# You Only Need Adversarial Supervision for Semantic Image Synthesis

## ICLR 2021

Edgar Schönfeld[1]*, Vadim Sushko[1]*, Dan Zhang[1],
Jürgen Gall[2], Bernt Schiele[3], Anna Khoreva[1]

*Equal contribution

[1]Bosch Center for Artificial Intelligence
[2]University of Bonn
[3]Max Planck Institute for Informatics
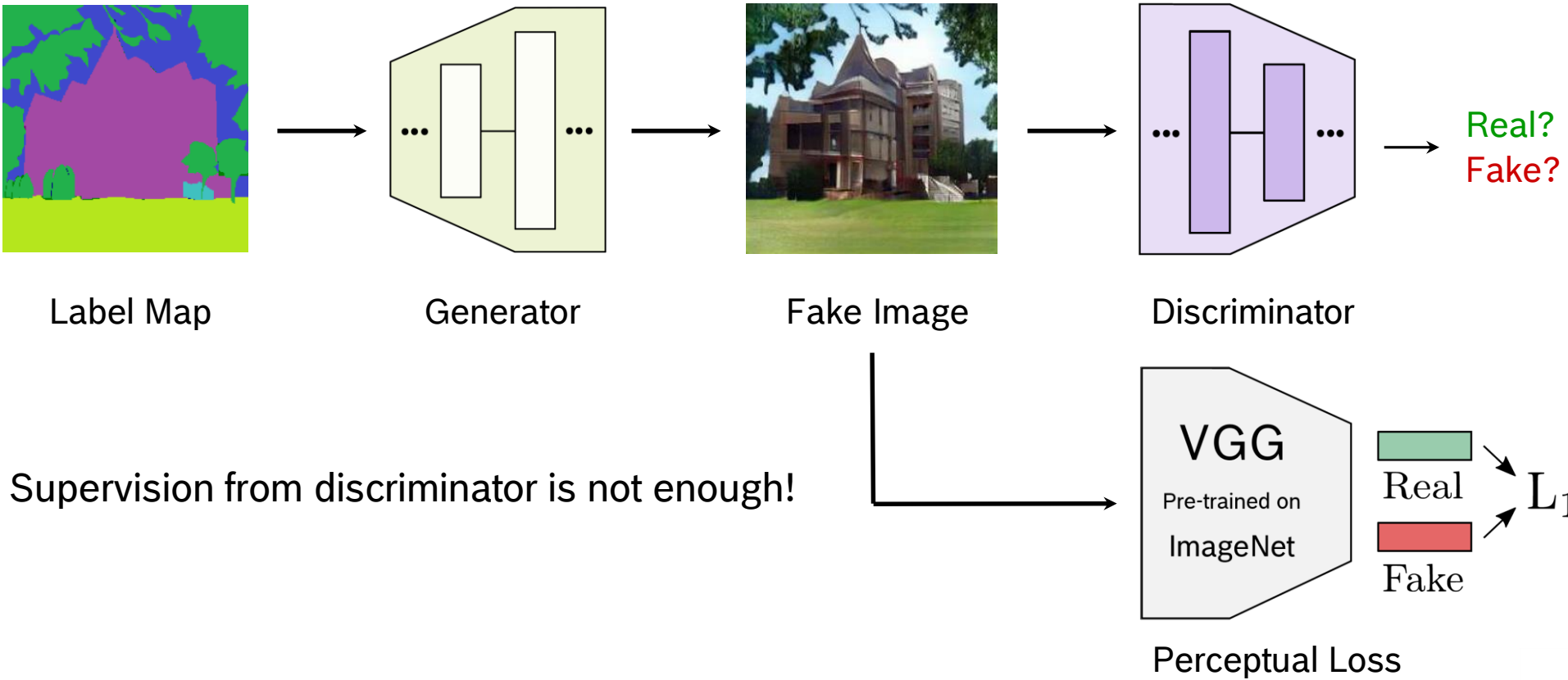
BOSCH

# Semantic image synthesis
## Problem statement

**Our goal:** multi-modal photorealistic image generation in alignment with a given semantic label map

BOSCH

# Limitations of previous GAN methods
## Perceptual loss



Label Map        Generator        Fake Image        Discriminator      Real? Fake?

Supervision from discriminator is not enough!

**VGG** Pre-trained on **ImageNet**

Real

Fake

$L_1$

Perceptual Loss

**BOSCH**

# Limitations of previous GAN methods
## Perceptual loss



**Real Image**

**Fake Image**

VGG
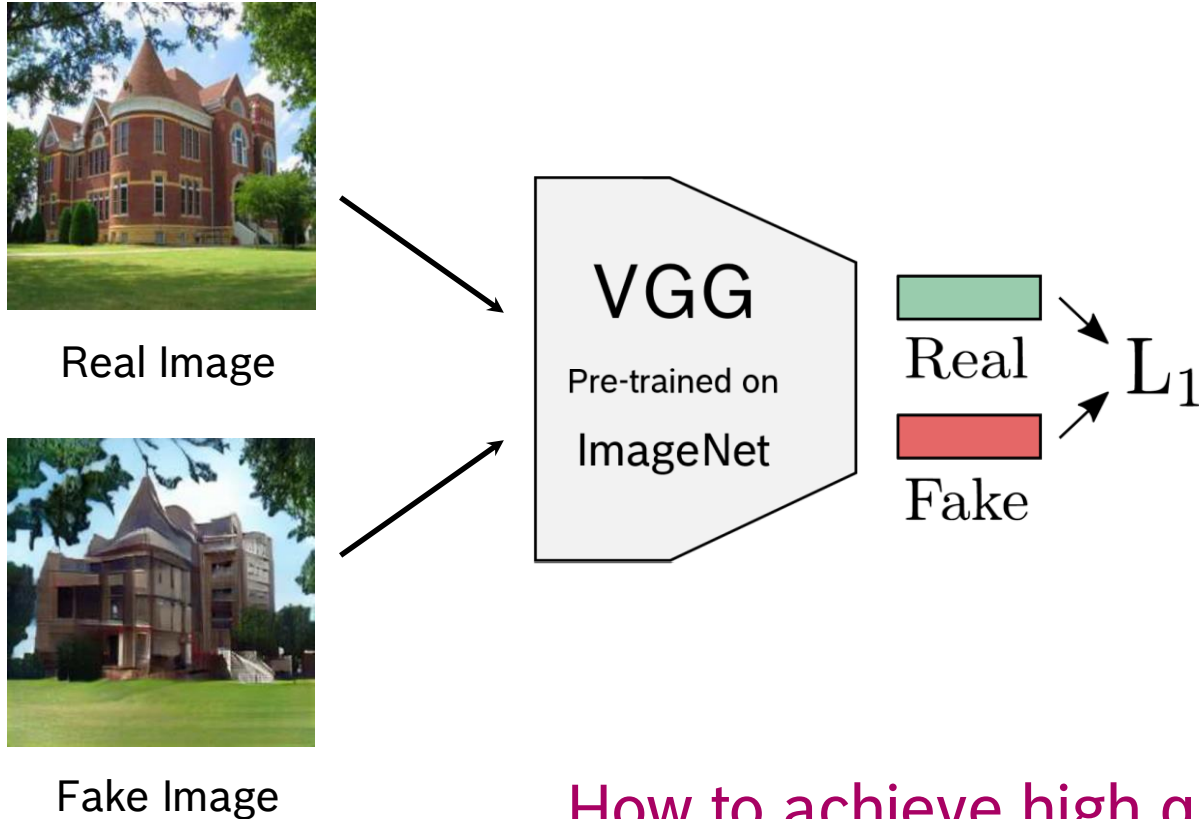Pre-trained on
ImageNet

Real

Fake

$L_1$

Used in:
**CRN** [Chen et. al, 2017]
**Pix2pixHD** [Wang et. al, 2018]
**SPADE** [Park et. al, 2019]
**CC-FPSE** [Liu et. al, 2020]

Label map　　W/o VGG　　With VGG

BOSCH

# Limitations of previous GAN methods
## Perceptual loss



Real Image

Fake Image

Effect of the perceptual loss:
- ▸ Stabilized training
- ▸ Improved quality of images

Drawbacks:
- ▸ Computational overhead
- ▸ Texture and color bias
- ▸ Constrained diversity
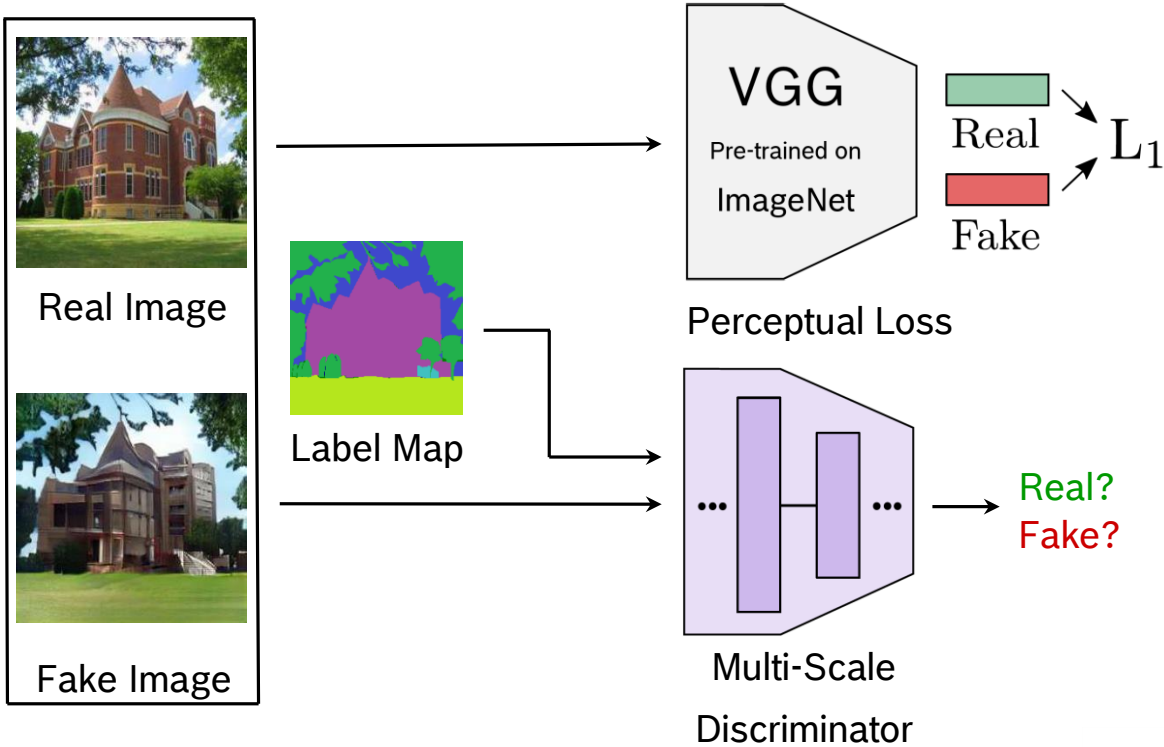
## How to achieve high quality without the perceptual loss?

BOSCH

# OASIS model
## Segmentation-based discriminator

Baseline: SPADE [Park et al., 2019]
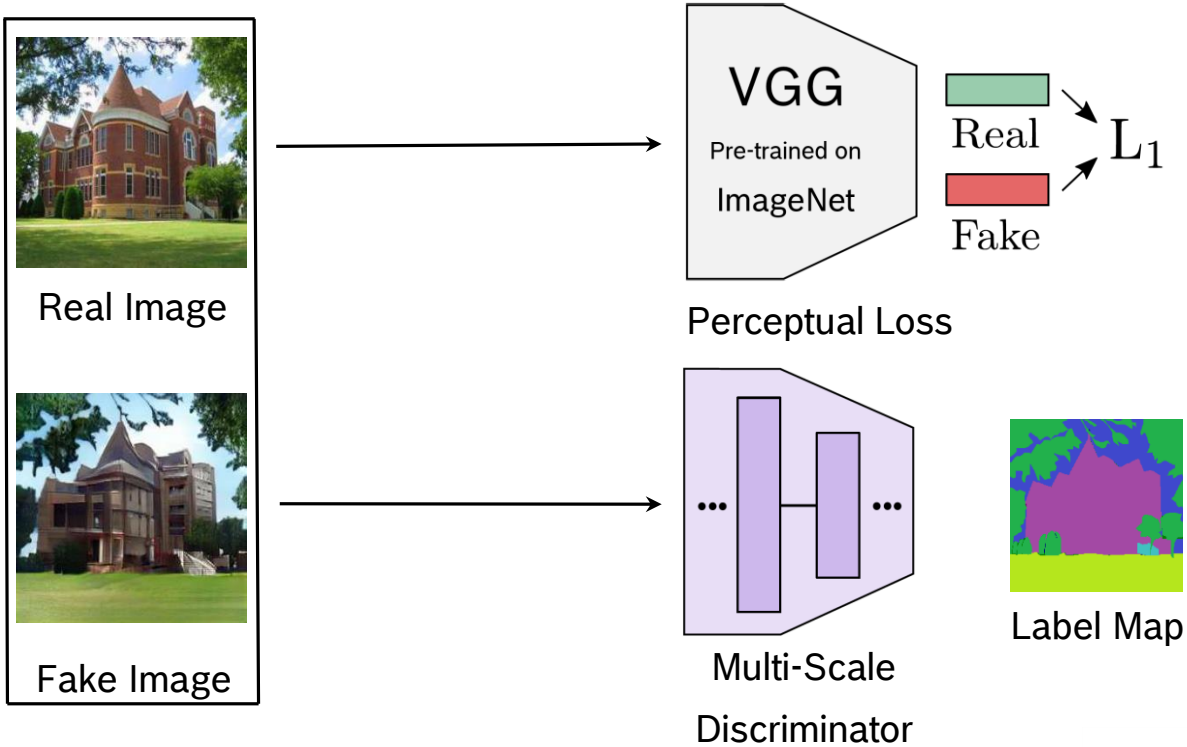
Our solution:
▶ Label map is a target, not input



Real Image

Fake Image

Label Map

VGG
Pre-trained on
ImageNet

Real
Fake
$L_1$

Perceptual Loss

Multi-Scale
Discriminator

Real?
Fake?

BOSCH

# OASIS model
## Segmentation-based discriminator

Baseline: SPADE [Park et al., 2019]

Our solution:

▶ Label map is a target, not input



Real Image

Fake Image

VGG
Pre-trained on
ImageNet

Real
Fake
$L_1$

Perceptual Loss
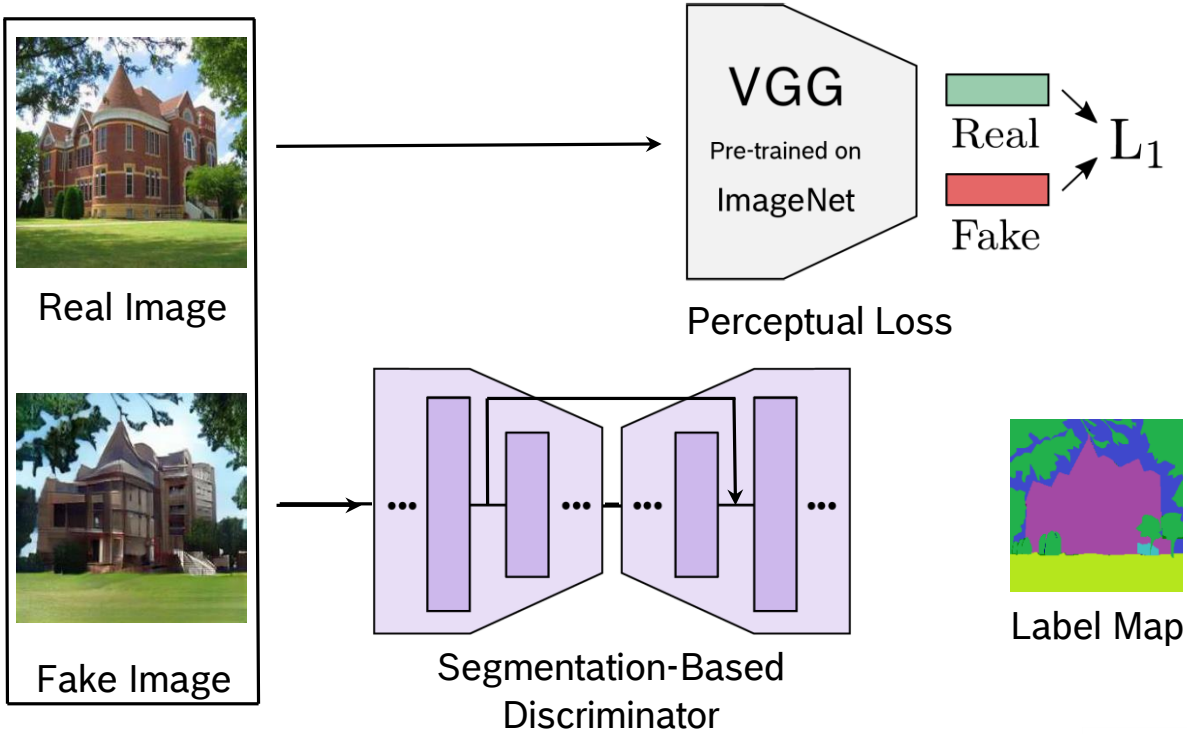
Multi-Scale

Discriminator

Label Map

BOSCH

# OASIS model
## Segmentation-based discriminator

Baseline: SPADE [Park et al., 2019]

Our solution:

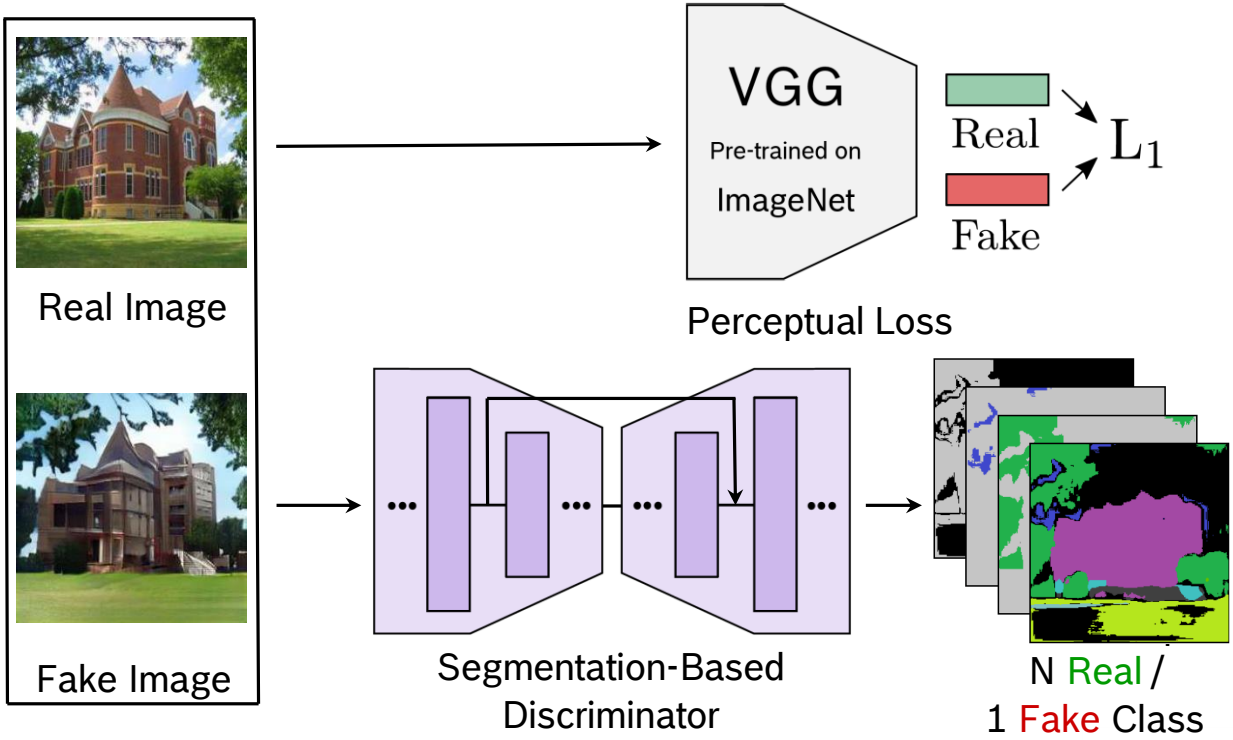▶ Label map is a target, not input

▶ Discriminator's task = segmentation



Real Image

Fake Image

VGG
Pre-trained on
ImageNet

Real

Fake

$L_1$

Perceptual Loss

Segmentation-Based
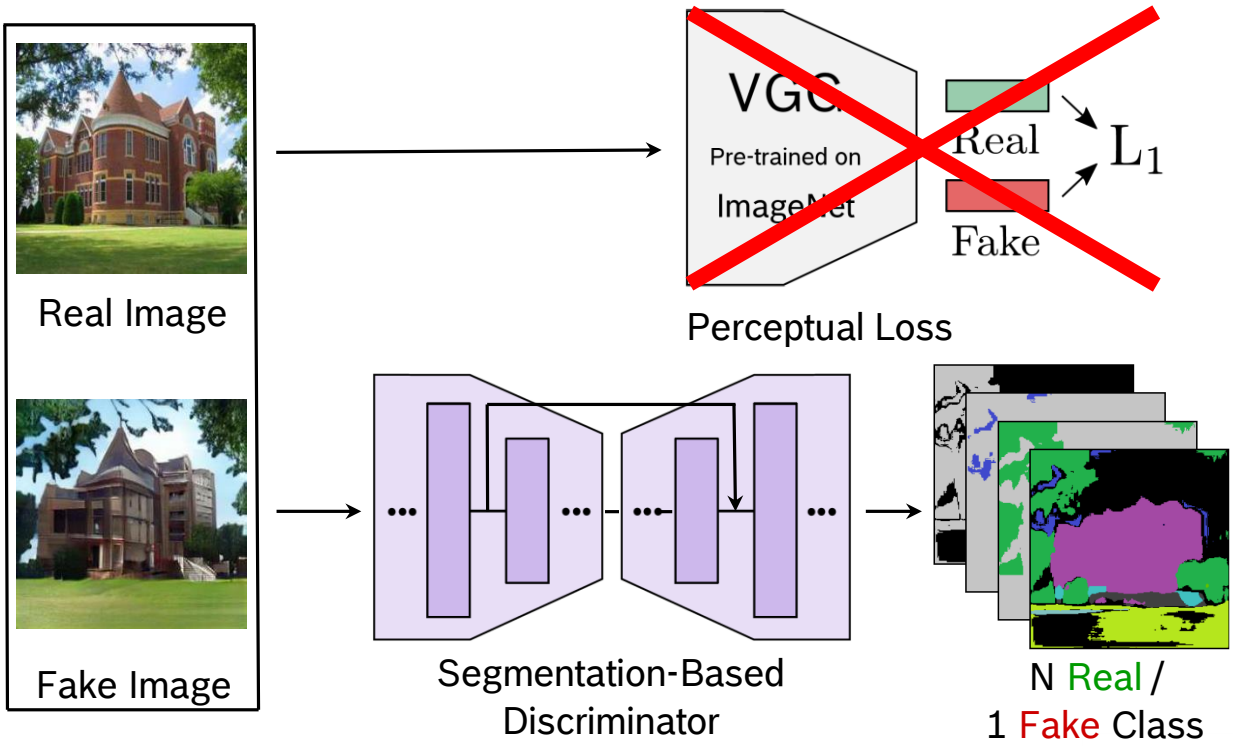Discriminator

Label Map

BOSCH

# OASIS model
## Segmentation-based discriminator

Baseline: SPADE [Park et al., 2019]

Our solution:

▶ Label map is a target, not input

▶ Discriminator's task = segmentation

▶ N+1 loss = adversarial loss



Real Image

Fake Image

VGG Pre-trained on ImageNet

Real Fake

$L_1$

Perceptual Loss

Segmentation-Based Discriminator

N Real / 1 Fake Class

BOSCH

# OASIS model
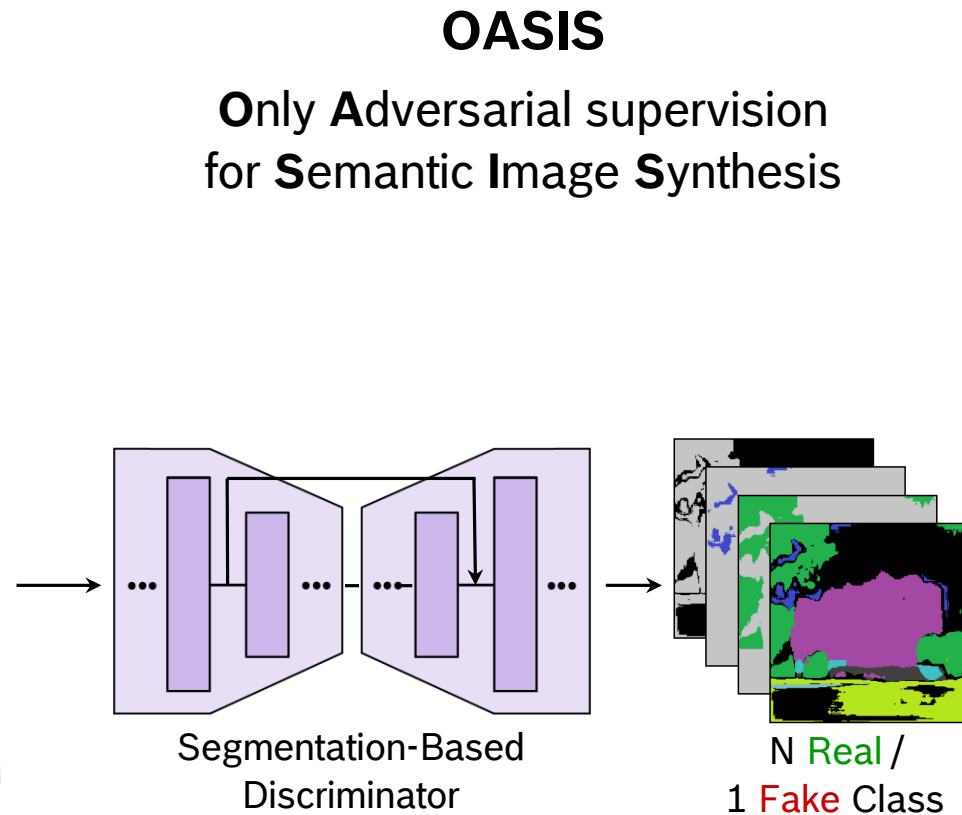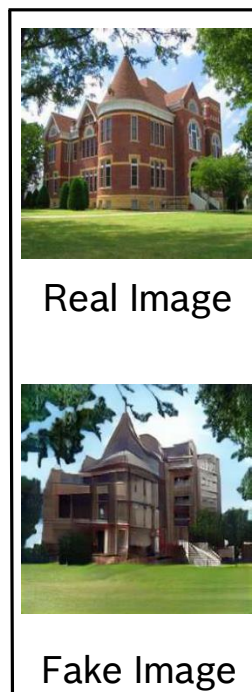## Segmentation-based discriminator

Baseline: SPADE [Park et al., 2019]

Our solution:
- ▶ Label map is a target, not input
- ▶ Discriminator's task = segmentation
- ▶ N+1 loss = adversarial loss

VGG loss becomes unnecessary!



Real Image

Fake Image

Perceptual Loss

Segmentation-Based Discriminator

N Real / 1 Fake Class

BOSCH

# OASIS model
## Segmentation-based discriminator

Baseline: SPADE [Park et al., 2019]

Our solution:

▶ Label map is a target, not input

▶ Discriminator's task = segmentation

▶ N+1 loss = adversarial loss

VGG loss becomes unnecessary!

| $D$ architecture | w/o VGG | | with VGG | |
|---|---|---|---|---|
| | FID↓ | mIoU↑ | FID↓ | mIoU↑ |
| SPADE | 60.7 | 21.0 | 32.9 | 42.5 |
| OASIS | **29.3** | **51.6** | **29.2** | **51.1** |

Real Image

Fake Image

**OASIS**

**O**nly **A**dversarial supervision
for **S**emantic **I**mage **S**ynthesis

Segmentation-Based
Discriminator

N Real /
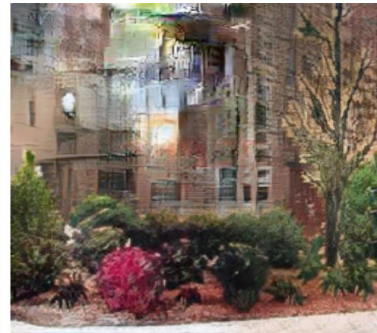1 Fake Class

**BOSCH**

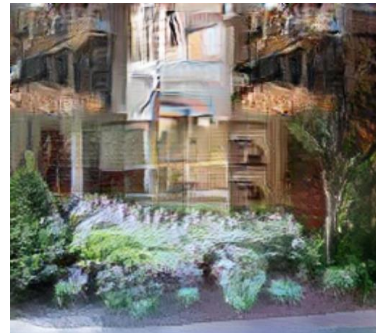# Results
## Comparison to prior art
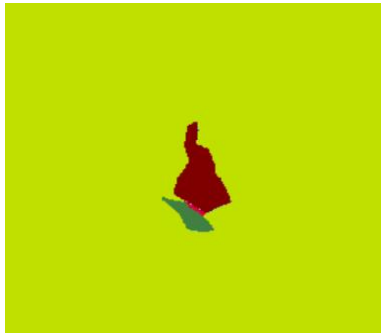


| Label map | Ground truth | SPADE [Park et. al, 2019] | CC-FPSE [Liu et. al, 2020] | OASIS |

BOSCH

# Semantic image synthesis
## Multi-modality



Images from [Park et. al, 2019]

**BOSCH**

# Problems of previous GAN methods
## Limited diversity



Label Map

Noise is ignored!

First reported in:

**Pix2pix** [Isola et. al, 2017]

Noise

Generator

Fake Image

How to achieve high diversity through noise sampling?

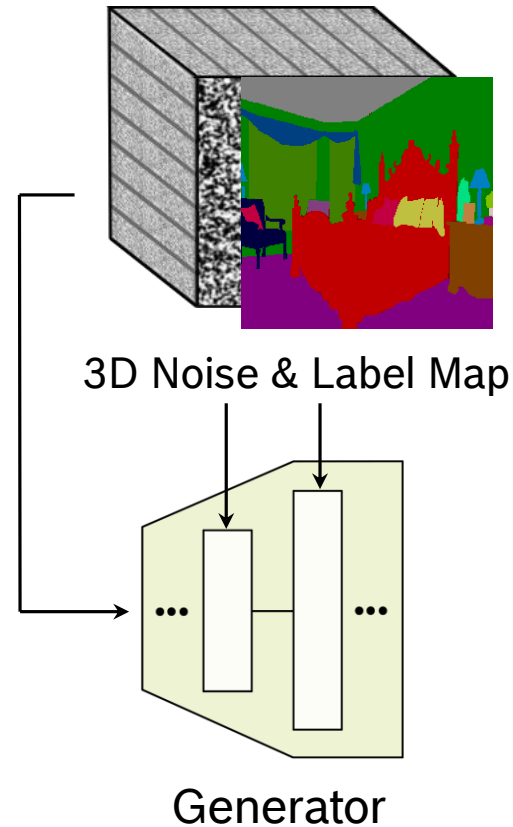**BOSCH**

# OASIS model
## 3D noise injection

**Step 1: Create a composite tensor**

1. Sample a 3D noise tensor

2. Concatenate 3D noise with the (3D) label map

**Step 2: Inject the 3D composite tensor**

1. Input to 1st generator layer

2. Input at *every* generator layer via the *spatially-adaptive* norm ("SPADE" layer)

3D Noise & Label Map

Generator
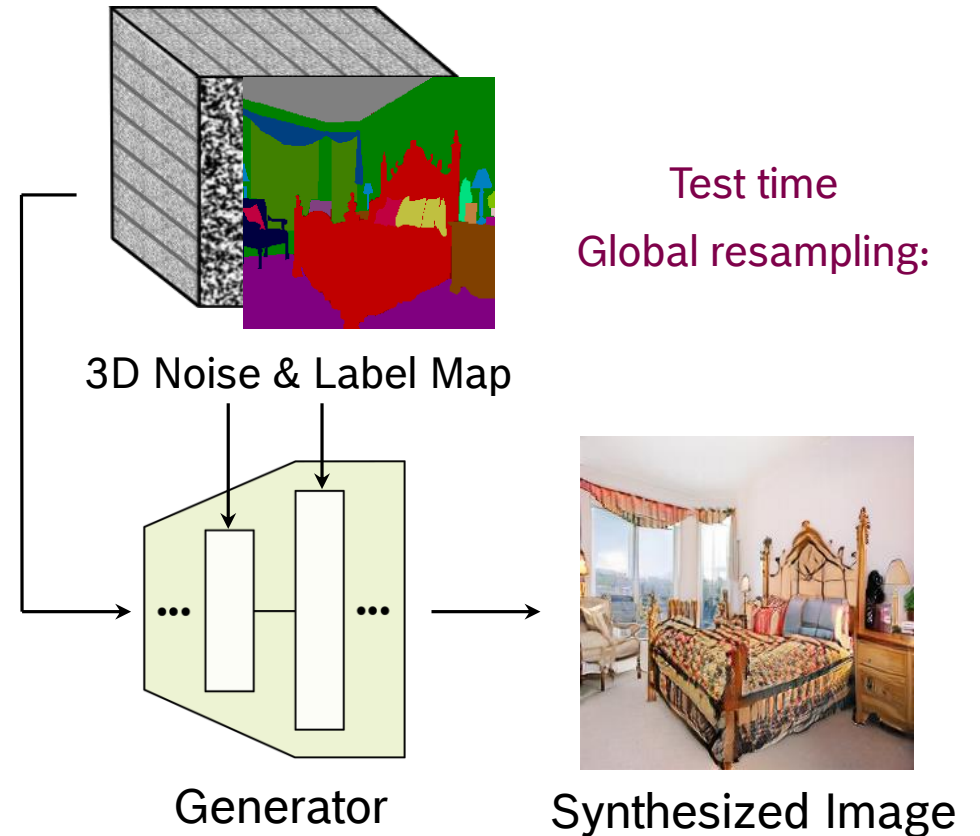
BOSCH

# OASIS model
## 3D noise injection

**Step 1: Create a composite tensor**

1. Sample a 3D noise tensor
2. Concatenate 3D noise with the (3D) label map

**Step 2: Inject the 3D composite tensor**

1. Input to 1st generator layer
2. Input at *every* generator layer via the *spatially-adaptive* norm ("SPADE" layer)

3D Noise & Label Map

Test time
Global resampling:
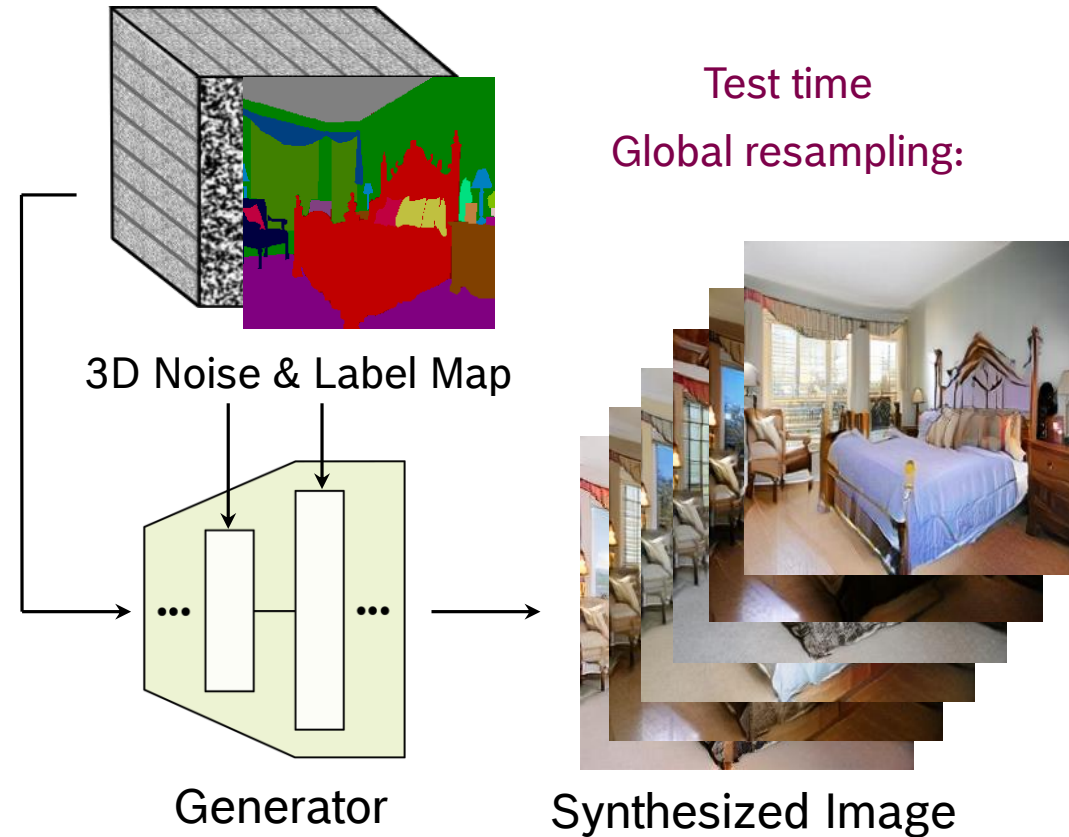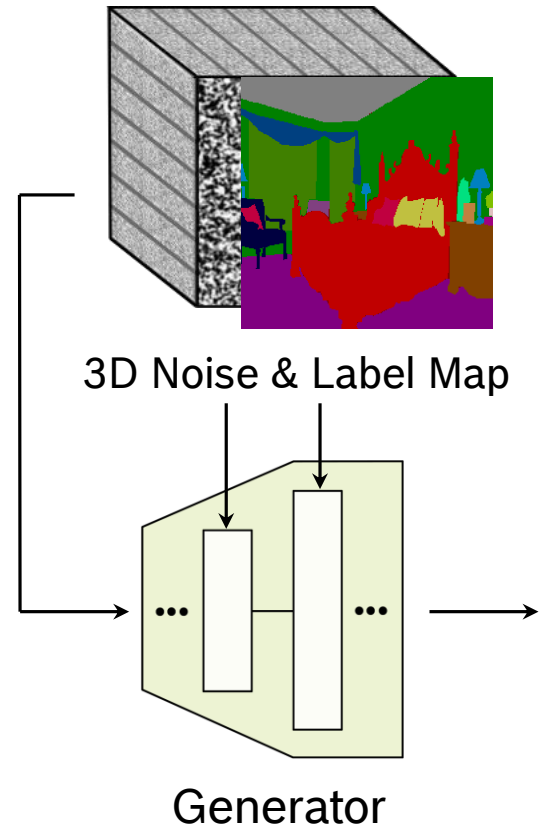
Generator

Synthesized Image

BOSCH

# OASIS model
## 3D noise injection

Step 1: Create a composite tensor

1. Sample a 3D noise tensor
2. Concatenate 3D noise with the (3D) label map

Step 2: Inject the 3D composite tensor

1. Input to 1st generator layer
2. Input at *every* generator layer via the *spatially-adaptive* norm ("SPADE" layer)



3D Noise & Label Map

Generator

Test time

Global resampling:

Synthesized Image

BOSCH

# OASIS model
## 3D noise injection

Local resampling:
Only for bed area

**Step 1: Create a composite tensor**

1. Sample a 3D noise tensor
2. Concatenate 3D noise with the (3D) label map

**Step 2: Inject the 3D composite tensor**

1. Input to 1st generator layer
2. Input at *every* generator layer via the *spatially-adaptive* norm ("SPADE" layer)

3D Noise & Label Map

Generator

Synthesized Image

**BOSCH**

# Results
## Multi-modal generation



Global

Local

BOSCH

# Results
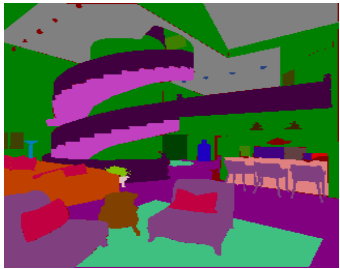## Multi-modal generation

Local
(shape)

BOSCH

# Summary
## Our contributions

1. New state of the art model
2. Segmentation-based discriminator with an N+1 adversarial loss
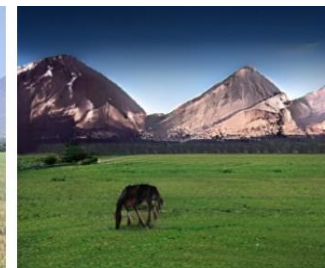3. 3D noise injection scheme

| Label map | With VGG loss | W/o VGG loss | W/o VGG loss, sampled with different 3D noise |
|---|---|---|---|



SPADE [Park et al., 2019]                OASIS (our model)

BOSCH

# Thank you!

**Paper:** https://arxiv.org/abs/2012.04781

**Code:** https://github.com/boschresearch/OASIS

**BOSCH**