

# Semantic Image Synthesis with Spatially-Adaptive Normalization

Taesung Park<sup>1,2</sup>, Ming-Yu Liu<sup>2</sup>, Ting-Chun Wang<sup>2</sup>, Jun-Yan Zhu<sup>2,3</sup>

<sup>1</sup>UC Berkeley, <sup>2</sup>NVIDIA , <sup>3</sup>MIT CSAIL

Alexander Koenig, Li Nguyen

Workshop in Machine Learning Applications for Computer Graphics  
Blavatnik School of Computer Science, Tel Aviv University

April 1, 2020

# Outline

1 Introduction

2 Related Work

3 Semantic Image Synthesis

4 Experiments

5 Conclusion

6 References

# Outline

1 Introduction

2 Related Work

3 Semantic Image Synthesis

4 Experiments

5 Conclusion

6 References

# Conditional Image Synthesis

- **Conditional Image Synthesis:** method to generate photorealistic images based on certain input (e.g. labels, text, ...)
- Semantic Image Synthesis: method to generate photorealistic images based on **semantic segmentation mask**

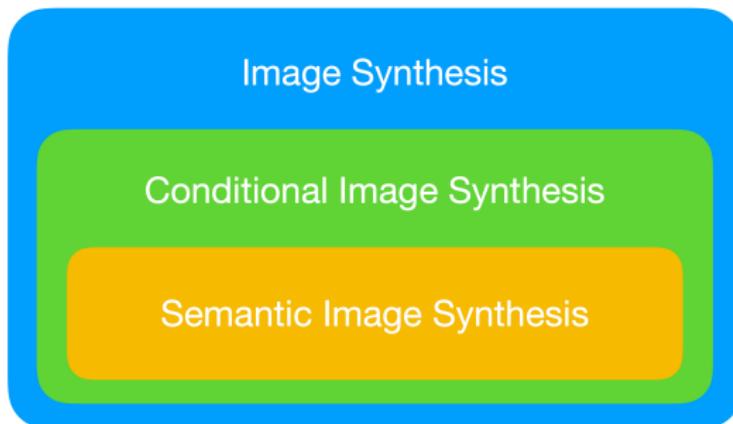


Figure: Euler diagram of image synthesis methods

# What is a Semantic Segmentation Mask?



Figure: Ground truth [4]



Figure: Segmentation mask [4]

- Semantic segmentation: clustering image pixels together which belong to the same object class [6]
- Goal: turn segmentation mask into a photorealistic image
- Application of Semantic Image Synthesis: content generation and image editing

# Outline

1 Introduction

2 Related Work

3 Semantic Image Synthesis

4 Experiments

5 Conclusion

6 References

# Related Work

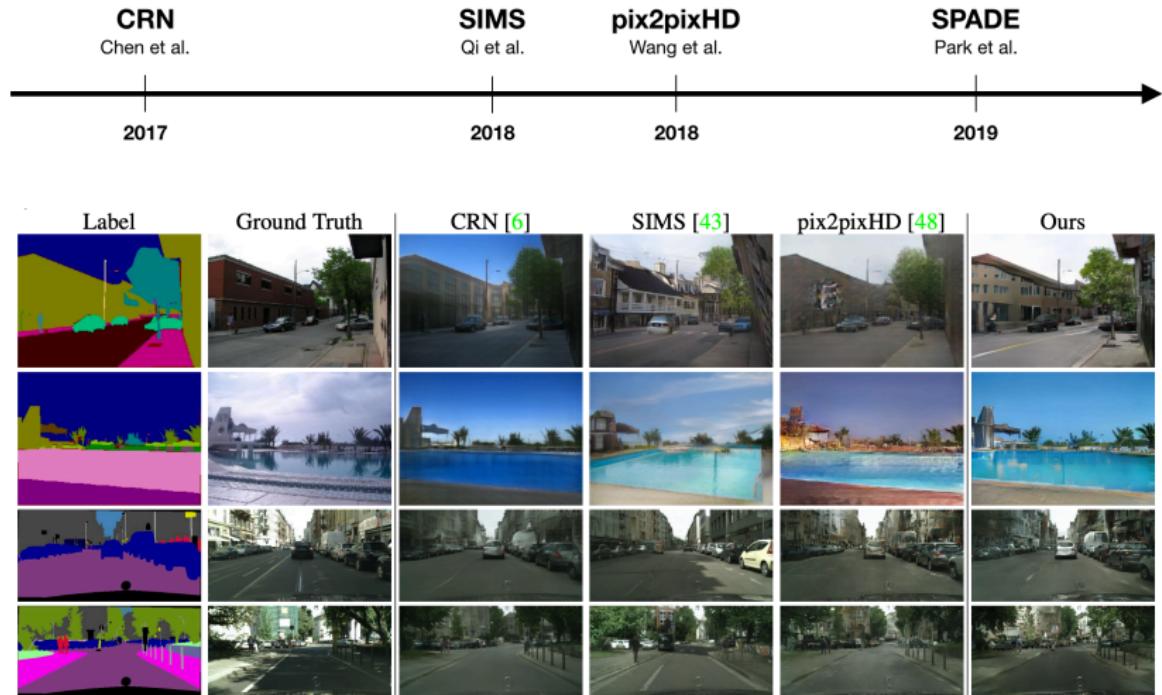


Figure: Visual Comparison of Park et al. to Related Works

## Related Work: Cascaded Refinement Network (CRN)

- The architecture consists of a **cascade** of refinement modules which operate at **different resolutions** each
- Each layer is followed by convolutions, normalization and a non-linearity [1]

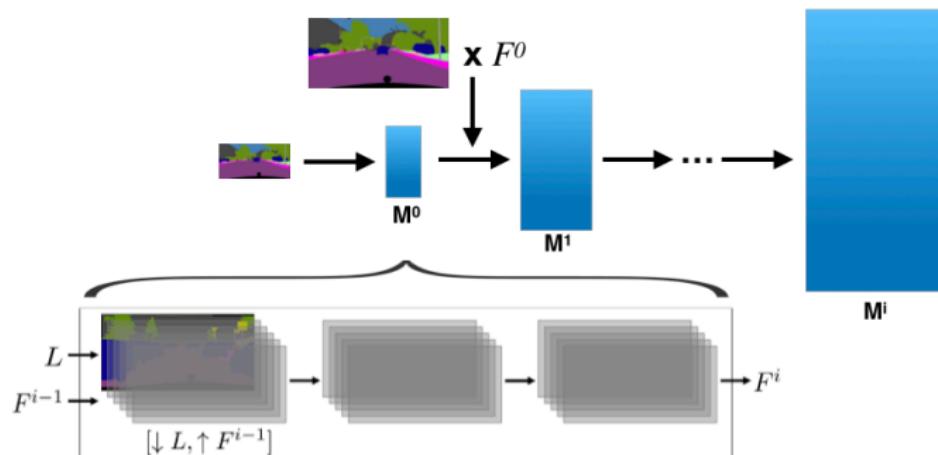
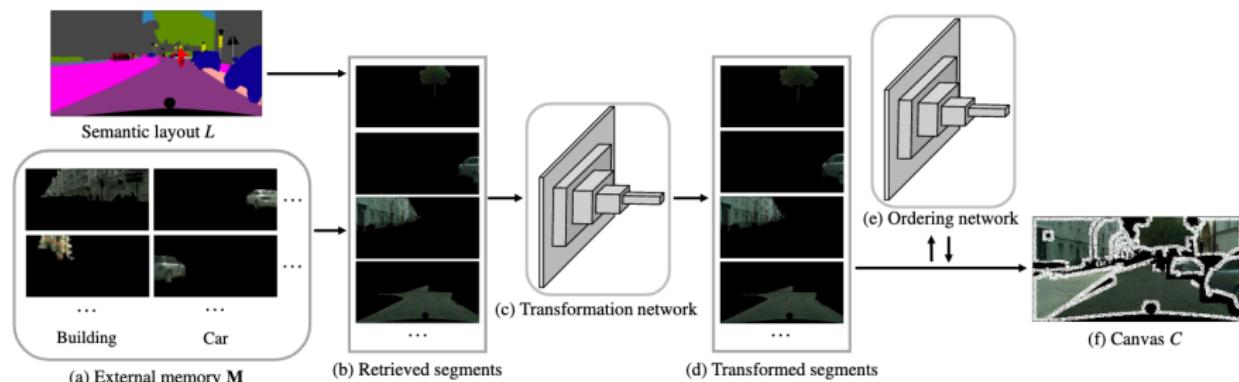


Figure: Network Architecture of CRN

## Related Work: SIMS

- SIMS = **S**emi-parametric **I**Mage **S**ynthesis
  - Image synthesis is performed by "**stitching**" parts of images together. The parts of the images stem from a memory bank of **image segments** which is created from a training set of images beforehand [5]



## Figure: Canvas Generator for SIMS

## Related Work: pix2pixHD

- Focus images with **high resolution** and **photorealism**
- **Approach:** using a coarse-to-fine generator and multi-scale discriminator architectures
- Decompose the generator into two **sub-networks**  $G_1$  and  $G_2$  to combine the **global** and **local** information [7]

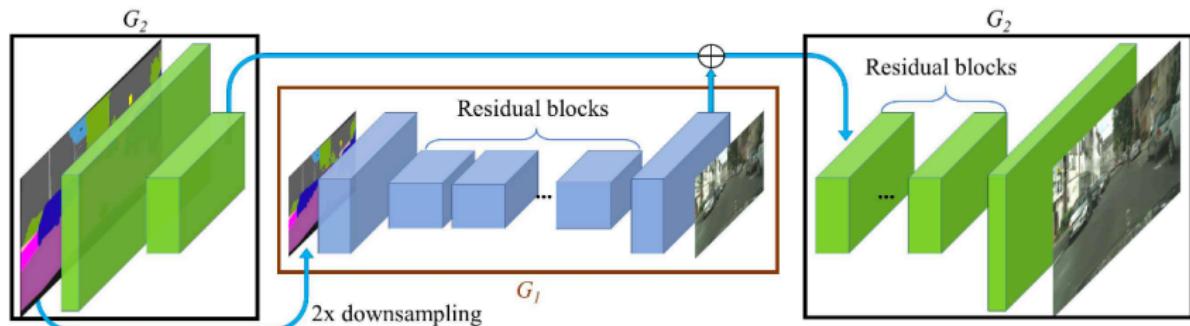


Figure: Network Architecture of pix2pixHD's Course-to-fine Generator

# Related Work: pix2pixHD

## Multi-scale Discriminator

- **Problem:** Discriminator needs **large receptive field** to differentiate between **high resolution** images. However, constructing a deeper network could lead to **overfitting** and a larger **memory** footprint
- **Solution:** Multi-scale discriminators: decompose into 3 identical discriminators ( $D_1, D_2, D_3$ ) with **different image scales**

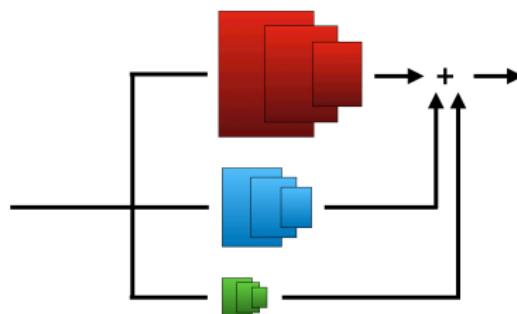


Figure: Architecture of the **pix2pixHD**'s Multi-Scale Discriminator

# Issues with Related Works

- **Current Approach:** Semantic information is direct input to neural network and processed through stacks of convolution, normalization, and non-linearity layers

# Issues with Related Works

- **Current Approach:** Semantic information is direct input to neural network and processed through stacks of convolution, normalization, and non-linearity layers
- **Problem:** Semantic information is not well preserved (normalization layers tend to "wash away" semantic information)

# Issues with Related Works

- **Current Approach:** Semantic information is direct input to neural network and processed through stacks of convolution, normalization, and non-linearity layers
- **Problem:** Semantic information is not well preserved (normalization layers tend to "wash away" semantic information)
- **Solution:** A novel conditional normalization method (SPADE) that modulates the activations using semantic layouts

# Outline

1 Introduction

2 Related Work

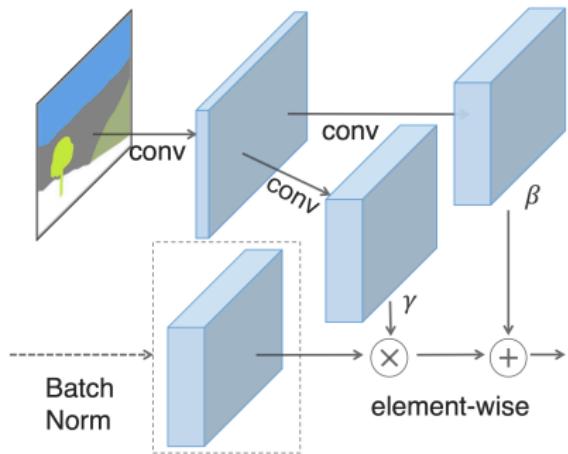
3 Semantic Image Synthesis

4 Experiments

5 Conclusion

6 References

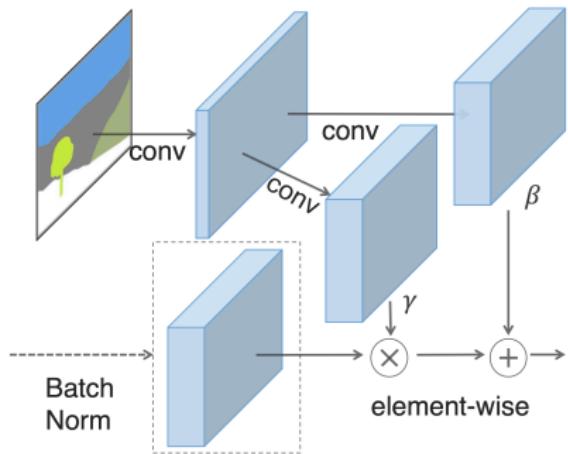
# Spatially-Adaptive Denormalization (SPADE) Layer



- ① Unconditional normalization of activations of previous layer with BatchNorm
- ② Denormalization with modulation parameters (scale  $\gamma$  and bias  $\beta$ )

Figure: The novel SPADE Layer

# Spatially-Adaptive Denormalization (SPADE) Layer



- ① Unconditional normalization of activations of previous layer with BatchNorm
- ② Denormalization with modulation parameters (scale  $\gamma$  and bias  $\beta$ )

Figure: The novel SPADE Layer

## Novelty

- $\gamma$  and  $\beta$  are learned and depend on location in segmentation mask!
- Modulation parameters encode semantic layout

# SPADE Generator

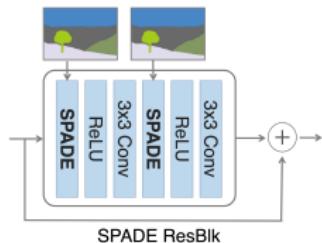


Figure: ResBlk

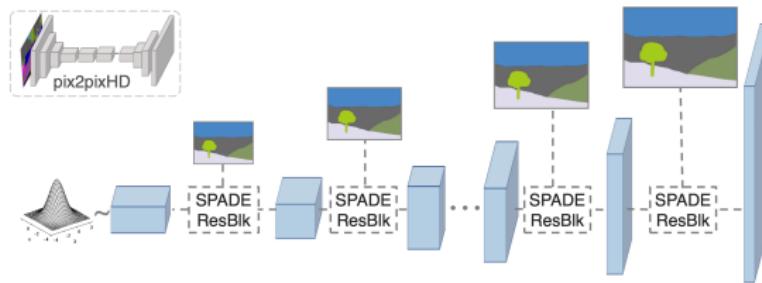


Figure: SPADE Generator

- ResBlk: residual block with skip connection

# SPADE Generator

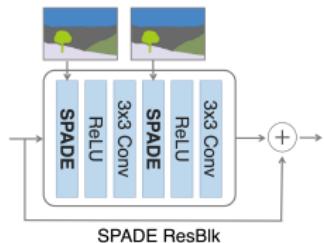


Figure: ResBlk

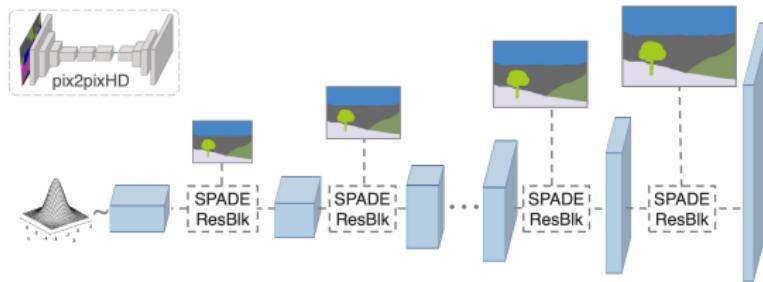


Figure: SPADE Generator

- ResBlk: residual block with skip connection
- Resized seg. masks influence generation through SPADE ResBlks

# SPADE Generator

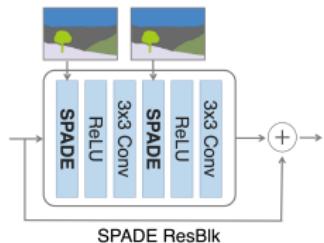


Figure: ResBlk

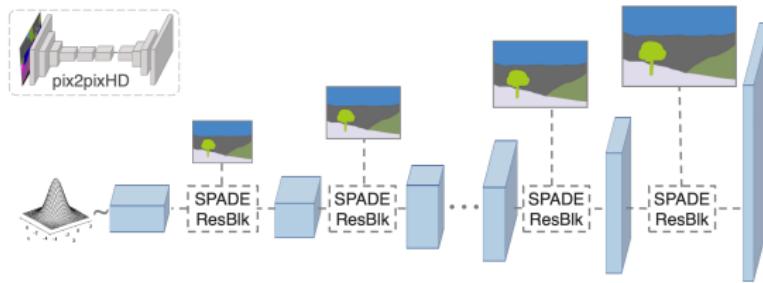


Figure: SPADE Generator

- ResBlk: residual block with skip connection
- Resized seg. masks influence generation through SPADE ResBlks
- Nearest neighbor upsampling

# SPADE Generator

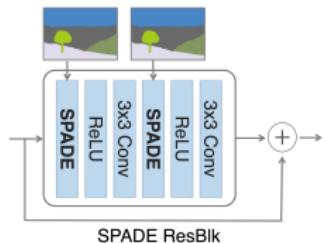


Figure: ResBlk

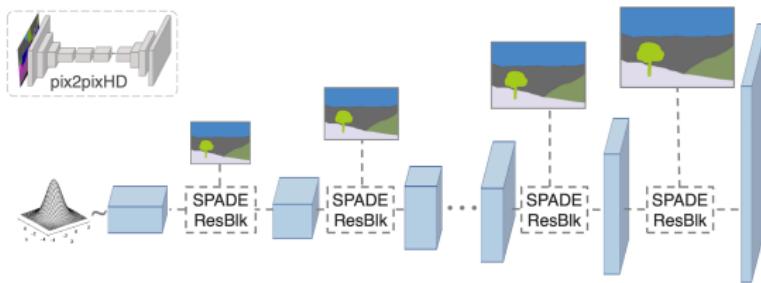


Figure: SPADE Generator

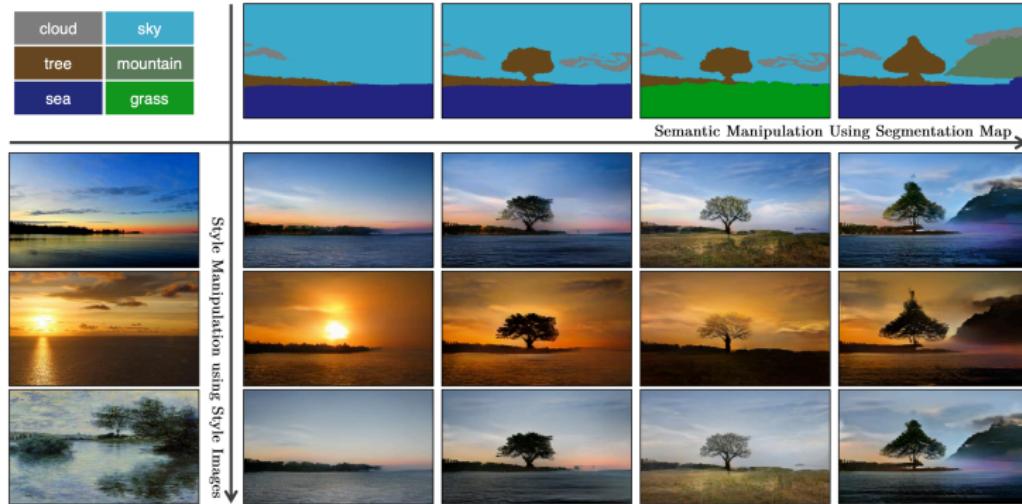
- ResBlk: residual block with skip connection
- Resized seg. masks influence generation through SPADE ResBlks
- Nearest neighbor upsampling
- Random noise fed to first layer instead of segmentation mask

# Multi-Modal Synthesis



- Different random inputs with the same segmentation mask lead to different appearances but same semantic layout

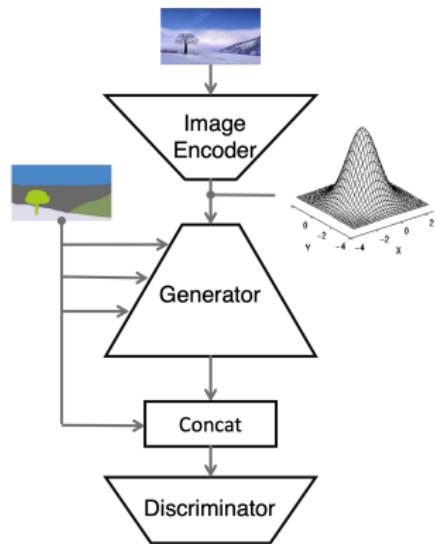
# Guided Image Synthesis



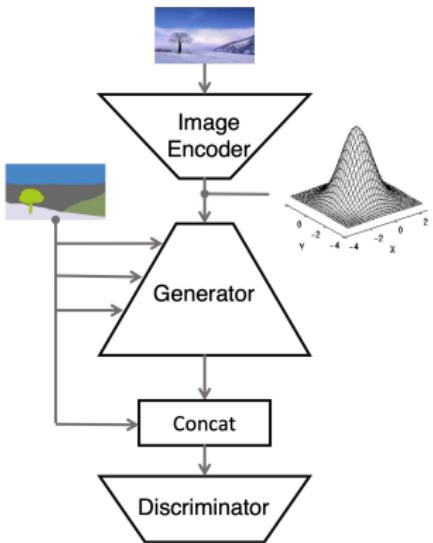
- Control semantics with segmentation mask and appearance with style image (style and semantics disentanglement)
- Interactive web application [GauGAN](#)

# Network Architecture

- **Image encoder** captures style of a real image in a latent representation.  
Outputs a mean vector  $\mu$  and a variance vector  $\sigma^2$

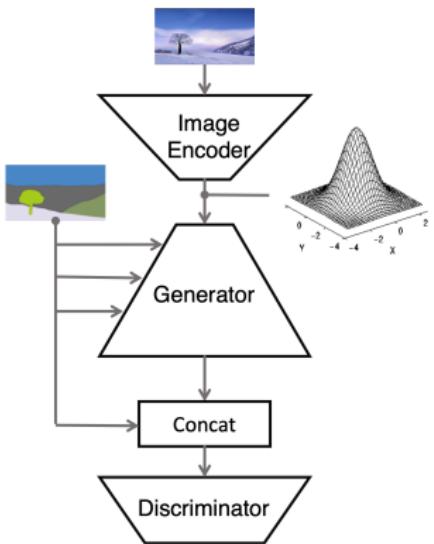


# Network Architecture



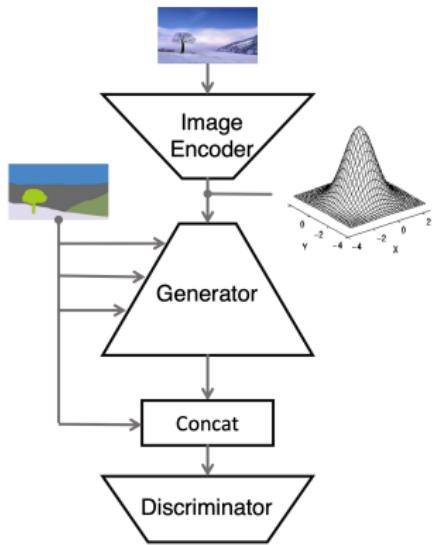
- **Image encoder** captures style of a real image in a latent representation.  
Outputs a mean vector  $\mu$  and a variance vector  $\sigma^2$
- **Generator** combines encoded style and seg. mask to reconstruct original image  
(Encoder + Generator = VAE)

# Network Architecture



- **Image encoder** captures style of a real image in a latent representation.  
Outputs a mean vector  $\mu$  and a variance vector  $\sigma^2$
- **Generator** combines encoded style and seg. mask to reconstruct original image  
(Encoder + Generator = VAE)
- **Concat** concatenates segmentation mask and generated image for comparison

# Network Architecture



- **Image encoder** captures style of a real image in a latent representation.  
Outputs a mean vector  $\mu$  and a variance vector  $\sigma^2$
- **Generator** combines encoded style and seg. mask to reconstruct original image  
(Encoder + Generator = VAE)
- **Concat** concatenates segmentation mask and generated image for comparison
- **Discriminator** same architecture and learning objective as pix2pixHD, but replace LS-GAN loss with Hinge loss.

# Outline

1 Introduction

2 Related Work

3 Semantic Image Synthesis

4 Experiments

5 Conclusion

6 References

# Main Datasets



Figure: COCO-Stuff



Figure: ADE20K

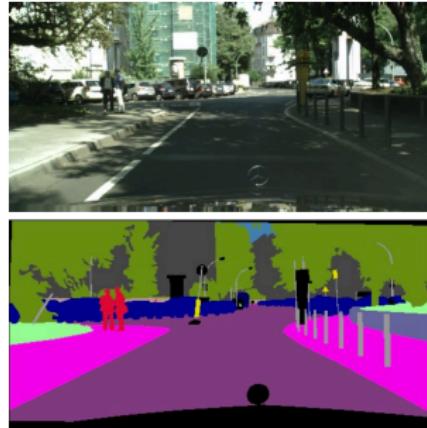


Figure: Cityscapes

Name	Train	Val	Classes	Description
COCO-Stuff	118k	5k	182	Challenging due to diversity
ADE20K	≈20k	2k	150	Similar to COCO, very diverse
Cityscapes	3k	0.5k	30	Street scene images

# Comparison of Qualitative Results

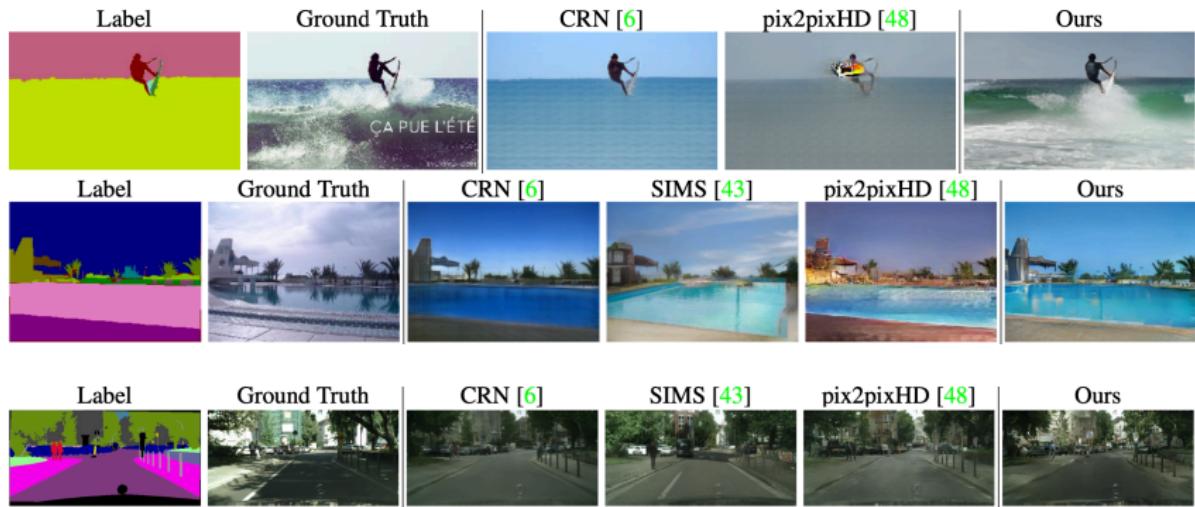


Figure: Top: COCO-Stuff, Middle: ADE20K, Bottom: Cityscapes

# Comparison of Quantitative Results

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	<b>37.4</b>	<b>67.9</b>	<b>22.6</b>	<b>38.5</b>	<b>79.9</b>	<b>33.9</b>	<b>30.8</b>	<b>82.9</b>	<b>63.3</b>	<b>62.3</b>	<b>81.9</b>	71.8

- Synthesized images are segmented with well-trained models and evaluated with performance metrics
- ➊ **Mean Intersection over Union:** What is the percentage overlap between predicted and ground truth mask?
  - ➋ **Pixel accuracy:** What is the percentage of correctly classified pixels?
  - ➌ **Fréchet Inception Distance:** What is the distance between distributions of feature vectors?

# Why does the SPADE work better?

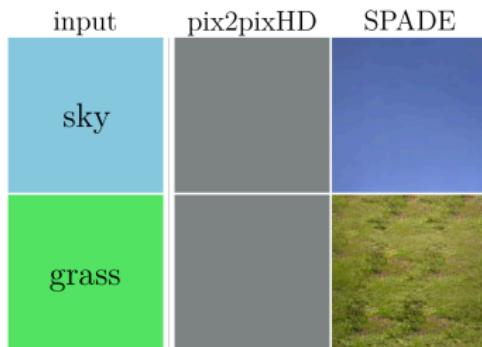


Figure: Semantic information loss after normalization layer

- Unconditional normalization layers (e.g. InstanceNorm) loose semantic info of uniform masks as **normalized activations are zero**

# Why does the SPADE work better?

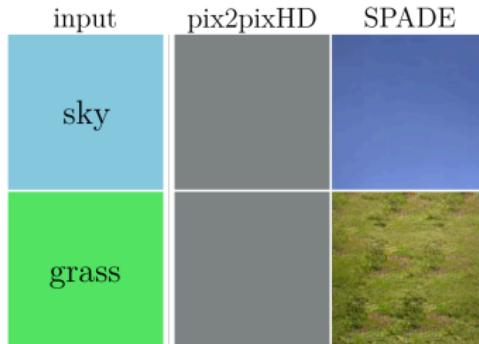


Figure: Semantic information loss after normalization layer

- Unconditional normalization layers (e.g. InstanceNorm) loose semantic info of uniform masks as **normalized activations are zero**
- SPADE better **preserves semantic information** because segmentation mask is not normalized but only modulated

# Why does the SPADE work better?

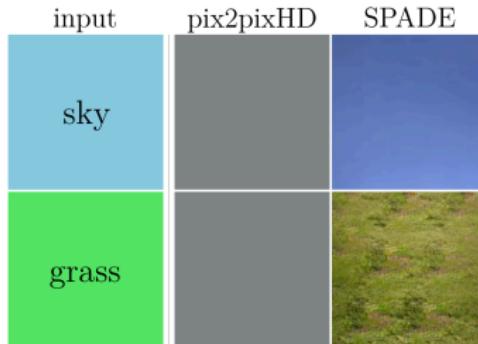


Figure: Semantic information loss after normalization layer

- Unconditional normalization layers (e.g. InstanceNorm) loose semantic info of uniform masks as **normalized activations are zero**
- SPADE better **preserves semantic information** because segmentation mask is not normalized but only modulated
- SPADE also improves performance of traditional architectures!

# Outline

1 Introduction

2 Related Work

3 Semantic Image Synthesis

4 Experiments

5 Conclusion

6 References

# Conclusion

Task	Dataset	Model	Metric Name	Metric Value	Global Rank	Compare
Image-to-Image Translation	ADE20K Labels-to-Photos	SPADE	mIoU	38.5	# 2	<a href="#">See all</a>
Image-to-Image Translation	ADE20K Labels-to-Photos	SPADE	Accuracy	79.9%	# 2	<a href="#">See all</a>
Image-to-Image Translation	ADE20K Labels-to-Photos	SPADE	FID	33.9	# 2	<a href="#">See all</a>
Image-to-Image Translation	ADE20K-Outdoor Labels-to-Photos	SPADE	mIoU	30.8	# 1	<a href="#">See all</a>
Image-to-Image Translation	ADE20K-Outdoor Labels-to-Photos	SPADE	Accuracy	82.9%	# 1	<a href="#">See all</a>
Image-to-Image Translation	ADE20K-Outdoor Labels-to-Photos	SPADE	FID	63.3	# 1	<a href="#">See all</a>
Image-to-Image Translation	Cityscapes Labels-to-Photo	SPADE	Per-pixel Accuracy	81.9%	# 2	<a href="#">See all</a>
Image-to-Image Translation	Cityscapes Labels-to-Photo	SPADE	mIoU	62.3	# 2	<a href="#">See all</a>

Figure: SPADE ranking as of March 2020 [3]

- Introduced spatially-adaptive normalization (SPADE) layer
- SPADE network outperforms the 2019 state-of-the-art methods by a large margin and is still top-performing (#1 is [2])

# Outline

1 Introduction

2 Related Work

3 Semantic Image Synthesis

4 Experiments

5 Conclusion

6 References

# References I

- [1] Q. Chen and V. Koltun.  
Photographic image synthesis with cascaded refinement networks.  
In *Proceedings of the IEEE international conference on computer vision*, pages 1511–1520, 2017.
- [2] X. Liu, G. Yin, J. Shao, X. Wang, and H. Li.  
Learning to predict layout-to-image conditional convolutions for semantic image synthesis, 2019.
- [3] Papers With Code.  
Evaluation Results SPADE.  
<https://paperswithcode.com/paper/semantic-image-synthesis-with-spatially>, 2020.

## References II

- [4] T. Park, M.-Y. Liu, T.-C. Wang, and J.-Y. Zhu.  
Semantic image synthesis with spatially-adaptive normalization.  
In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019.
- [5] X. Qi, Q. Chen, J. Jia, and V. Koltun.  
Semi-parametric image synthesis.  
In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8808–8816, 2018.
- [6] M. Thoma.  
A survey of semantic segmentation, 2016.

## References III

- [7] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro.  
High-resolution image synthesis and semantic manipulation with conditional gans.  
In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018.