

Semantic Image Synthesis with Spatially-Adaptive Normalization

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

Presented by - Ayushi Bansal, Bhargav Sundararajan, Mridula Gupta

Problem - Semantic Image Synthesis

Generate photorealistic image based on given input semantic layout

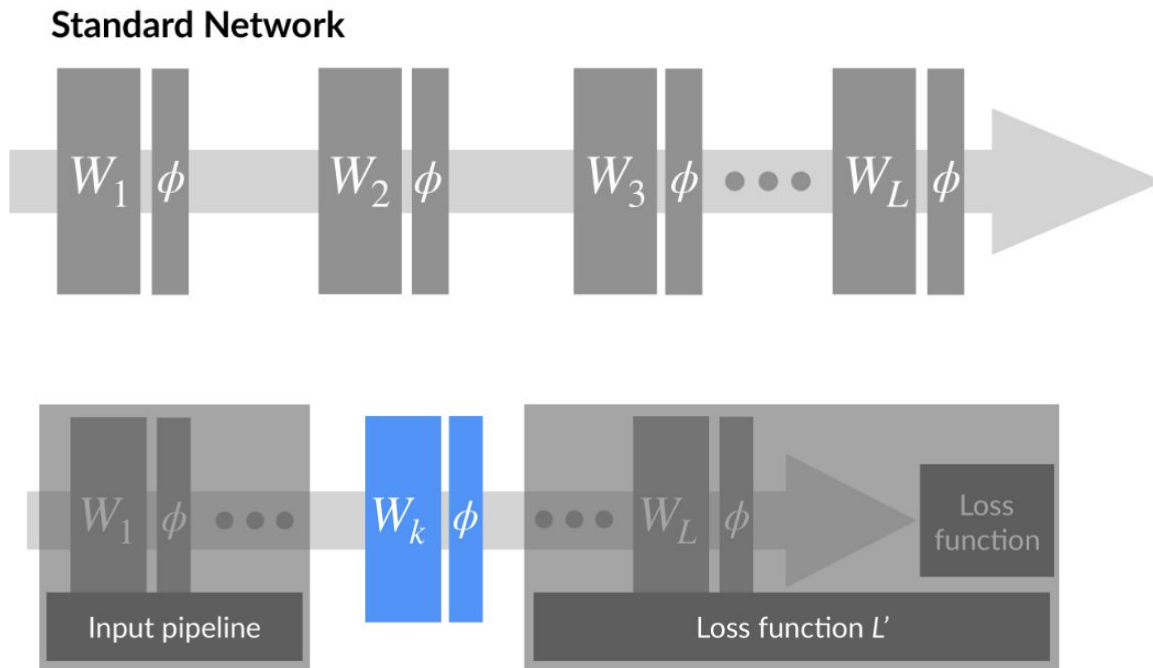


Normalization

- Adjusting values measured on different scales to a notionally common scale.
- Done by subtracting the mean and dividing by the standard deviation.
- Usually a data preprocessing step.

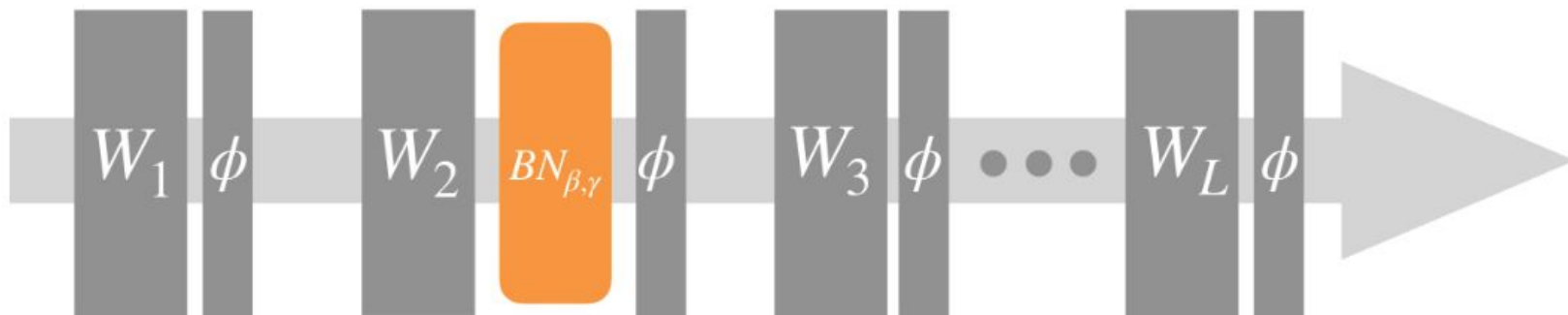
Normalization

It is basically used to adjust and scale activations in neural network



Normalization (Batch Norm)

Adding a BatchNorm layer (between weights and activation function)



$$\bar{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

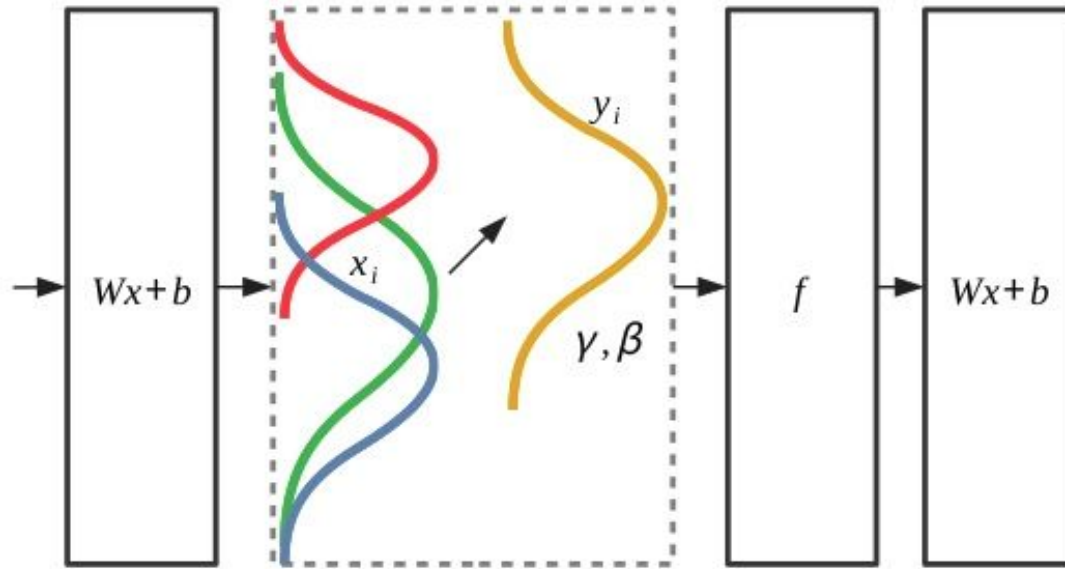
Mean

Variance

$$y_i = \gamma \bar{x}_i + \beta$$

Trainable Parameters (Scale and Shift)

Why Normalization between layers?



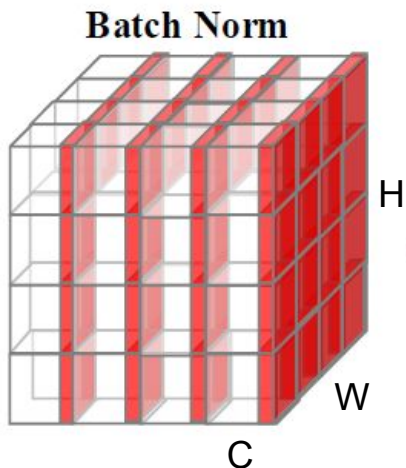
Ensures that output statistics of the layer are fixed

Unconditional Normalization

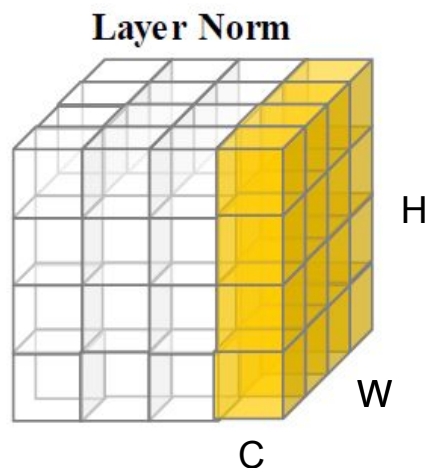
- No dependence on external data during normalization process
- Types
 - Batch Normalization
 - Instance Normalization
 - Layer Normalization
 - Group Normalization
 - Weight Normalization

Unconditional Normalization

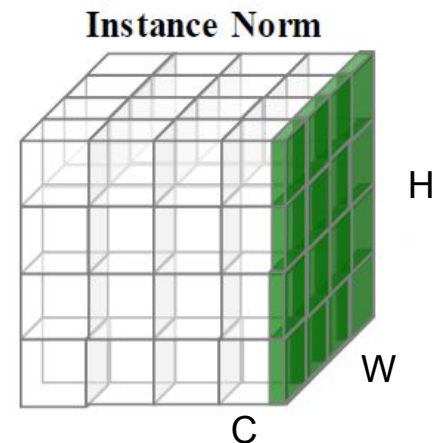
H	Height of Image
W	Width of Image
C	Channels
N	Number of Instances in a batch



Single Channel
over Multiple
training instances



Multiple Channel
over single training
instance



Single Channel
over single training
instance

Conditional Normalization

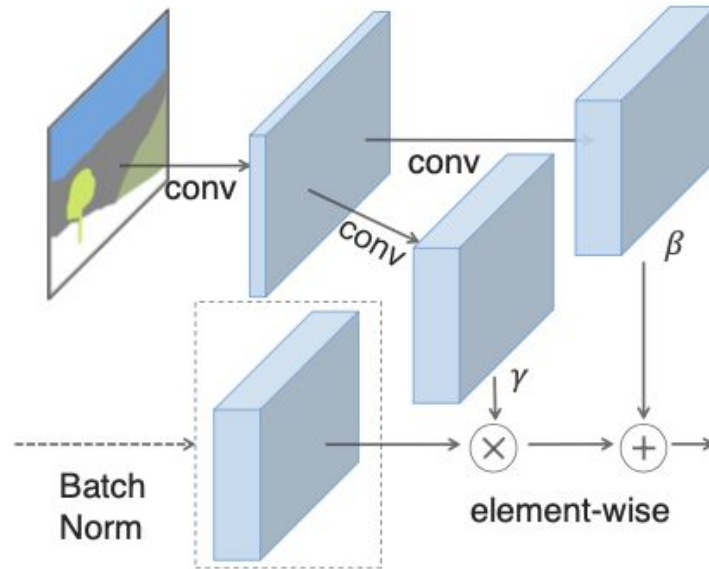
- External data is used to condition the normalization.
- Used in Style transfer and Visual Question Answering
- Example: Conditional Instance Normalization and Adaptive Instance Normalization

$$\text{CIN}(x; s) = \gamma^s \left(\frac{x - \mu(x)}{\sigma(x)} \right) + \beta^s$$

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

SPADE Model - Conditional Normalization

Unlike prior conditional normalization methods, γ and β are not vectors, but tensors with spatial dimensions. Hence, the name Spatially Adaptive Normalization



SPADE Model - Conditional Normalization

Previous Activation

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

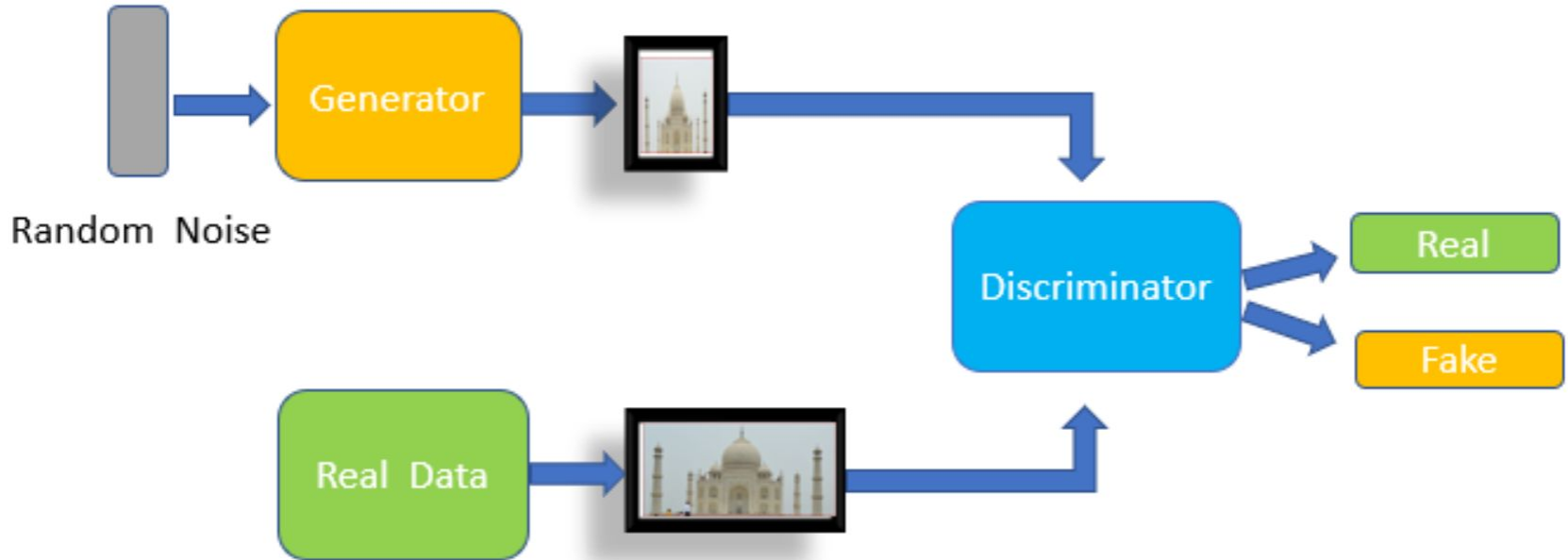
($n \in N, c \in C_i, y \in H_i, x \in W_i$)

Learned modulation parameters

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

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

GANs (Generative Adaptive Networks)

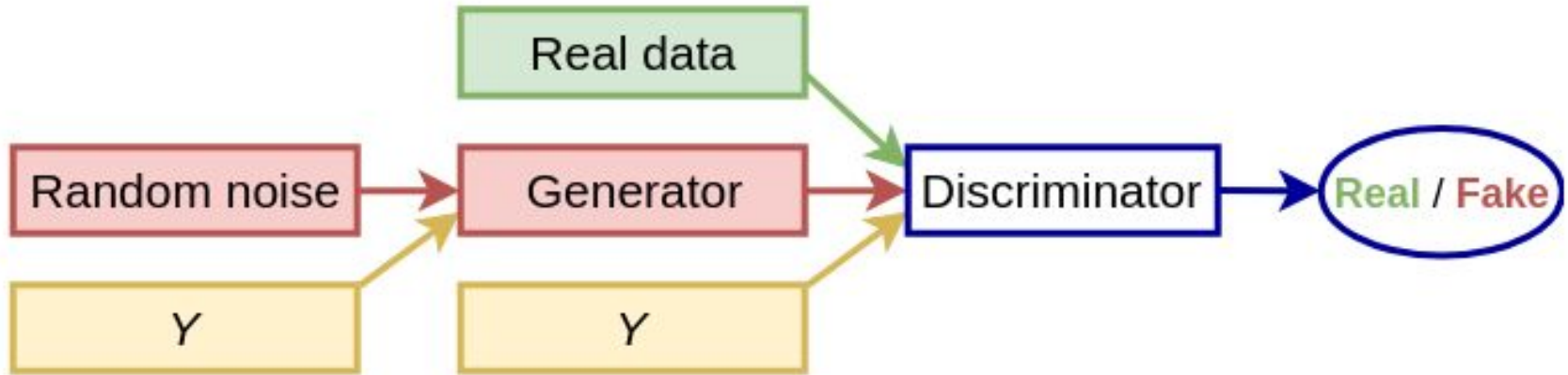


GANs (Generative Adaptive Networks)

- **Generator:** Learns to generate plausible image
- **Discriminator:** Learns to distinguish generator's fake image from real image
- **Training**
 - Generator and Discriminator alternatively trained (1 epoch each)
 - Converge the results

Conditional GANs

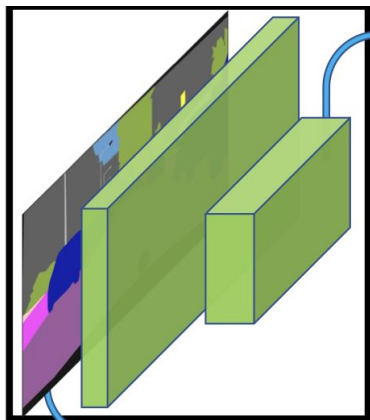
In addition to random noise (for generator) and input image (for discriminator), we are giving **one-hot encoded vector** as input (Y)



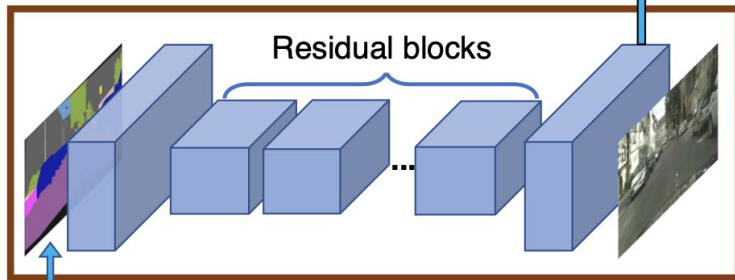
Pix2PixHD Model

Local Enhancer Generator

G_2



Global Generator G_1

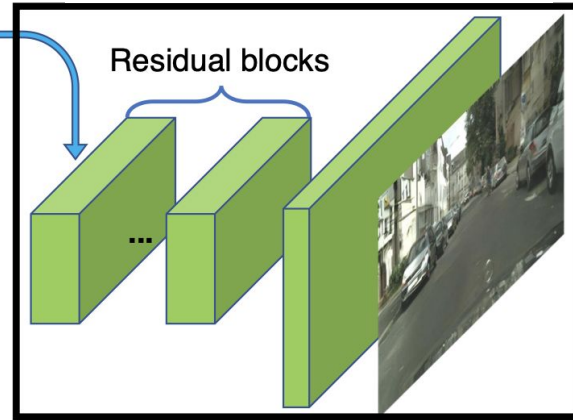


2x downsampling

G_1

Local Enhancer Generator

G_2



Residual blocks

Training

- Train global generator G_1 separately first
- Fine-tune whole network with G_1 and G_2

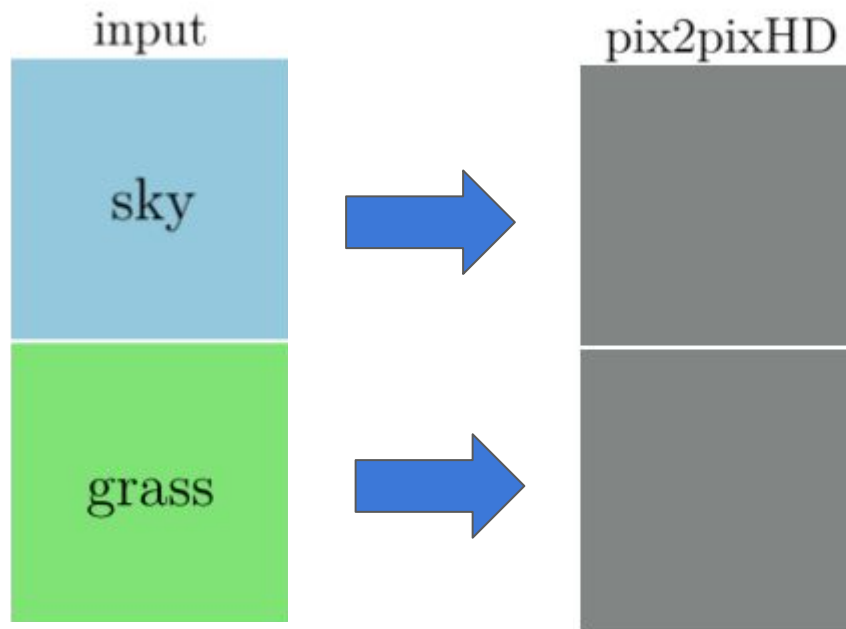
Pix2PixHD Model

Two main contributions:

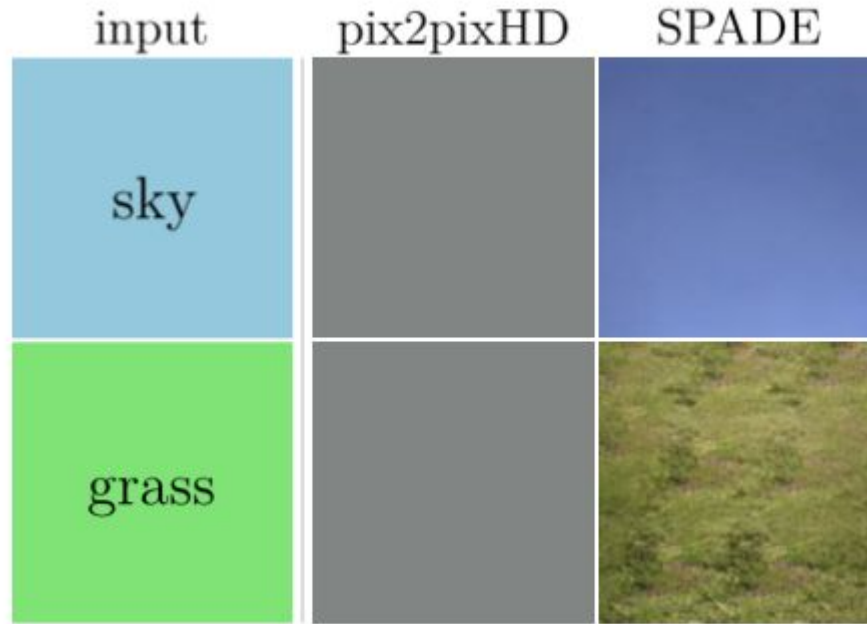
1. Instance-level object segmentation information is used, which can separate different object instances within the same category.
2. Generate diverse results given the same input label map, allowing the user to edit the appearance of the same object interactively.

Motivation - Pix2PixHD

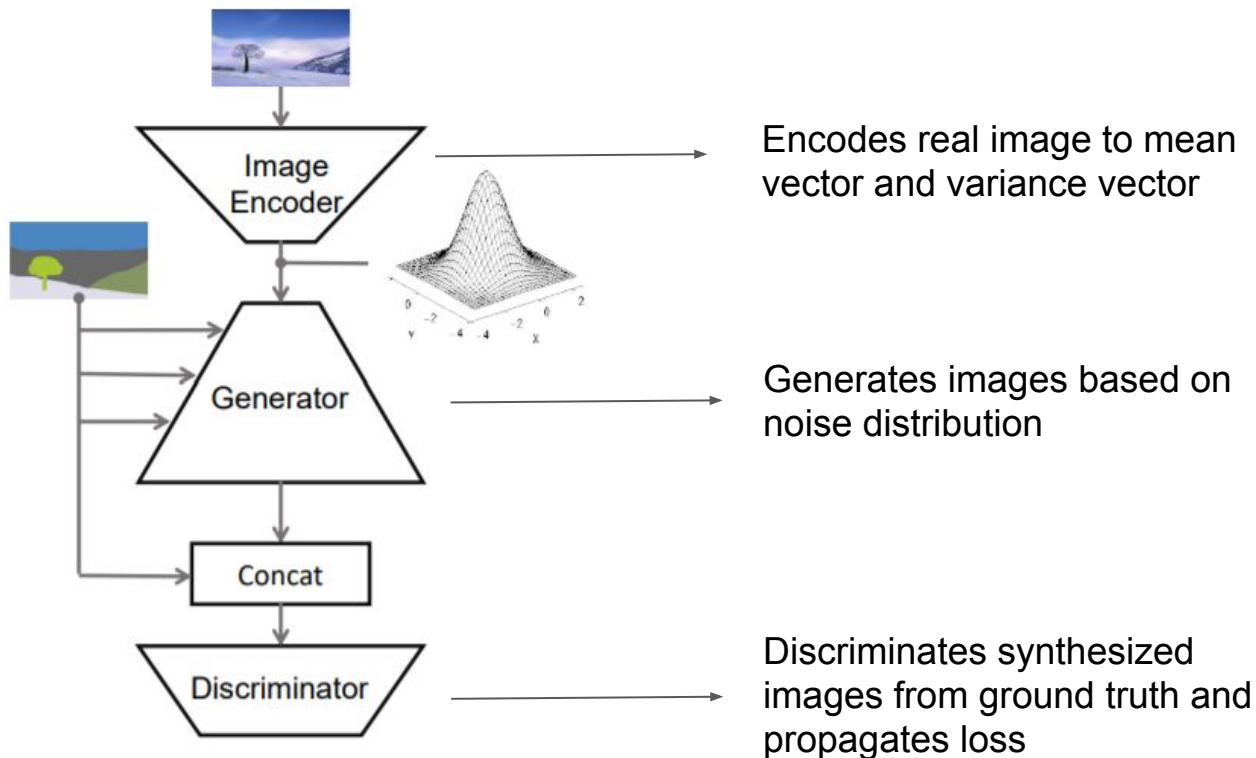
Normalization “washes”
away semantic information
when applied to uniform or
flat segmentation mask



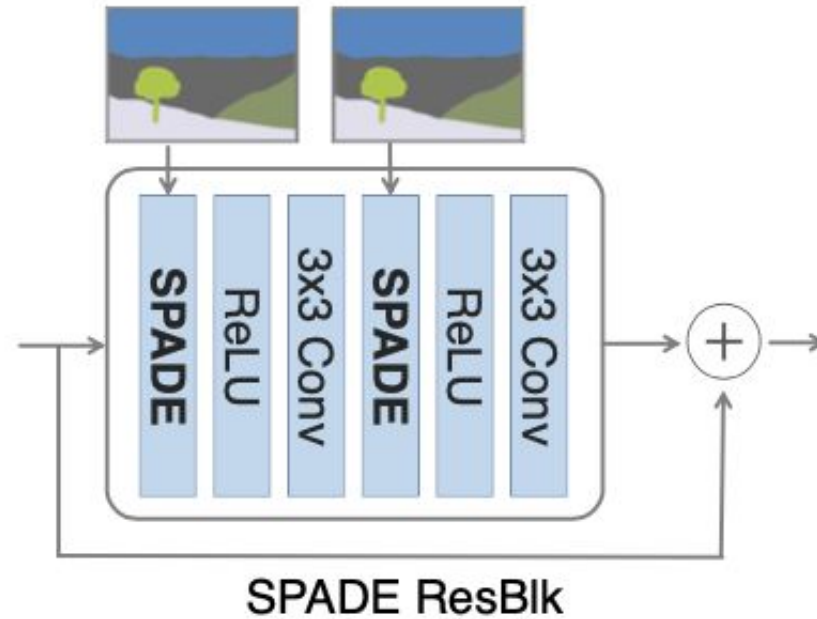
Resolves Pix2PixHD problem



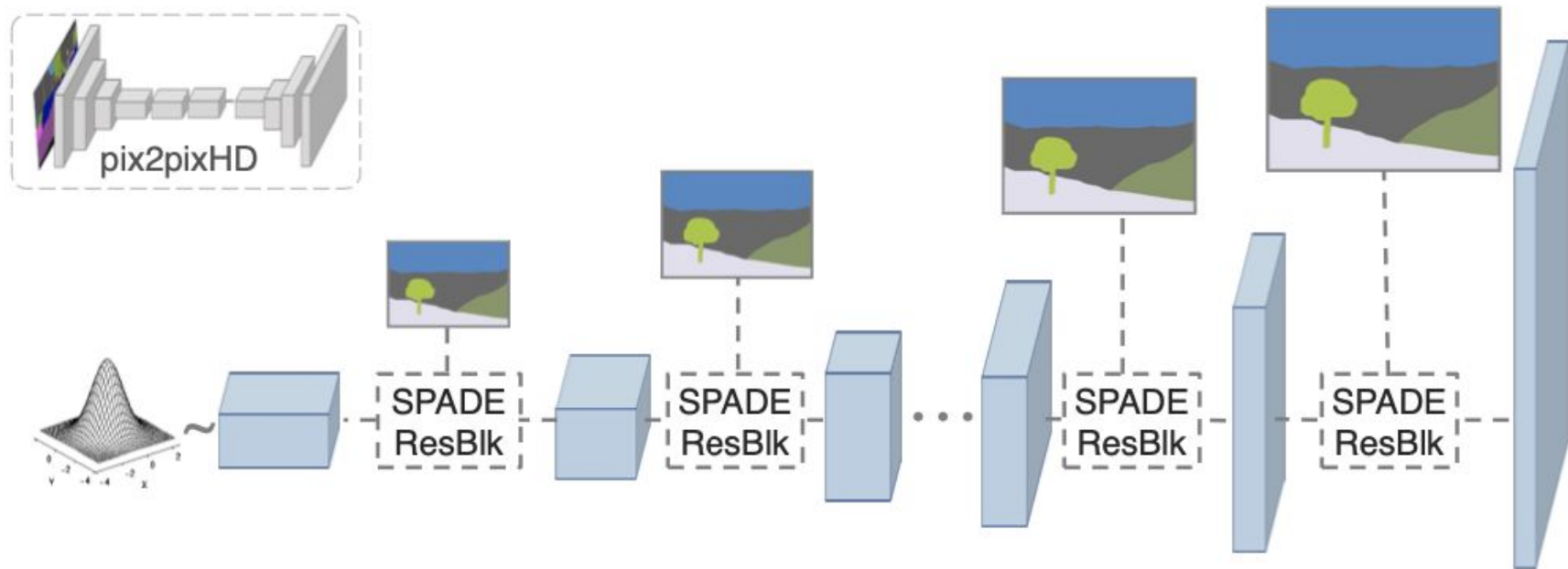
The Model Architecture



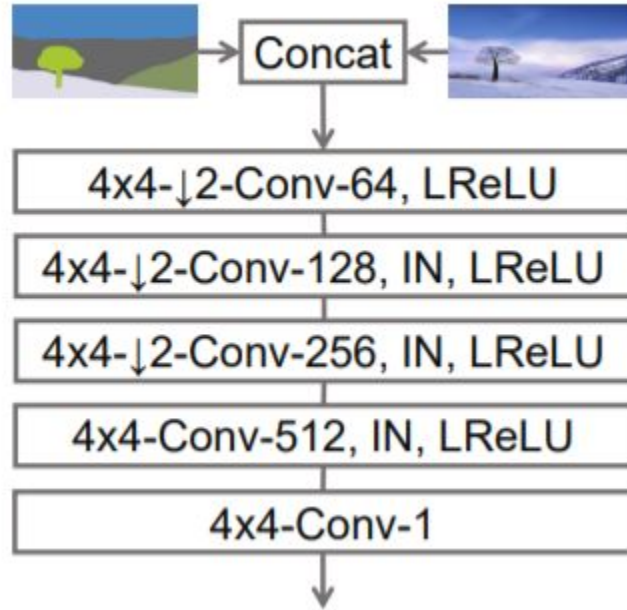
SPADE Model - Residual Block



SPADE Model - Generator



Discriminator



Loss Function

\mathcal{L}_{FM} compares the feature map of every intermediate layer of the discriminator.

$$\min_G \left(\left(\max_{D_1, D_2, D_3} \sum_{k=1,2,3} \mathcal{L}_{\text{GAN}}(G, D_k) \right) + \lambda \sum_{k=1,2,3} \mathcal{L}_{\text{FM}}(G, D_k) \right)$$

$$\mathcal{L}_{\text{FM}}(G, D_k) = \mathbb{E}_{(\mathbf{s}, \mathbf{x})} \sum_{i=1}^T \frac{1}{N_i} [\|D_k^{(i)}(\mathbf{s}, \mathbf{x}) - D_k^{(i)}(\mathbf{s}, G(\mathbf{s}))\|_1],$$

↓
Number of
Intermediate Layers

Loss - Hinge Loss

Prediction for
real Image


$$V_D(\hat{G}, D) = \mathbb{E}_{\mathbf{x} \sim q_{\text{data}}(\mathbf{x})} [\min(0, -1 + D(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\min(0, -1 - D(\hat{G}(\mathbf{z})))]$$

$$V_G(G, \hat{D}) = - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\hat{D}(G(\mathbf{z}))],$$

Loss - Hinge Loss

$$V_D(\hat{G}, D) = \mathbb{E}_{\mathbf{x} \sim q_{\text{data}}(\mathbf{x})} [\min(0, -1 + D(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\min(0, -1 - D(\hat{G}(\mathbf{z})))]$$

Prediction for
Generated Image


$$V_G(G, \hat{D}) = - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\hat{D}(G(\mathbf{z}))],$$

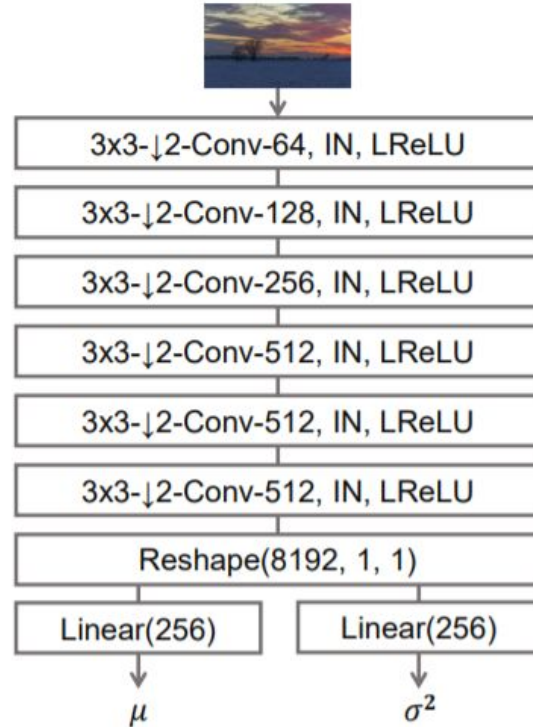
Loss - Hinge Loss

$$V_D(\hat{G}, D) = \mathbb{E}_{\mathbf{x} \sim q_{\text{data}}(\mathbf{x})} [\min(0, -1 + D(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\min(0, -1 - D(\hat{G}(\mathbf{z})))]$$

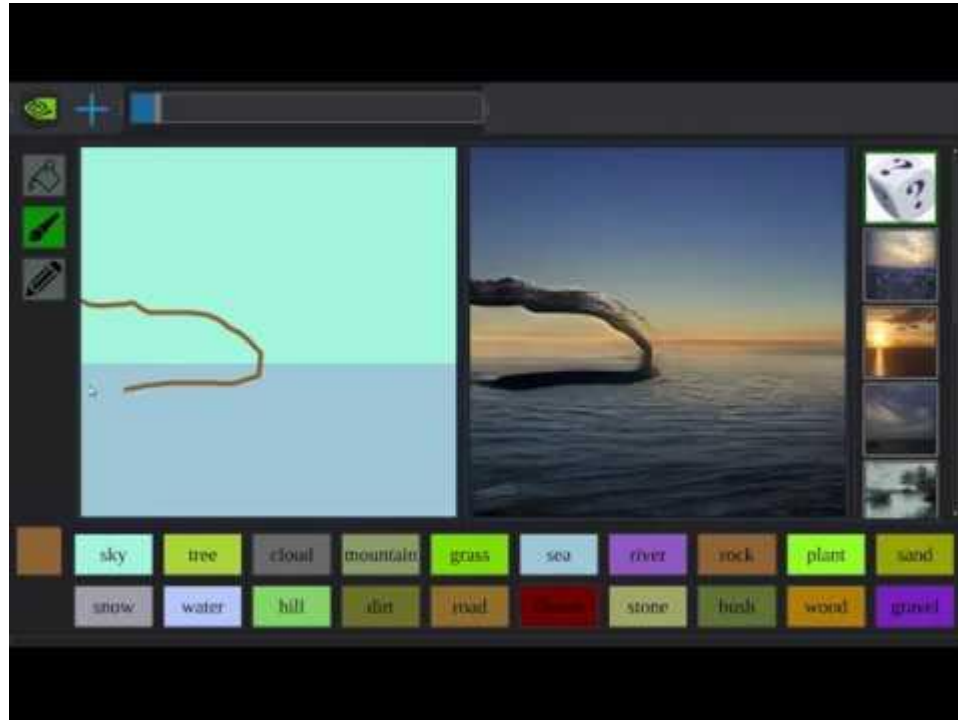
$$V_G(G, \hat{D}) = - \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\hat{D}(G(\mathbf{z}))],$$

↓
Loss of
Generator

Image Encoder - Style Transfer



GauGAN - Interactive Demo



Datasets

COCO-stuff

Derived from COCO
of classes: 182
of training: 11.8k
of validation: 5k

ADE20K

of classes: 150
of training: 20.2k
of validation: 2k

ADE20K-outdoor

Only outdoor images

Cityscapes dataset

Street scenes in
German cities
of training: 3k
of validation: 500

Flickr Landscapes

of training: 41k
of validation: 1k

Experimental Setup

Hyperparameters:

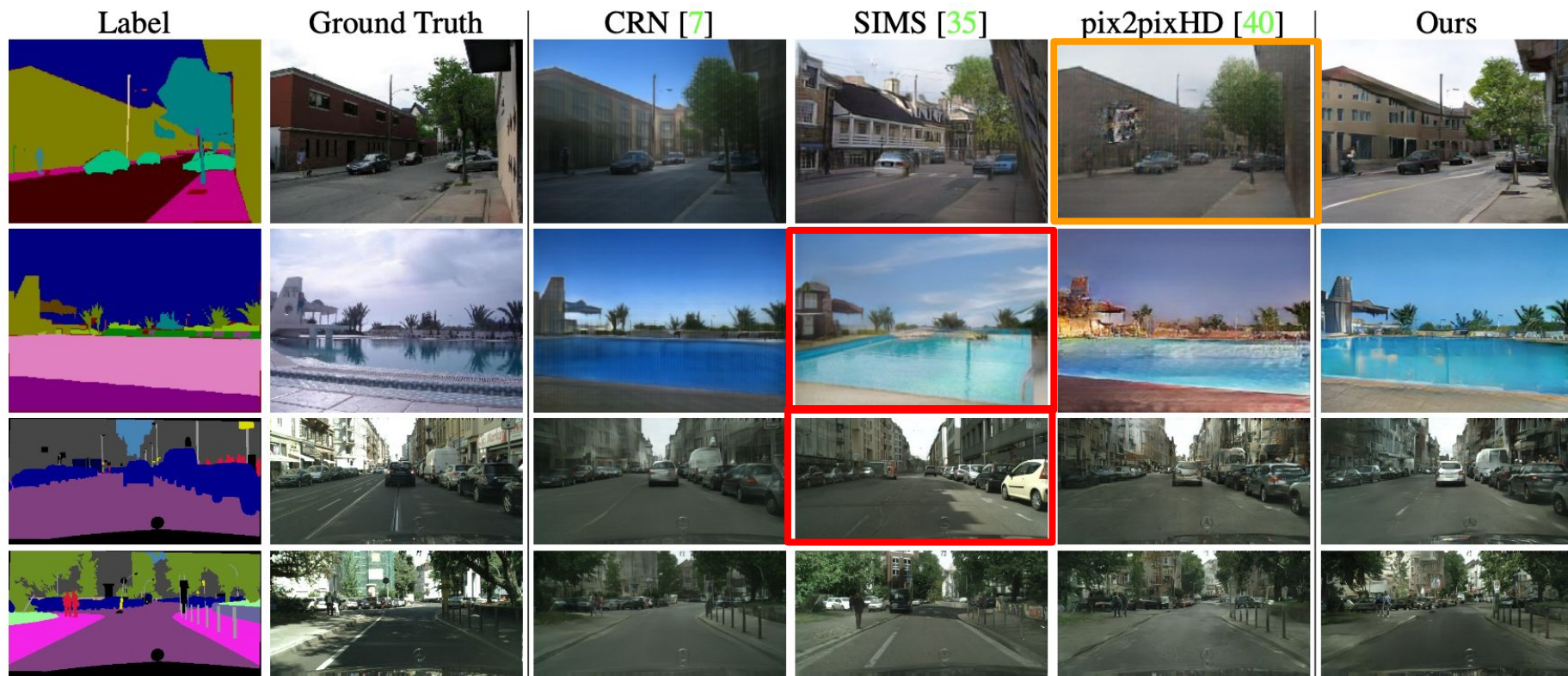
- Learning Rate (Generator & Discriminator): 0.0001 and 0.0004
- ADAM Optimizer: $\beta_1 = 0$, $\beta_2 = 0.999$

System Setup: NVIDIA DGX1 with 8 V100 GPUs

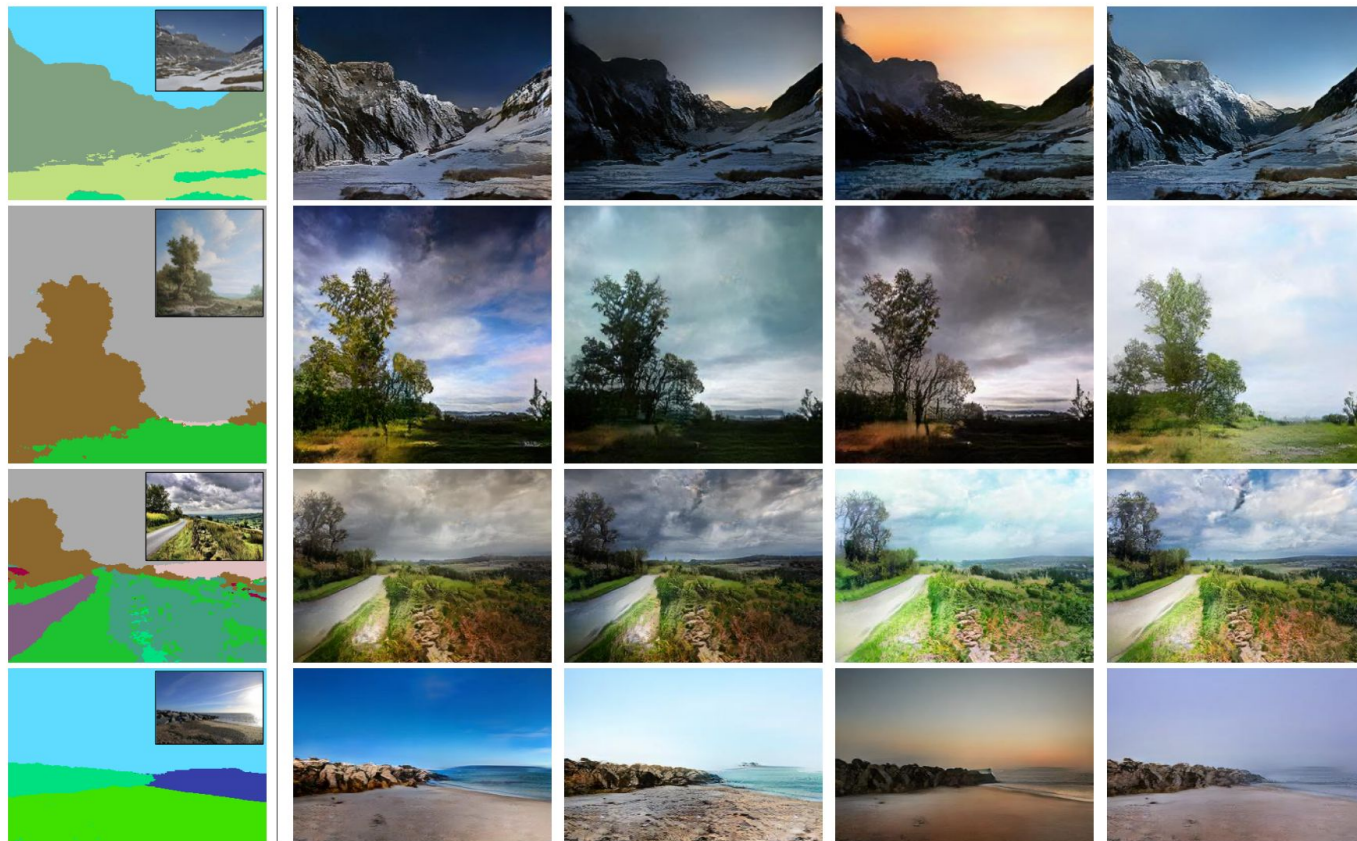
Baseline Models:

1. pix2pixHD: State-of-the-art GAN-based model
2. CRN (Cascaded Refinement Network): Deep Network refines the output from low to high resolution
3. SIMs (Semi-parametric IMage synthesis): Composites real segments from training set and refines boundaries

Results - Qualitative on ADE20K outdoor & Cityscapes

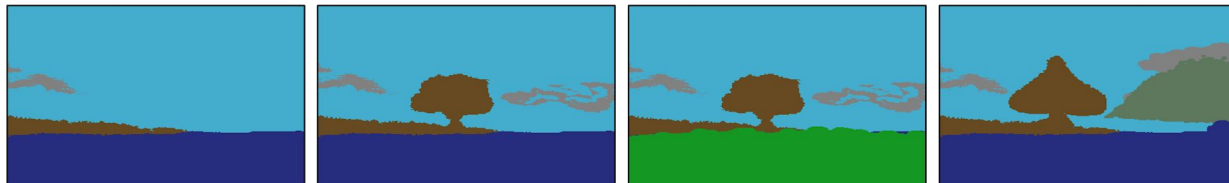


Results - Multimodal

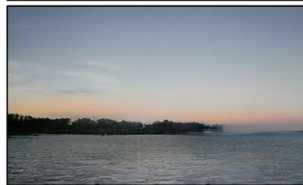
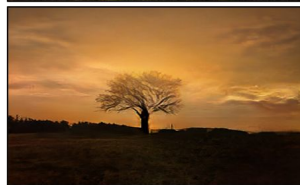
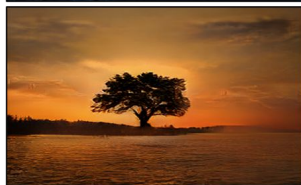
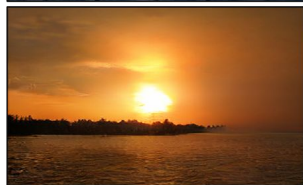
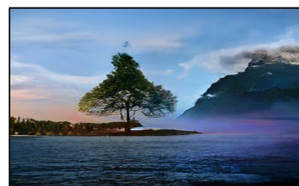
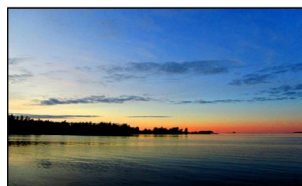


Results - Semantic and Style Control

cloud	sky
tree	mountain
sea	grass



Semantic Manipulation Using Segmentation Map →



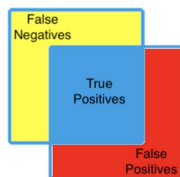
Stylization using Guide Images ↓

Experiment - Metrics



Segment Label Comparison:

1. mIoU (mean Intersection over Union)
IoU: Area of Overlap/Area of Union
mIoU: Mean IoU for each class
2. Accu (Pixel Accuracy): % of pixels classified correctly i.e. TP



$$IoU = \frac{TP}{(TP + FP + FN)}$$

Images Comparison:

- FID (Fréchet Inception Distance) - Distance between the distributions of synthesized results and the distribution of real images.

For mIoU and pixel accuracy, higher is better. For FID, lower is better.

Result - Quantitative

Method	COCO-Stuff			ADE20K			ADE20K-outdoor			Cityscapes		
	mIoU	accu	FID	mIoU	accu	FID	mIoU	accu	FID	mIoU	accu	FID
CRN [7]	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 [35]	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 [40]	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 Theirs	37.4	67.9	22.6	38.5	79.9	33.9	30.8	82.9	63.3	62.3	81.9	71.8

Outperforms current leading methods in semantic segmentation scores

Result - Human Evaluation

Randomly generated 500 questions for each dataset, and each question answered by 5 different human evaluators

Dataset	Ours vs. CRN	Ours vs. pix2pixHD	Ours vs. SIMS
COCO-Stuff	79.76	86.64	N/A
ADE20K	76.66	83.74	N/A
ADE20K-outdoor	66.04	79.34	85.70
Cityscapes	63.60	53.64	51.52

Users preferred the results of the proposed method over the competing methods

Result - Ablation Study

- **pix2pixHD++**: Including all the techniques with pix2pixHD except SPADE
- **pix2pixHD++ w/ Concat**: Concat the segment mask input at all the intermediate layers
- **pix2pixHD++ w/ SPADE**: Strong baseline with SPADE
- Also, compare models with different capacity generators

Method	#param	COCO.	ADE.	City.
decoder w/ SPADE (Ours)	96M	35.2	38.5	62.3
compact decoder w/ SPADE	61M	35.2	38.0	62.5
decoder w/ Concat	79M	31.9	33.6	61.1
pix2pixHD++ w/ SPADE	237M	34.4	39.0	62.2
pix2pixHD++ w/ Concat	195M	32.9	38.9	57.1
pix2pixHD++	183M	32.7	38.3	58.8
compact pix2pixHD++	103M	31.6	37.3	57.6
pix2pixHD [40]	183M	14.6	20.3	58.3

Comparison done with respect to parameters used and mIOU scores

Result - SPADE generator variations

- Two inputs - Segmentation map, random noise input
- Varying kernel size, different capacity, and types of normalization

Method	COCO	ADE20K	Cityscapes
segmap input	35.2	38.5	62.3
random input	35.3	38.3	61.6
kernelsize 5x5	35.0	39.3	61.8
kernelsize 3x3	35.2	38.5	62.3
kernelsize 1x1	32.7	35.9	59.9
#params 141M	35.3	38.3	62.5
#params 96M	35.2	38.5	62.3
#params 61M	35.2	38.0	62.5
Sync Batch Norm	35.0	39.3	61.8
Batch Norm	33.7	37.9	61.8
Instance Norm	33.9	37.4	58.7

All scores are mIOU scores, Standard used bold ones

Contributions

- Semantics of image is captured by adding **SP**atially-**A**daptive **DE**normalization in the pix2pixHD architecture.
- Users can control both semantic and style for image synthesis

Strengths

- Diverse results for the same input, depending on the user's interaction with objects (Style and Semantic support)
- SPADE resblk can be integrated with any existing architecture
- Reduced the number of trainable parameters and improved efficiency, when compared to previous state-of-the-art pix2pixHD
- Created semantic image synthesis interface(Live Demo) to interact and play on a canvas

Weaknesses

- Doesn't train well for fine-grained details like faces
- Limited number of classes in demo
- Can be easily forced to generate un-natural shapes.



Future Work

- Can be extended to video synthesis
- Technique can be used for image super-pixelation task also
- In-the-wild technique can be integrated to put a constraint on the shapes of the object generated

Questions?