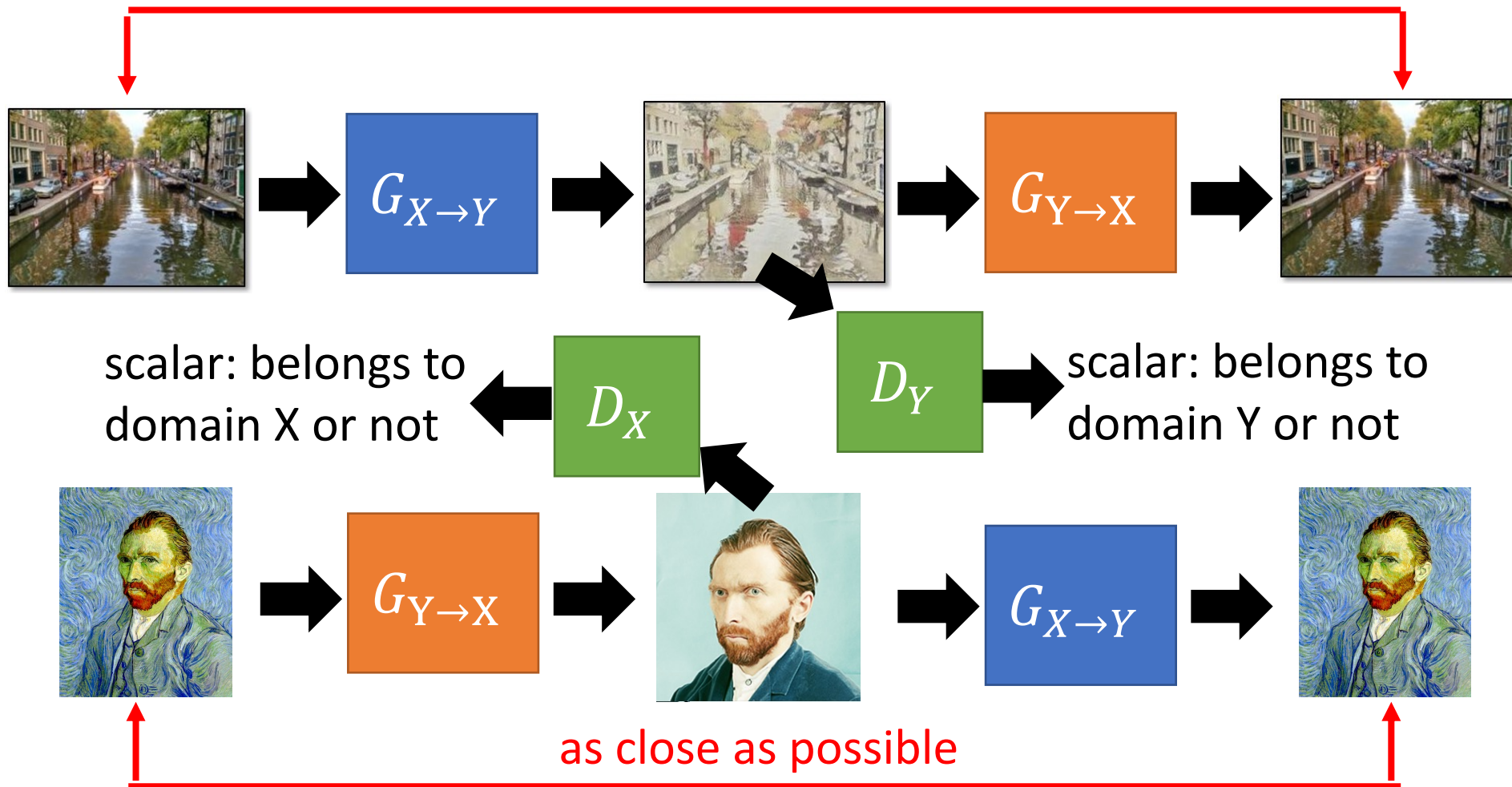
Domain X



Domain Y



as close as possible



Results -- Cycle GAN



photo



van Gogh

Domain X



Domain Y



Monet ↔ Photos



Monet → photo

Zebras ↔ Horses



zebra → horse

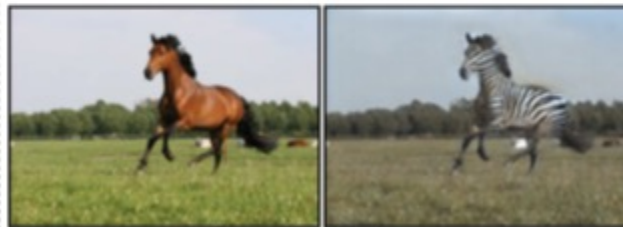
Summer ↔ Winter



summer → winter



photo → Monet



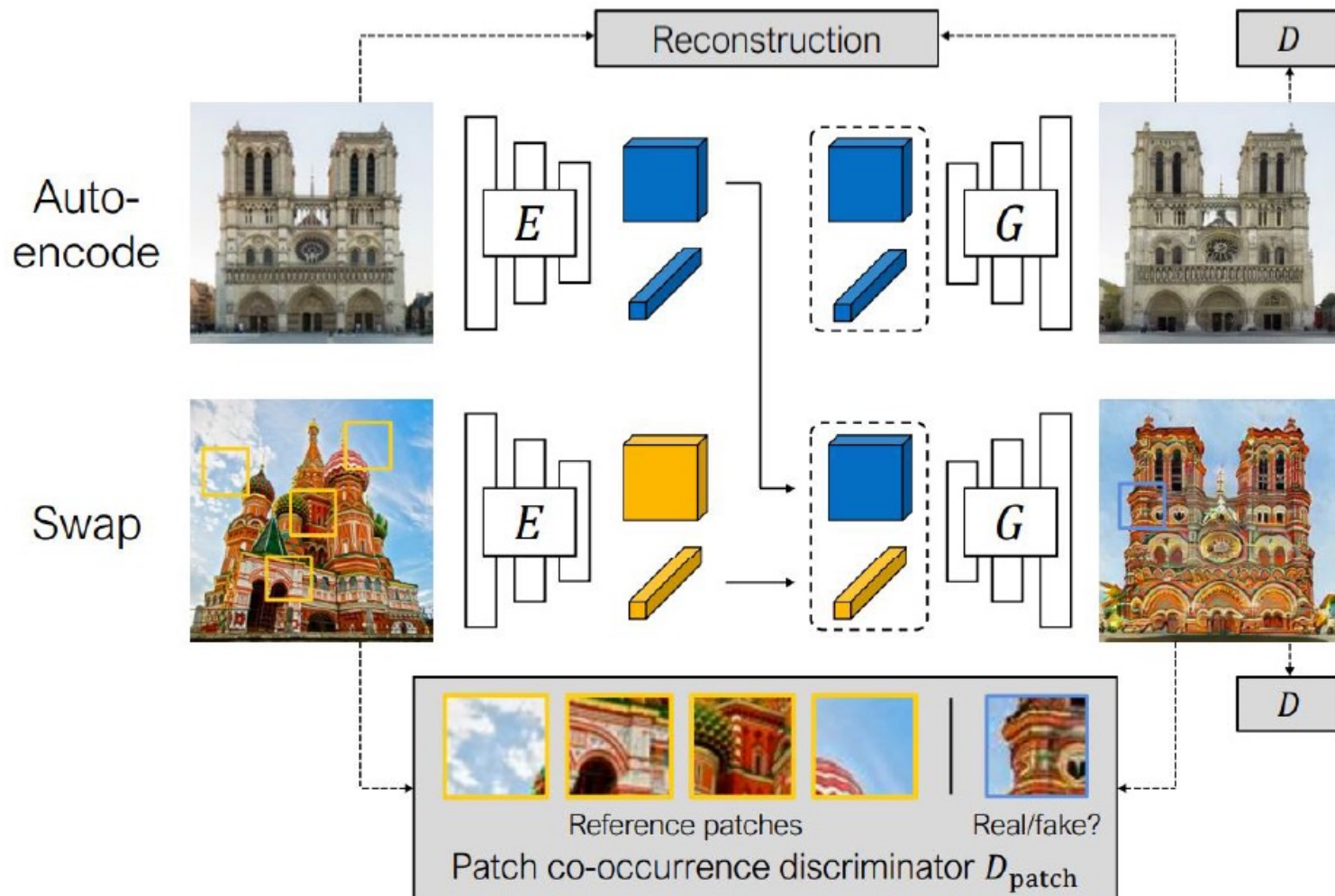
horse → zebra



winter → summer

GANs: works in progress

A lot of things to better understand, to use, adapt, ...

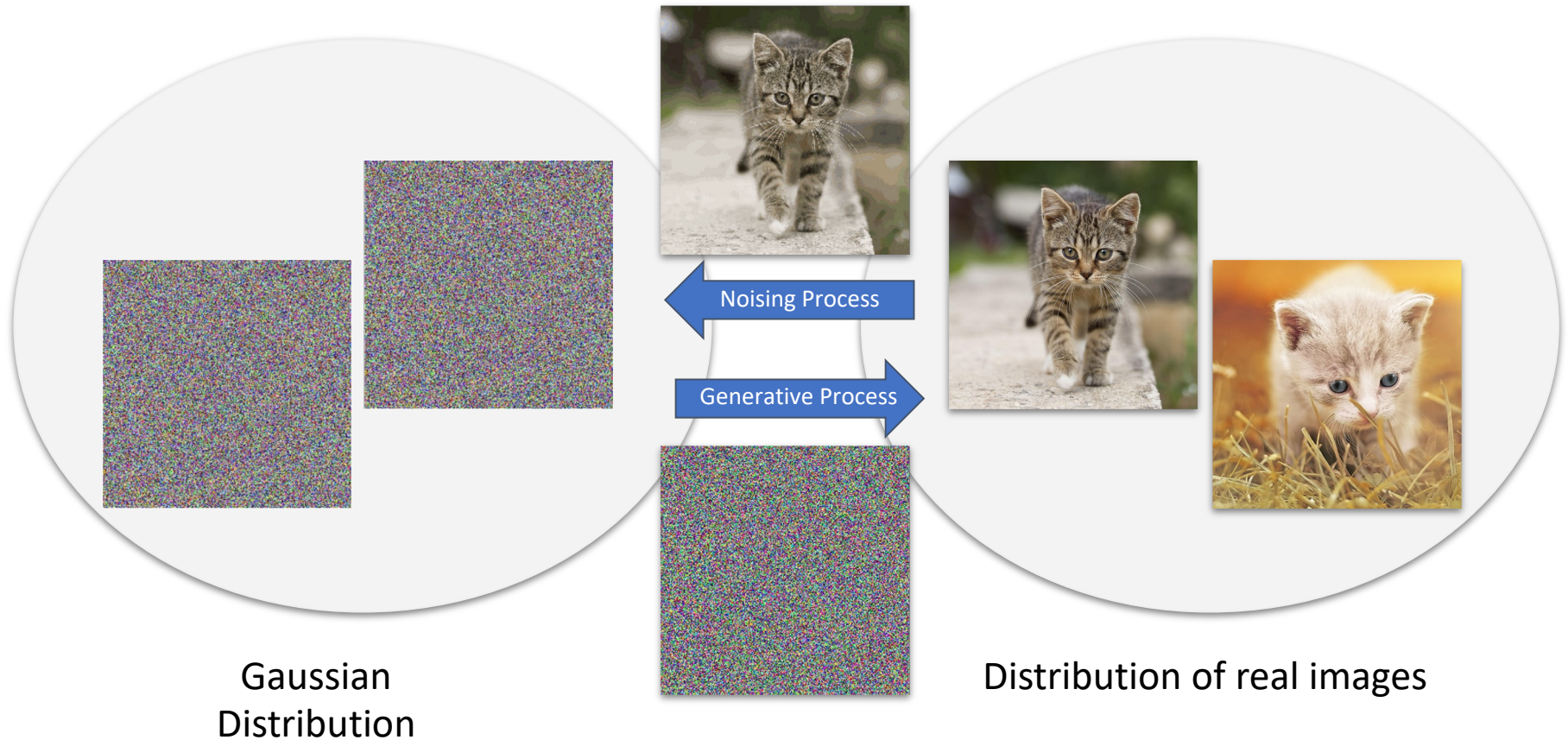


Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
4. Editing
5. Conditional GANs
- 6. Diffusion models**

Generative Modelling with Diffusion models



Generative Modelling with Diffusion models

DDPM: Denoising Diffusion Probabilistic Models
In context with other generative Models:

