



UNIVERSITÀ
DI TRENTO

Semantic Image Synthesis with Spatially-Adaptive Normalization

Taesung Park, Ming-Yu Liu, Ting-Chun Wang, Jun-Yan Zhu

UC Berkeley, NVIDIA, MIT CSAIL

Presented by:

Matteo Destro, Zihadul Azam

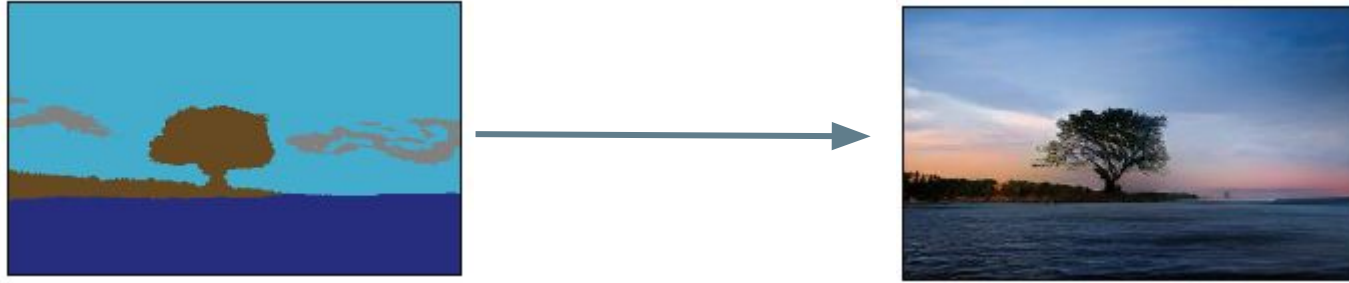
A decorative network graph pattern in the top-left corner, featuring a complex web of interconnected nodes and edges. Some nodes are highlighted with blue circles, and others with solid blue dots. The pattern is composed of light gray lines and dots.

1.

Introduction

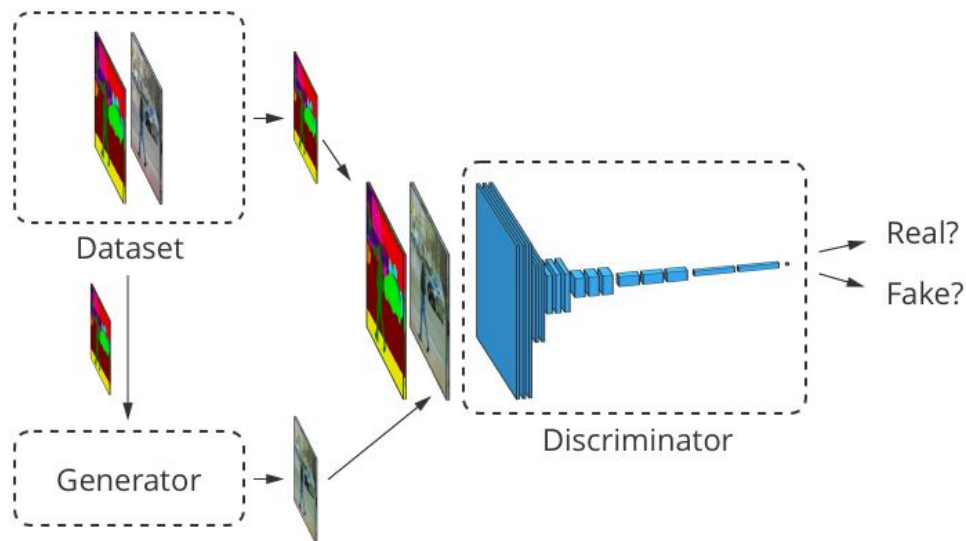
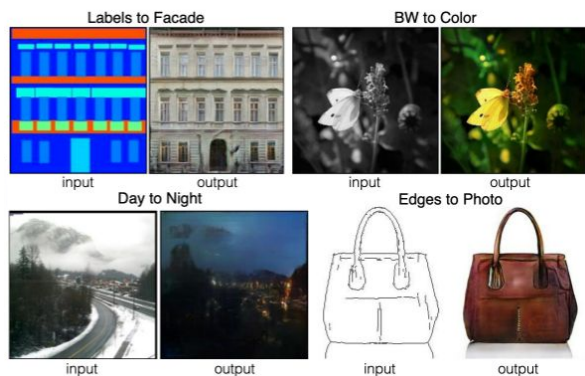
A decorative network graph pattern in the bottom-right corner, featuring a complex web of interconnected nodes and edges. Some nodes are highlighted with blue circles, and others with solid blue dots. The pattern is composed of light gray lines and dots.

What is this paper about?

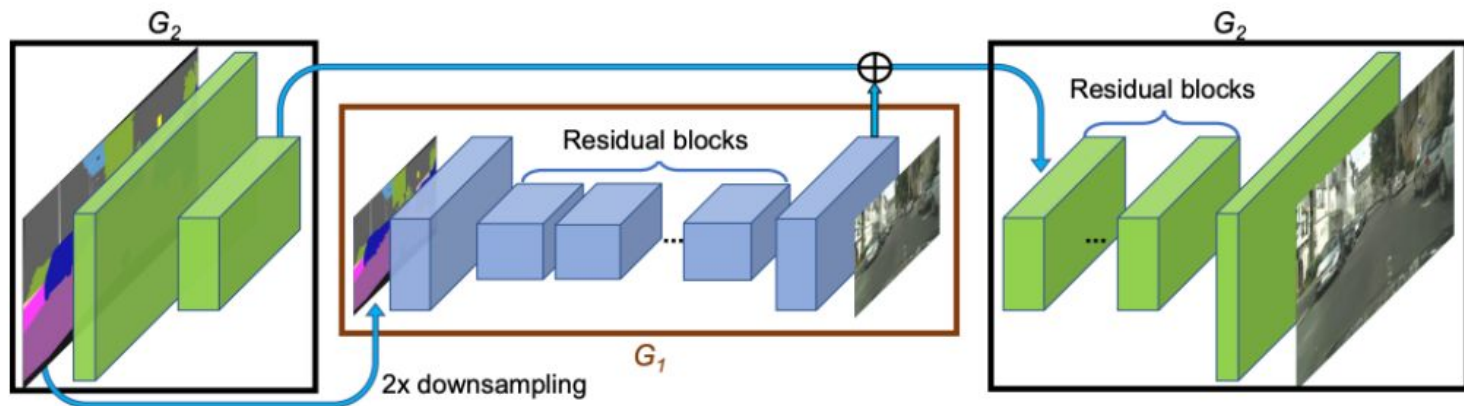


- ◉ Converting a **semantic segmentation** mask to a **photorealistic image**
- ◉ Preserve semantic information
- ◉ Produce images with multiple styles

Image-to-image translation



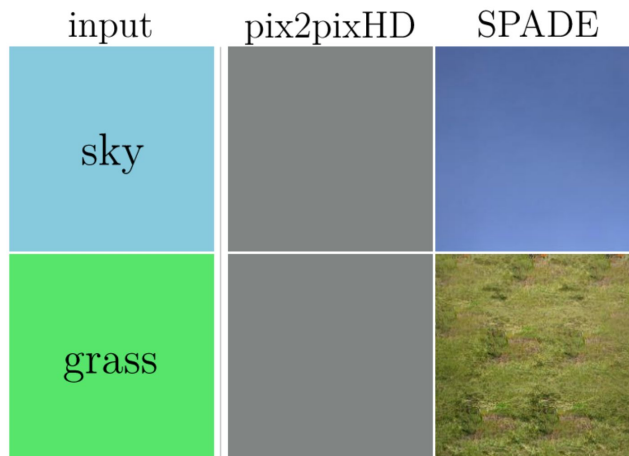
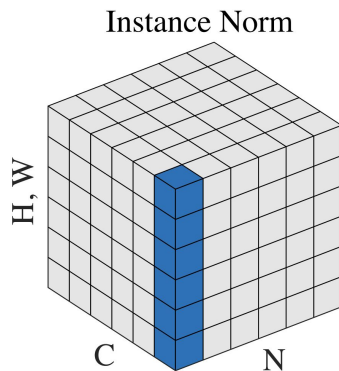
Pix2pixHD: Generator



- Encoder for downscale
- Decoder for upscale

What is wrong with pix2pixHD ?

- ◉ Use **Instance Normalization**
- ◉ **“Wash away information”** problem: in pix2pixHD, instance normalization tends to throw away information from the segmentation map. For single-class images, it produces the same image regardless of the class.
- ◉ Instance normalization **produce 0** as output of the normalization with uniform semantic mask.



Solution:



SPADE

SPatially-Adaptive (DE)normalization

SPADE: normalization

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

Where:

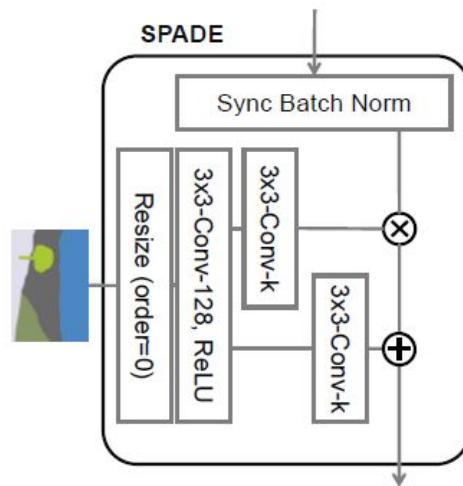
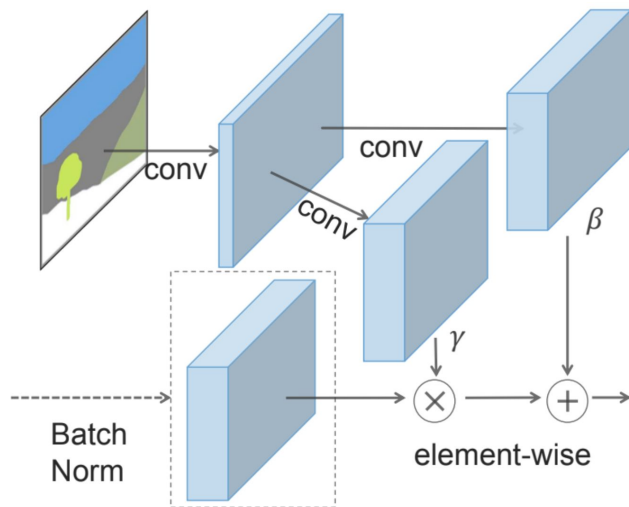
- ◉ $\mathbf{h}_{n;c;y;x}^i$ is the activation at site before normalization
- ◉ μ_c and σ_c are the mean and standard deviation of the activations in channel \mathbf{c}
- ◉ \mathbf{m} is the segmentation mask

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

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

SPADE

- The mask is first projected onto an **embedding space**
- Then it's convolved to produce the modulation parameters γ and β
- The produced γ and β are multiplied and added to the normalized activation element-wise

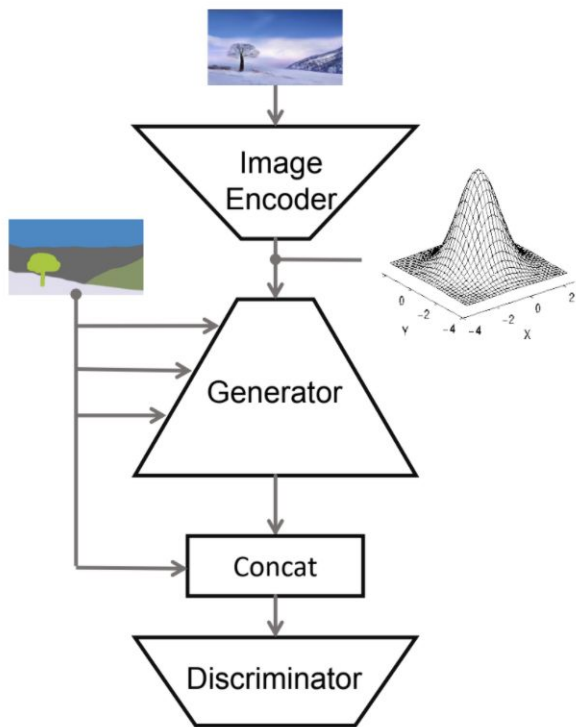


A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. Some nodes are highlighted with blue circles, and others with blue dots. The lines are thin and gray, creating a mesh-like structure.

2. **Architecture**

A decorative network diagram in the bottom-right corner, similar to the one in the top-left. It shows a network of nodes and lines, with some nodes highlighted in blue. The overall style is clean and modern, with a focus on connectivity and structure.

GauGAN architecture



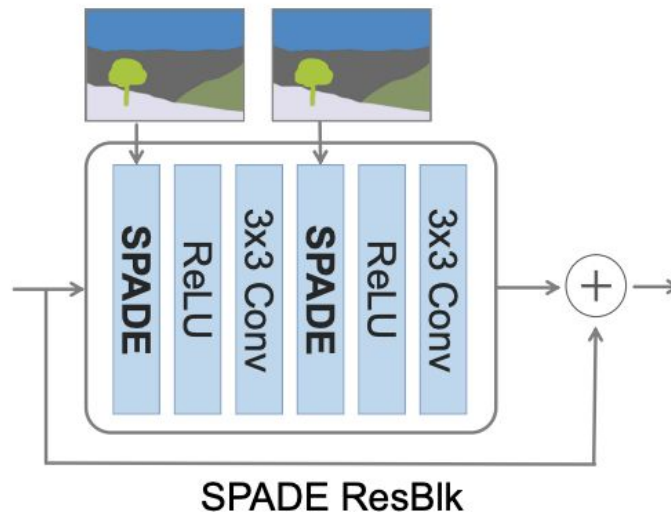
Three main components:

- ◉ **Generator:** generate the final synthetic image.
- ◉ **Discriminator:** discriminate between real and synthetic images. Used only during training.
- ◉ **Encoder:** convert an input image to its latent space representation z . Can be used to transfer the style of an image to the generated one.

Generator: SPADE Residual Block

The generator is composed by a sequence of **SPADE Residual Blocks**:

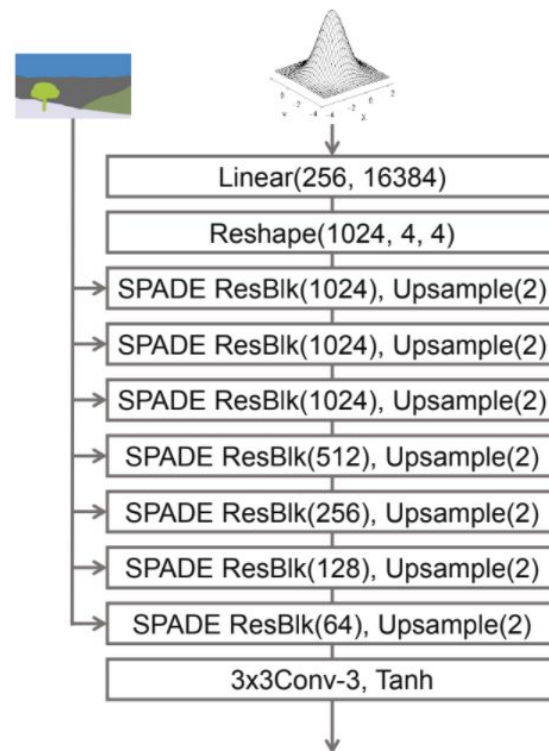
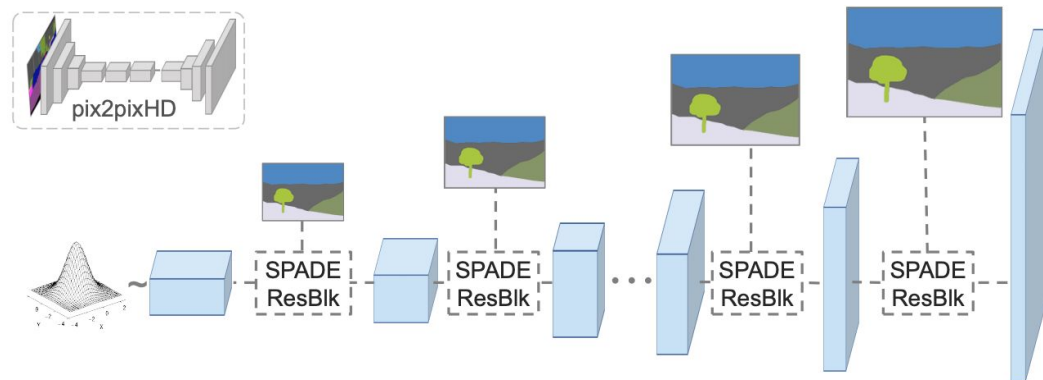
- ◉ **SPADE** normalizations:
 - the semantic mask is downsampled to match the spatial resolution
- ◉ **ReLU** activations
- ◉ **Convolutions:**
 - 3x3 kernel with padding 1
 - k channels in output
 - input and output spatial dimensions are equal
- ◉ **Residual connection**



Generator

The generator is composed by a sequence of **SPADE Residual Blocks** interleaved by **upsampling layers**.

Less parameters: no need for downsampling layers like Pix2pixHD to feed the semantic map to the model (SPADE modulations parameters are enough).

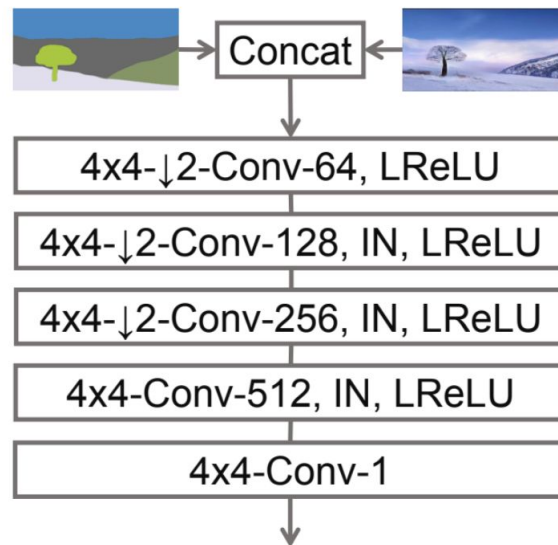


Discriminator: PatchGAN

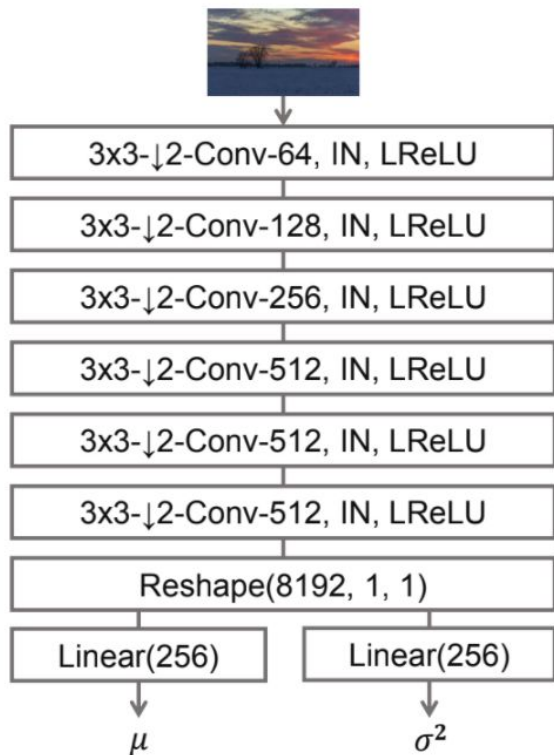
The discriminator:

- ◉ Takes in input the concatenation of the semantic map (one-hot encoded) and the image
- ◉ Is composed by a series of **convolutional layers**:
 - 4x4 kernel with stride 2
 - Leaky ReLU activations
 - Instance normalization
- ◉ A final **convolutional layer** computes the final prediction

Two discriminators are used, working at different scales.



Encoder: VAE



The encoder is composed by:

- Series of **convolutional layers**:
 - 3x3 kernel with stride 2
 - Leaky ReLU activations
 - Instance normalization
- Two **linear layers** to compute:
 - mean vector (μ)
 - variance vector (σ^2)

The final latent space representation \mathbf{z} is computed via a reparameterization trick:

$$z = \mu + \mathcal{N}(0, 1) \cdot e^{\frac{1}{2}\sigma^2}$$

Loss functions

- ◉ **Adversarial loss:** uses the **Hinge loss** instead of the Least Squares loss of Pix2pixHD:

$$L_D = -\mathbb{E} [\min(0, -1 + D(x, y))] - \mathbb{E} [\min(0, -1 - D(G(z, \mathbf{m}), y))]$$

$$L_G = -\mathbb{E} [D(G(z, \mathbf{m}), y)]$$

- ◉ **Feature matching loss:**

$$L_{FM} = \mathbb{E} \left[\sum_{i=1}^L \left\| D_k^{(i)}(x) - D_k^{(i)}(G(z, \mathbf{m})) \right\|_1 \right]$$

- ◉ **Perceptual loss:** similar concept to Feature loss, but uses activations of a VGG-19 pretrained model

$$L_{VGG} = \mathbb{E} \left[\sum_{i=1}^L \left\| VGG^{(i)}(x) - VGG^{(i)}(G(z, \mathbf{m})) \right\|_1 \right]$$

- ◉ **Encoder loss** (if used):

$$L_{KLD} = \mathcal{D}_{KL}(q(z|x) || p(z))$$

A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. Some nodes are highlighted with blue circles, and others with blue dots.

3. Results

A decorative network diagram in the bottom-right corner, featuring a complex web of interconnected nodes and lines. Some nodes are highlighted with blue circles, and others with blue dots.

Datasets

- ◉ **COCO-Stuff:**
 - 118'000 training and 5'000 validation images
 - 182 semantic classes
- ◉ **ADE20K (-outdoor):**
 - 20'120 training and 2'000 validation images
 - 150 semantic classes
- ◉ **Cityscapes:**
 - 3'000 training and 500 validation images
 - 19 semantic classes
- ◉ **Flickr Landscapes:**
 - 41'000 training and 1'000 validation images
 - Semantic maps computed using a pre-trained DeepLabV2 model

Evaluation metrics

Compare:

- ⦿ **Human evaluation**
- ⦿ **Fréchet Inception Distance** (FID)
- ⦿ **Mean Intersection Over Union** (mIoU)
- ⦿ **Pixel accuracy** (accu)

Idea: run a **semantic segmentation** model on the generated images and compare the resulting segmentation map with the original input.

Use state-of-the-art models for each dataset: **DeepLabV2** for COCO-Stuff, **UperNet101** for ADE20K, **DRN-D-105** for Cityscapes

Experiment results

Results compared with state-of-the-art semantic image synthesis models:

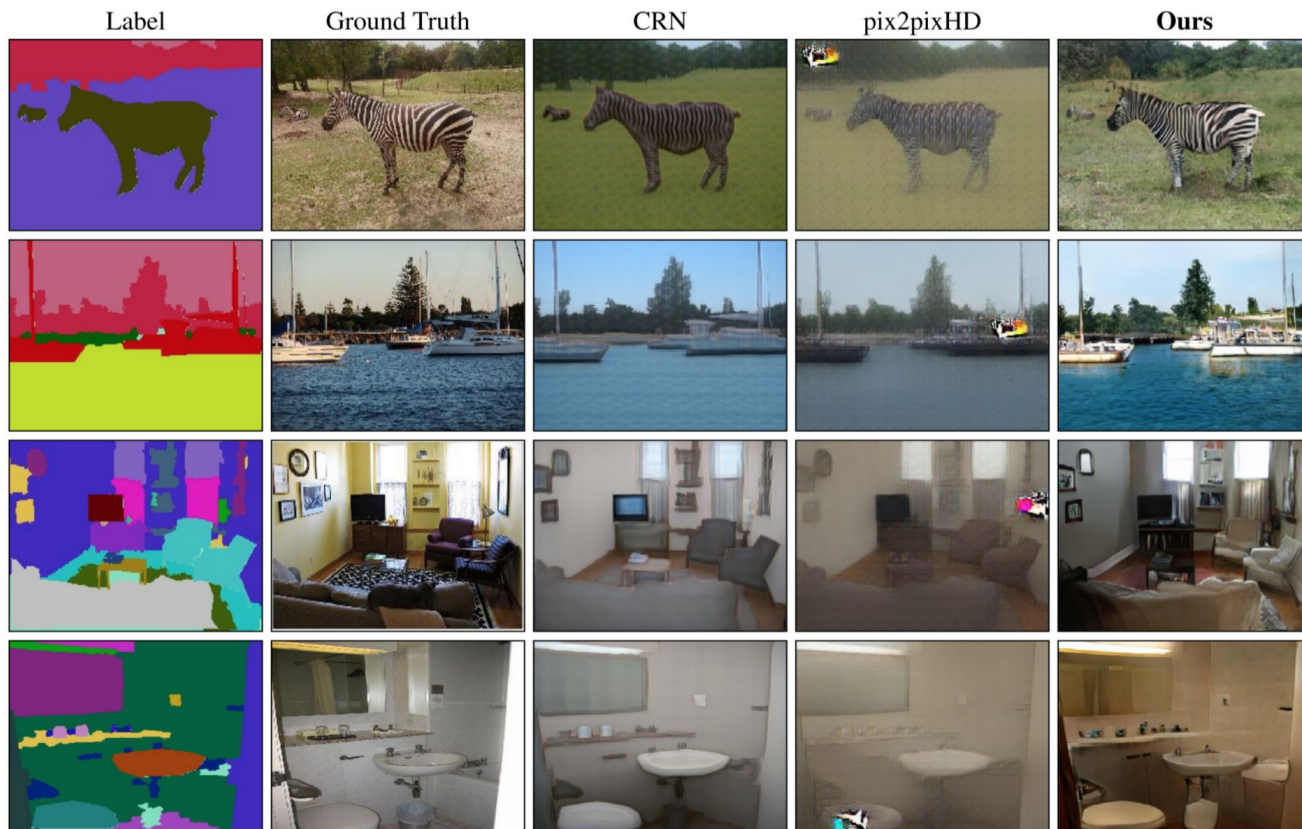
- ◉ **Cascaded Refinement Network**^[1] (CRN)
Uses a deep network that repeatedly refines the output from low to high resolution
- ◉ **Semi-parametric Image Synthesis Method**^[2] (SIMS)
Semi-parametric approach that composites real segments from a training set and refines the boundaries.
- ◉ **Pix2pixHD**
State-of-the-art GAN-based conditional image synthesis framework

Method	COCO-Stuff			ADE20K			ADE20K-outdoor			Cityscapes		
	mIoU	accu	FID	mIoU	accu	FID	mIoU	accu	FID	mIoU	accu	FID
CRN [6]	23.7	40.4	70.4	22.4	68.8	73.3	16.5	68.6	99.0	52.4	77.1	104.7
SIMS [43]	N/A	N/A	N/A	N/A	N/A	N/A	13.1	74.7	67.7	47.2	75.5	49.7
pix2pixHD [48]	14.6	45.8	111.5	20.3	69.2	81.8	17.4	71.6	97.8	58.3	81.4	95.0
Ours	37.4	67.9	22.6	38.5	79.9	33.9	30.8	82.9	63.3	62.3	81.9	71.8

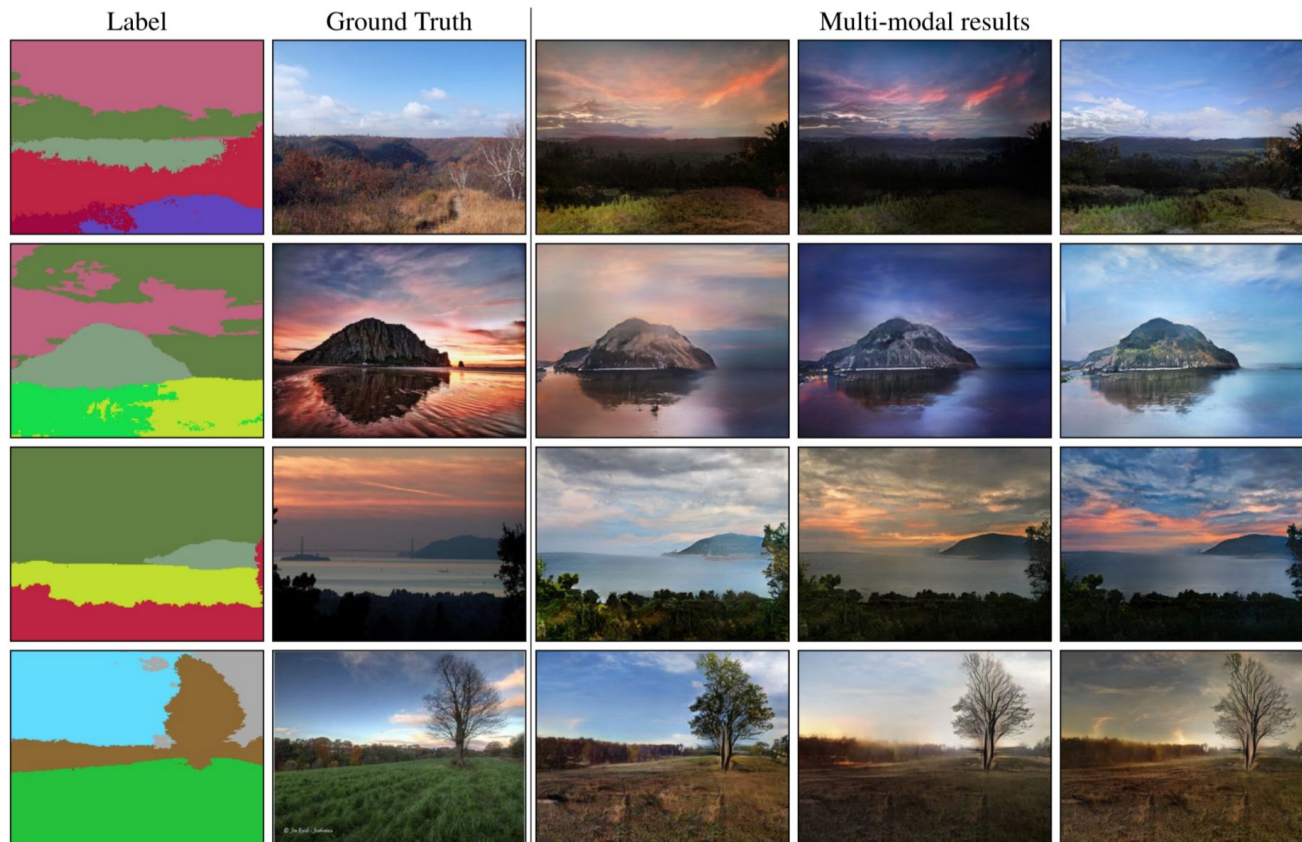
[1] Q. Chen and V. Koltun, “Photographic image synthesis with cascaded refinement networks”, ICCV 2017

[2] X. Qi, Q. Chen, J. Jia, and V. Koltun, “Semi-parametric image synthesis”, CVPR 2018

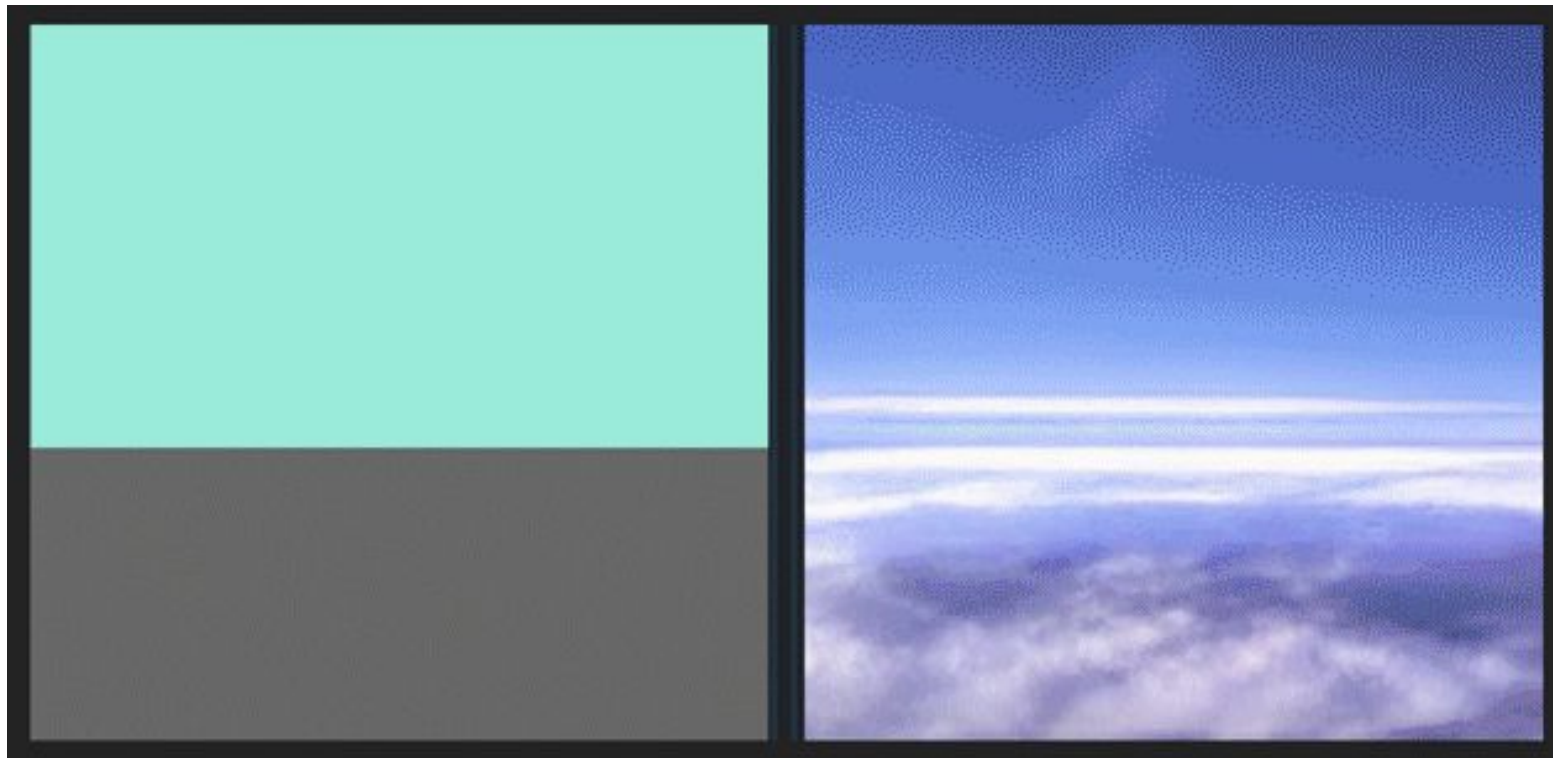
Experiment results: Image Synthesis



Experiment results: Multimodal Synthesis

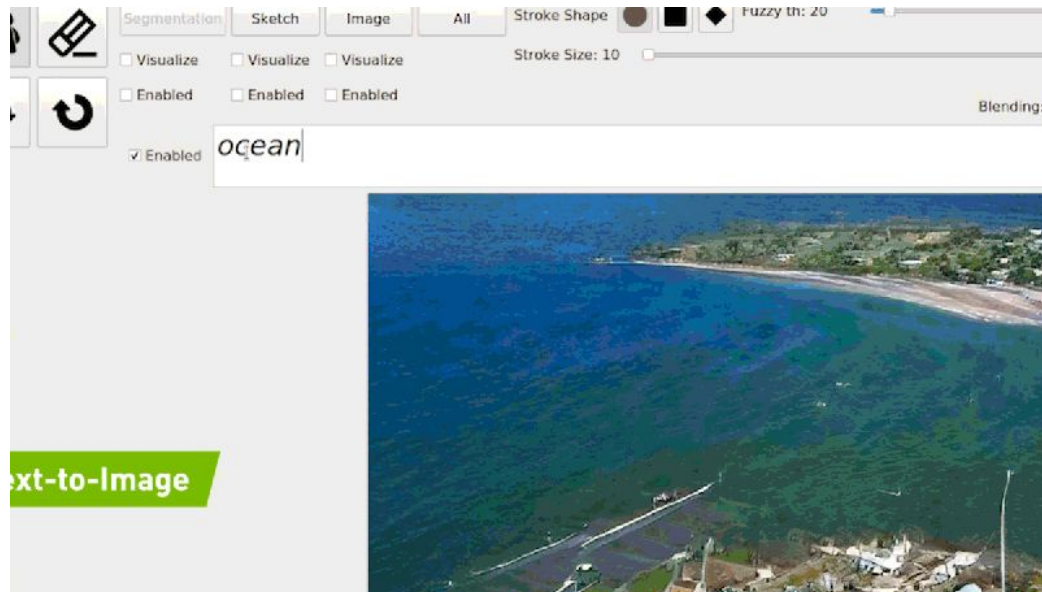


Demo: Flickr Landscapes dataset



Improvements from other works

- ◉ **CLADE**: Class-Adaptive (DE)normalization
- ◉ **SEAN**: Image Synthesis with Semantic Region-Adaptive Normalization
- ◉ **GauGAN2**



A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and lines. Some nodes are highlighted with blue circles, and others with blue dots. The lines are thin and grey, creating a mesh-like structure.

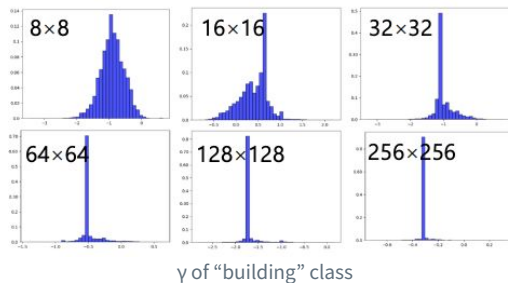
4.

Improvements From Other Works

A decorative network diagram in the bottom-right corner, similar to the one in the top-left, showing a web of interconnected nodes and lines with some nodes highlighted in blue.

CLADE: Class-Adaptive (DE)normalization

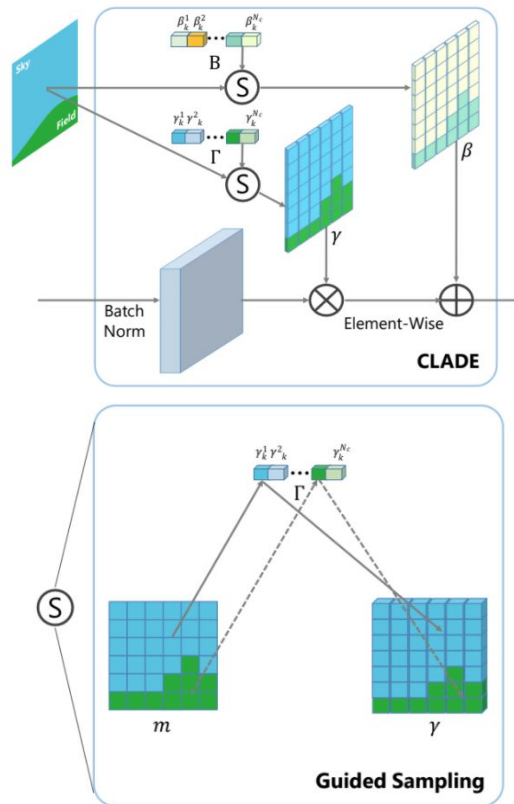
They showed that the spatial information is not very beneficial in SPADE, its effectiveness is given by its semantic awareness.



Idea:

- Learn a single γ and β for each semantic class, instead of modulating them through two a shallow CNN
- Create spatial-aware modulation parameters through **guided sampling**

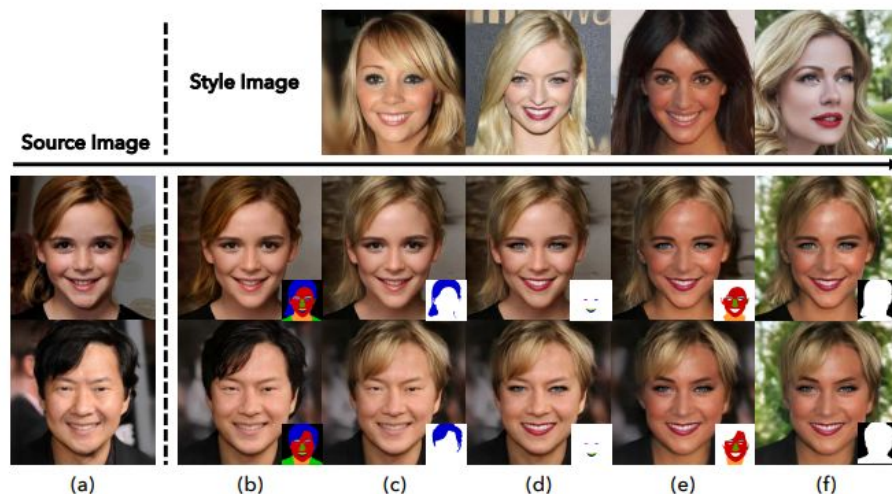
The modulation convolutional network of SPADE is removed



SEAN: Image Synthesis with Semantic Region-Adaptive Normalization

Face image editing controlled via style images and segmentation masks:

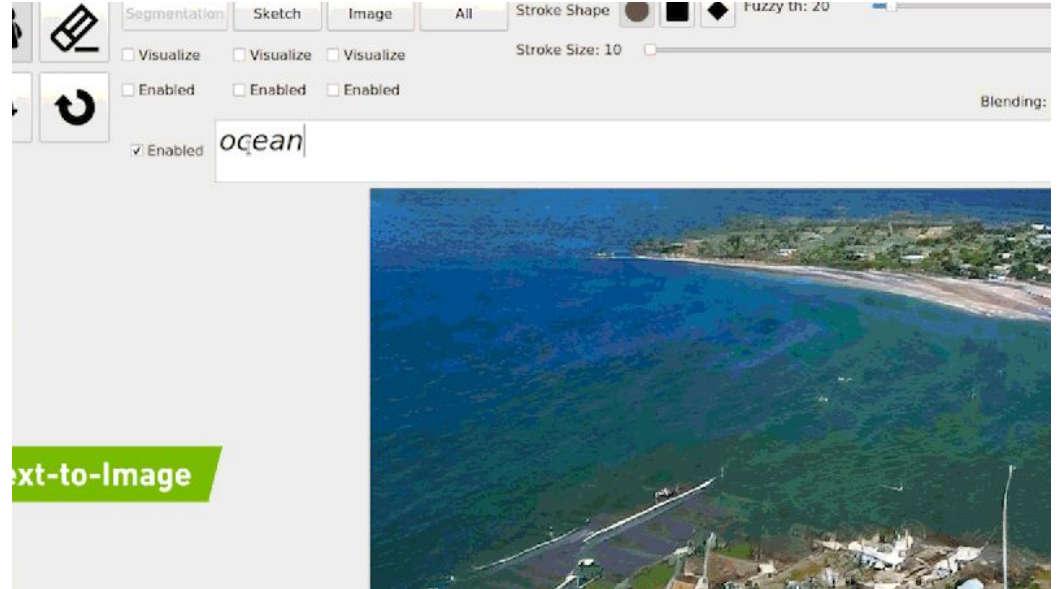
- ◉ This network architecture can control the **style** of each semantic region individually
- ◉ Here we can specify one style reference image per region



GauGAN2

GauGAN2 is designed to create photorealistic art with a **mix of words** and **drawings**:

- ◉ Trained on 10 million images
- ◉ Can translate natural language descriptions into landscape images



A decorative network diagram in the top-left corner, featuring a complex web of interconnected nodes and edges. Some nodes are highlighted with blue circles, and others with blue dots. The network is composed of various shapes, including circles and squares, connected by thin lines.

5. Unpaired SPADE

A decorative network diagram in the bottom-right corner, similar to the one in the top-left. It shows a network of nodes and edges, with some nodes highlighted in blue. The network is composed of various shapes, including circles and squares, connected by thin lines.

Unpaired Semantic Image Synthesis

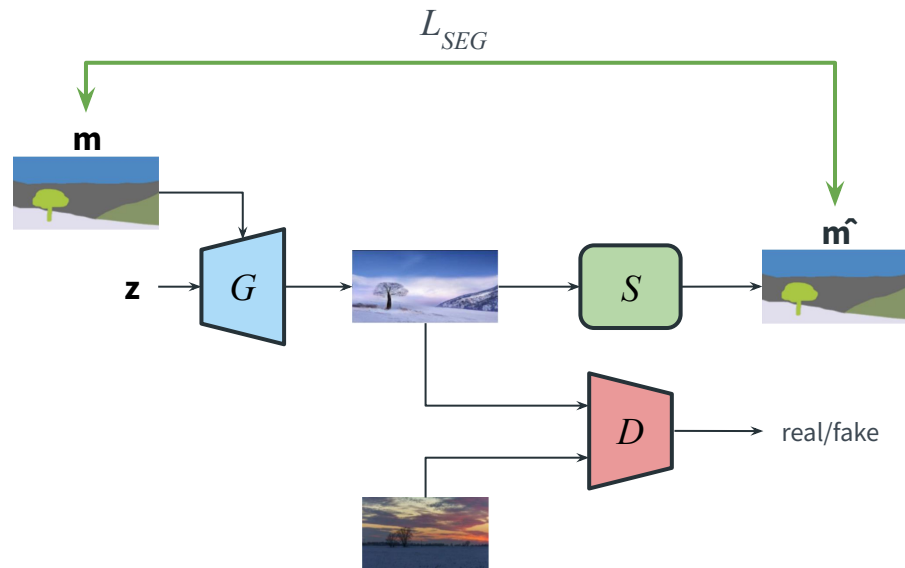
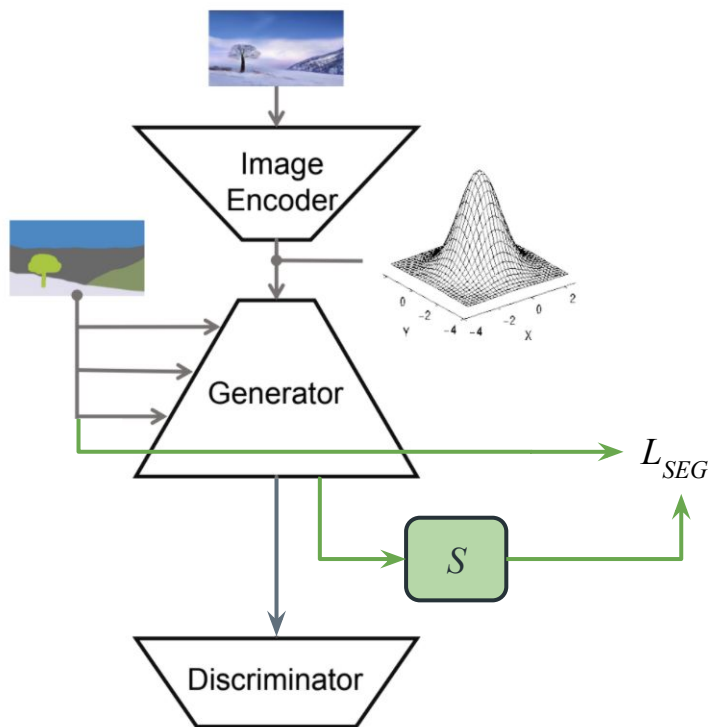
What if we don't have paired samples?

Idea: exploit a **semantic segmentation model** to train the model in an unpaired setting, introducing a reconstruction loss w.r.t the semantic map input.

We can use an approach similar to CycleGAN (unpaired image-to-image translation):

- ◉ Add a **segmentation model** in the loop (e.g. U-Net, DeepLabV3, ...)
- ◉ Add a **segmentation loss** based on the reconstructed segmentation

Unpaired Semantic Image Synthesis



Unpaired Semantic Image Synthesis

Introduce a class-balanced semantic segmentation loss in the generator:

- ◉ **Class-balanced** loss based on **effective number of samples**^[1]: $\alpha_c = \frac{1 - \beta}{1 - \beta^{n_c}}$

- ◉ **Focal loss**^[2]:
$$L_{SEG} = - \sum_{c,i,j}^{C \times H \times W} \alpha_c \cdot \mathbf{m}_{c,i,j} \cdot (1 - \hat{\mathbf{m}}_{c,i,j})^\gamma \log(\hat{\mathbf{m}}_{c,i,j})$$

Final generator loss: $L_G = -\mathbb{E} [D(G(z, \mathbf{m}))] + \lambda \mathbb{E} [L_{SEG}]$

[1] Y. Cui, M. Jia, T. Lin, Y. Song and S. Belongie, “Class-Balanced Loss Based on Effective Number of Samples”, CVPR 2019

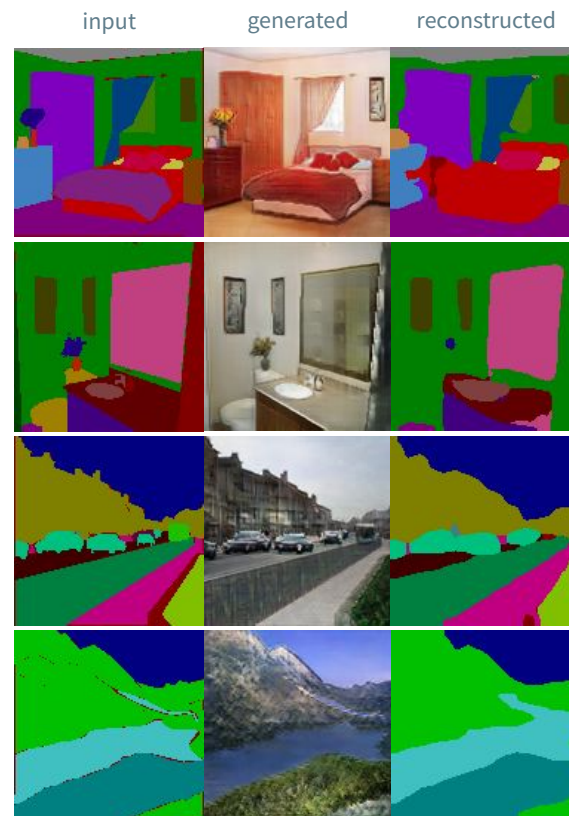
[2] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, “Focal loss for dense object detection”, PAMI 2018

Unpaired Semantic Image Synthesis: Results

Performances obtained over the ADE20K dataset. To reduce training time, we used:

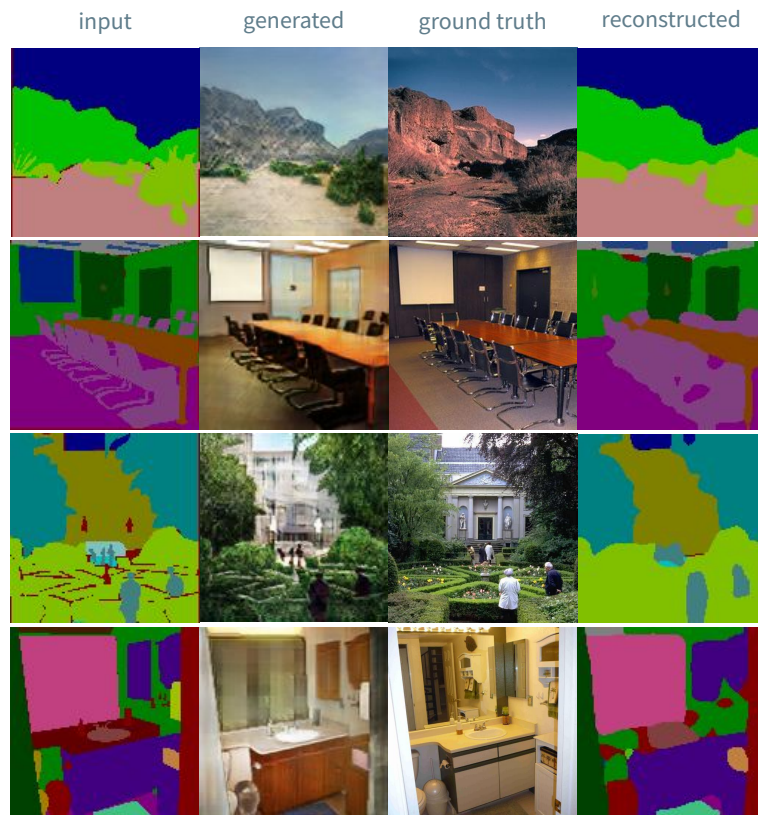
- ◉ **Smaller model:** 14.1M parameters instead of 96.5M
- ◉ **Smaller input images:** 64x64 instead of 256x256
- ◉ **Less epochs:** 50 instead of 200
- ◉ UPerNet18^[1] as **segmenter** network

Method	mIoU	Accuracy	FID
SPADE (baseline)	27.03	58.35	87.13
SPADE + pretrained UPerNet	29.74	59.96	80.10
SPADE + untrained UPerNet	29.54	59.73	83.67



[1] T. Xiao, Y. Liu, B. Zhou, Y. Jiang and J. Sun, “Unified Perceptual Parsing for Scene Understanding”, ECCV 2018

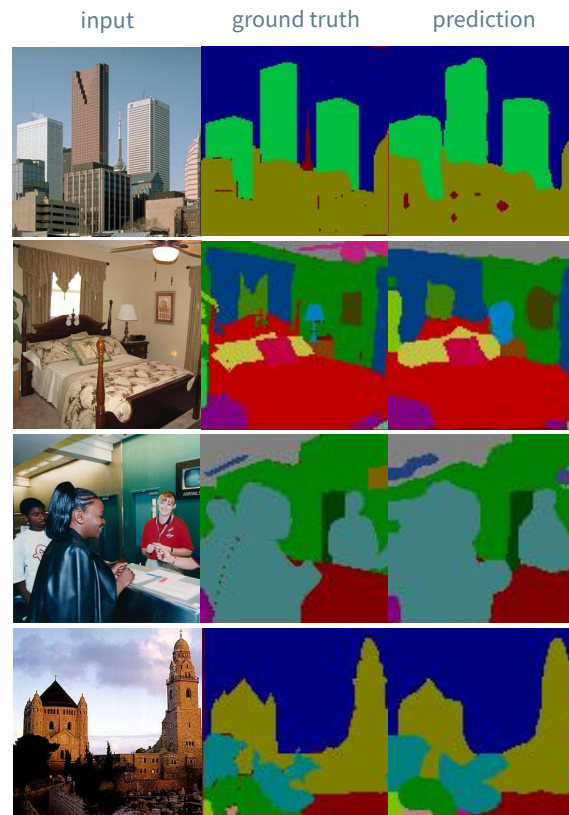
Unpaired Semantic Image Synthesis: Results



Unpaired Semantic Image Synthesis: Results

As side-effect, we train a **semantic segmentation** network in an **unsupervised** way.

Method	mIoU	Accuracy
Pretrained UPerNet	38.18	78.64
UPerNet trained with SPADE	26.76	57.27

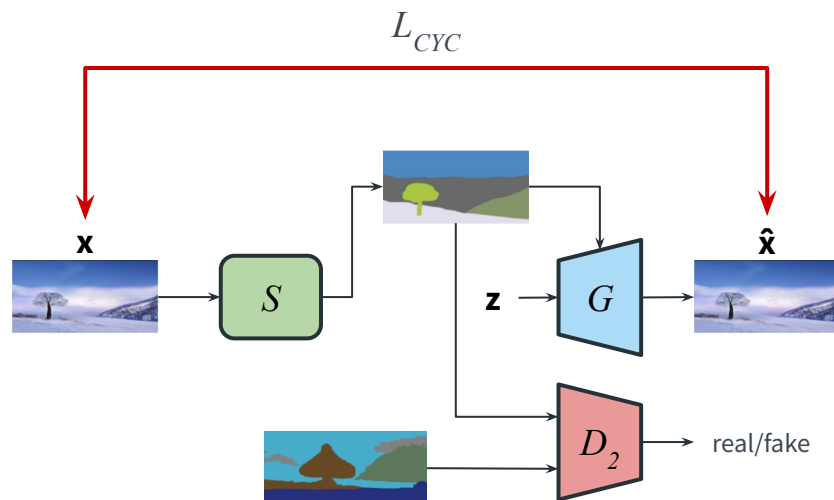
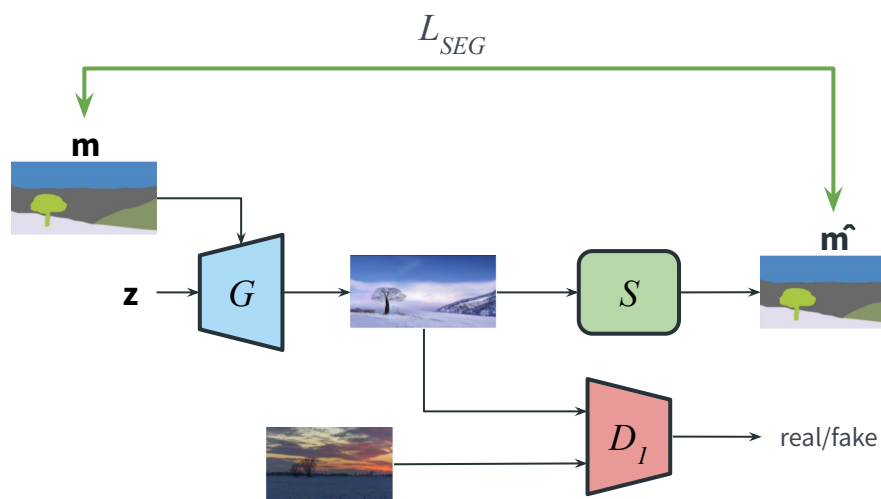


Unpaired Semantic Image Synthesis: Future works

We may also add a second discriminator that detect real/fake semantic maps in order to introduce a **cycle consistency loss**:

$$L_G = -\mathbb{E} [D_1(G(z, \mathbf{m}))] - \mathbb{E} [D_2(S(x))] + \lambda_1 \mathbb{E} [L_{SEG}] + \lambda_2 \mathbb{E} [L_{CYC}]$$

$$L_{CYC} = ||G(z, S(x)) - x||_1$$



Unpaired Semantic Image Synthesis: Recap

- ◉ No need for **paired** label-image samples (e.g. we could create **synthetic input labels**)
- ◉ Secondary result: **unsupervised semantic segmentation** model
- ◉ The **rescaled loss** give more importance to classes with small objects
- ◉ We can exploit pretrained semantic segmentation models, if available
- ◉ If a set of paired samples is available (even small), it can be exploited to further train the semantic segmentation network



Thank you for the attention