# Blending Generative Adversarial Image Synthesis with Rendering for Computer Graphics

**Ekim Yurtsever**      **Dongfang Yang**      **Ibrahim Mert Koc**      **Keith A. Redmill**

The Ohio State University
{yurtsever.2, yang.3455,redmill.1}@osu.edu

## Abstract

Conventional computer graphics pipelines require detailed 3D models, meshes, textures, and rendering engines to generate 2D images from 3D scenes. These processes are labor-intensive. We introduce Hybrid Neural Computer Graphics (HNCG) as an alternative. The contribution is a novel image formation strategy to reduce the 3D model and texture complexity of computer graphics pipelines. Our main idea is straightforward: Given a 3D scene, render only important objects of interest and use generative adversarial processes for synthesizing the rest of the image. To this end, we propose a novel image formation strategy to form 2D semantic images from 3D scenery consisting of simple object models without textures. These semantic images are then converted into photo-realistic RGB images with a state-of-the-art conditional Generative Adversarial Network (cGAN) based image synthesizer trained on real-world data. Meanwhile, objects of interest are rendered using a physics-based graphics engine. This is necessary as we want to have full control over the appearance of objects of interest. Finally, the partially-rendered and cGAN synthesized images are blended with a blending GAN. We show that the proposed framework outperforms conventional rendering with ablation and comparison studies. Semantic retention and Fréchet Inception Distance (FID) measurements were used as the main performance metrics.

## 1   Introduction

The visual fidelity of a conventional computer graphics pipeline depends on the quality of its models, textures, and rendering engine. High-quality 3D models and textures require artisanship, whereas the rendering engine must run complicated physical calculations for the realistic representation of lighting and shading [1]. These processes are labor-intensive. Here we investigate alternatives for alleviating the aforementioned costs.

The alternative to rendering is neural network based generative adversarial image synthesis. The advent of Generative Adversarial Networks (GAN) [2] enabled realization of photo-realistic image synthesis [3–9]. A particular sub-problem, conditional image synthesis [10–14], delves into the more specific task of mapping a pixel-wise semantic layout to a complying photo-realistic image. The conditional semantic layout is the key link between the 3D scene and the generative synthesizer in our framework. More recently, video-to-video synthesis [15] was proposed as an alternative to image synthesis. The temporal dimension was added to the generative process to reduce inconsistencies between synthesized frames.

In this paper, we propose to integrate generative adversarial image synthesis into computer graphics. For each time step, a physics engine determines the semantic layout of the scene with simple 3D models that are radiant with a unique class color without texture. Then, a 2D semantic image is formed with a virtual pinhole camera. This image is the equivalent of a pixel-wise semantic segmentation mask. Next, the GAN-based image synthesizer converts the 2D semantic image to a photo-realistic

Input: 3D Scene

Semantic layout with simple 3D models

Camera pose

Detailed 3D models of a few object-of-interest

Image formation*

Synthesis GAN

~ Generator

Partial Rendering

Blending GAN

Output: 2D Image

*all objects are opaque and radiant with unique class color in the 3D semantic layout scene. No ambient occlusion is considered. Then, a pinhole camera can easily convert the 3D scene into a corresponding upside-down semantic image.

Figure 1: Overview of HNCG. We introduce a novel neural information processing system to form 2D image representations of virtual 3D scenes. Most of the scene is generated with very simple 3D models without texture except for a few partially rendered objects of interest. We blend the cGAN synthesized image with a physics-based partial render for increasing visual fidelity *and* have full control over the appearance of objects of interest.

image. Conditional GAN (cGAN) [16] and CYcle GAN (Cy-GAN) [10] are the main enablers for this step. Simultaneously, a few objects of interest are partially rendered using a conventional pipe. The final frame is a mixture of the generated image with the partial renders. Overview of our approach is shown in Figure 1.

Main contributions:

- The formal introduction of the Hybrid Neural Computer Graphics (HNCG) pipeline
- A novel image formation strategy: Blending generative adversarial image synthesis with physics-based partial rendering
- Reducing the need for texturing and 3D model making artistry in computer graphics pipelines

## 2 Related work

**Rendering.** Physics-based rendering [1] has been used at the end of conventional computer graphics pipelines to form 2D imagery from virtual 3D scenes for a long time. The most common approaches, rasterization and ray-tracing, require a full pipeline of detailed 3D models, their surface textures and materials, and a physics engine such as Unreal Engine 4 [17] to run complicated calculations for representing light and shading. Here, we propose to partially replace this pipe with much simpler 3D models and remove the need for light, texture, and material information for most of the objects in the scene. We also show that visual fidelity can be increased with the proposed method.

**Neural rendering.** Recent work demonstrated that 2D image formation could be achieved given a camera pose and light position in a 3D scene using differentiable convolutional networks [18]. The key enabler here is the formulation of the discrete rasterization problem as a differentiable process [19]. With a differentiable rendering framework, a neural network can be trained with backpropagation. There is more work [20–22] focusing on the different aspects of differentiable rendering formulation and approximations. Neural rendering is an interesting take on an old problem. However, this approach still requires detailed 3D models and is incapable of generating texture information, which reduces the visual fidelity of the output. Here, we propose using generative

adversarial processes and a partial rendering strategy to reduce the complexity of 3D models and textures significantly.

**Generative adversarial image synthesis.** The main difference between generative adversarial image synthesis and neural rendering is the lack of physics. Physical phenomena such as lighting and reflectivity are completely ignored by GAN based neural image synthesizers [3–9]. Instead, the photo-realism is achieved by training the GAN with real-world data. In other words, the network learns to generate photo-realistic images holistically and in an end-to-end fashion. This approach has one major drawback: there is no constraint on the semantic layout of the generated 2D image. Hence, no association with 3D scenery can be constructed. As such, this methodology cannot be applied for our image formation purposes.

**Conditional generative adversarial image synthesis.** On the other hand, conditional GANs [10–13, 15, 14] have been effectively used for image synthesis while retaining a semantic constraint. Typically, this constraint is a pixel-wise semantic segmentation mask, but other modalities such as text [23] have also been used. One limiting factor for cGANs is the paired data requirement. The dataset must contain semantic segmentation masks and the corresponding real-world images together. Building such annotated and paired datasets are labor-intensive.

**Cycle-consistency and domain adaptation.** Cycle consistent GANs and unsupervised domain adaptation techniques remove the paired dataset requirement [24–29, 26]. These works have illustrated that high fidelity image-to-image translation and style transfer can be realized with unpaired data also. Style transfer is very promising and has a huge application range. For example, CyCADA [25] can translate a game engine generated image into a corresponding photo-realistic image. However, the fully-rendered game engine generated image must still *exist* for CyCADA to translate it into the photo-realism domain.

As far as we know, the aforementioned GAN-based image synthesis techniques have not been integrated into computer graphics pipelines until now. Our contribution is novel in this regard. We propose to use simple 3D models radiant with unique class color-codes without textures to form a 2D semantic image first. This image is analogous to a 2D semantic segmentation mask. Then, state-of-the-art GAN-based image synthesizers trained on real-world datasets can be used to generate RGB imagery. We tried both cGAN and Cy-GAN variants. Additionally, we render certain important objects of interest, such as cars in an urban scene, with Unreal Engine 4. The blended image is much more realistic and retains the semantic layout of the scene better.

## 3  Method

### 3.1  Problem formulation

We define a virtual 3D scene $S$ with a 6-tuple $(O_1, O_2, P_1, P_2, T_2, \mathbf{x})$. Where $O_1 = (\mathbf{o}_1, \mathbf{o}_2, \cdots, \mathbf{o}_n)$ is a list of object pose vectors, $\mathbf{o} \in \mathbb{R}^6$, and $P_1 = (M_1, M_2, \cdots, M_n)$ is the list of corresponding simple object meshes. We assume $P_1$ is radiant with unique class color-codes. $O_2$ is a sublist of $O_1$ for certain objects of interest, and it has a corresponding list of more complicated object meshes $P_2$. $P_2$ is not radiant. $T_2$ is a list of texture maps that corresponds to $P_2$. $\mathbf{x} \in \mathbb{R}^6$ is the pose vector of a virtual camera. It should be noted that a corresponding $T_1$ to $O_1$ does not exist.

We follow the formal definition of a triangular mesh given in [30]. $M := (V, Q)$ is a triangular mesh defined with faces $Q \subseteq \{1, \cdots, |V|\}^3$ and vertices $V \subseteq \mathbb{R}^3$. Where $q = (q_1, q_2, q_3) \in Q$ is a triangular face with corresponding vertices $v_{q_1}$, $v_{q_2}$, and $v_{q_3}$. $E(Q)$, edges between the vertices are defined by faces implicitly.

**Problem 1.** Given $S$, we are interested in finding a mapping function $U : \mathbf{x} \rightarrow \mathbb{R}^{H \times W \times 3}$ that will convert the camera pose vector $\mathbf{x}$ to a photo-realistic RGB image with height $H$ and width $W$.

The overview of our solution is shown in Figure 1, and the formal description starts below.

### 3.2  Semantic Image formation

A semantic image formation function $h$ can be obtained with $O_1$, $P_1$ and a pinhole camera model.

Figure 2: Key enablers. $G$ and $F$ are the generators. cGAN can give more realistic results at the cost of a paired dataset. Whereas Cy-GAN completely removes the paired dataset requirement. Details can be found in Appendix A.

$h : \mathbf{x} \to \mathbb{M}^{H \times W}$ maps $\mathbf{x}$ to an integer subspace ($\mathbb{M} \subset \mathbb{Z}$) using the pinhole camera model [31] given below, where $\mathbf{m} \in \mathbb{M}^{H \times W}$ is a pixel-wise semantic image whose entries correspond to the semantic classes of the scene.

$$\begin{pmatrix} m_1 \\ m_2 \end{pmatrix} = -\frac{d}{p_3} \begin{pmatrix} p_1 \\ p_2 \end{pmatrix} \tag{1}$$

Where $(p_1, p_2, p_3)$ is the 3D coordinates of point $\mathbf{p}$ in $\mathbb{R}^3$. $(m_1, m_2)$ is the corresponding pixel coordinates in $\mathbf{m}$. $d$ is the distance between the focal point and image formation plane. $\mathbf{m}$ is an upside-down image as shown in Figure 1. $\mathbf{m}$ is rotated $180°$ for the next step. For simplicity, we use the same notation $\mathbf{m}$ for the rotated image in the remainder of the paper.

Then, the problem narrows down to finding $f : \mathbf{m} \to \mathbb{R}^{H \times W \times 3}$. This is the exact same goal of the well-studied [10–13] conditional image synthesis problem.

### 3.3   cGANs and Cy-GANs: Generative Adversarial Image Synthesis

We propose to use the generator networks of cGANs or Cy-GANs to learn $f : \mathbf{m} \to \mathbb{R}^{H \times W \times 3}$. Training is done on a real-world paired dataset (segmentation mask, real-world image) for cGAN, while Cy-GANs can be trained with an unpaired dataset. Figure 2 gives an overview of enablers. Details of datasets, cGANs, and Cy-GANs are discussed in Appendix A and Section 4.

#### 3.3.1   Baseline: SPADE

The baseline cGAN employed in this study is a SPatially-Adaptive-DE- normalization (SPADE) [12] network, which is a state-of-the-art cGAN based image synthesizer. SPADE outperforms other image-to-image synthesizers by retaining semantic information against conventional normalization operations [12]. This is achieved through the following de-normalization operation where the activation value at layer $i$ is given by:

$$\gamma^i_{c,y,x}(\mathbf{m}) \frac{h^i_{n,c,y,x} - \mu^i_c}{\sigma^i_c} + \beta^i_{c,y,x}(\mathbf{m}). \tag{2}$$

Where $h^i_{n,c,y,x}$ is the activation before normalization, and $\mu^i_c$ and $\sigma^i_c$ are the mean and standard deviation in channel $c$. $\gamma^i_{c,y,x}(\mathbf{m})$ and $\beta^i_{c,y,x}(\mathbf{m})$ are learned varaibles that modulates the normalization process. We refer the readers to the original SPADE paper [12] for more details.

We use a pre-trained SPADE on the Cityscapes dataset [32] as the mapping function $f_s$, and obtain the synthesized image with it $\mathbf{I} = f_s(\mathbf{m})$.

4

### 3.4 Partial rendering

To increase visual fidelity and have full control over certain objects of interest, we propose using physics-based rendering to obtain partially-rendered images $\mathbf{I}_r$. Besides $O_2, P_2, T_2$ and $\mathbf{x}$, a light source is also needed for rendering. Here we assume the properties and location of the light source are fixed and known relative to $\mathbf{x}$. Then the rendering equation [1] can be used to render objects of interest.

$$L_0(\mathbf{p}, \omega, \lambda, t) = L_e(\mathbf{p}, \omega_0, \lambda, t) + \int_\Omega f_r(\mathbf{p}, \omega_i, \omega_0 \lambda, t) L_i(\mathbf{p}, \omega_i, \lambda, t)(\omega_i.\mathbf{n}) d\omega_i \tag{3}$$

Where $L_0(\mathbf{p}, \omega, \lambda, t)$ is the total spectral radiance, $\lambda$ is wavelength, $\omega_0$ is the outgoing light direction, $\omega_i$ is the incoming light direction, $t$ is time and $\mathbf{p}$ is a point in 3D space. $L_e(\mathbf{p}, \omega_0, \lambda, t)$ is the emitted spectral radiance, $\Omega$ is a unit hemisphere with the surface normal center $\mathbf{n}$ of $\mathbf{p}$ and it contains all values for $\omega_i$. $f_r(\mathbf{p}, \omega_i, \omega_0 \lambda, t)$ is the bidirectional reflectance function and finally $L_i(\mathbf{p}, \omega_i, \lambda, t)$ is the spectral radiance of the incoming wavelength.

With Equation 3, the spectral radiance of each 3D point on a few objects of interest is obtained. Then, the partially rendered image $\mathbf{I}_r$ is formed with the same pinhole camera model introduced in Equation 1.

### 3.5 Blending

Here we propose to blend the synthesized image $\mathbf{I}$ with the partially rendered image $\mathbf{I}_r$ to get a hybrid image $\mathbf{I}_h$. The hybrid image is defined as:

$$\mathbf{I}_h := b(\mathbf{I}, \mathbf{I}_r). \tag{4}$$

Where the blending function $b : (\mathbf{I}, \mathbf{I}_r) \to \mathbb{R}^{HxW \times 3}$ maps the synthesized and partially rendered images to a new hybrid RGB image. We compared three different blending functions $b$ in this study.

**Alpha blending.** Taking $\mathbf{I}$ as the background image and $\mathbf{I}_r$ as the foreground image, the alpha blended image $\mathbf{I}_h$ can be obtained with:

$$\mathbf{I}_h = \alpha \mathbf{I} + (1 - \alpha)\mathbf{I}_r \tag{5}$$

**Pyramid blending.** With the gaussian pyramid mask $G_R$ [33], $L_a$ the laplacian pyramid of the foreground $\mathbf{I}_r$, and $L_b$ laplacian pyramid of background $\mathbf{I}$, the laplacian blended pixel $b(i, j)$ can be obtained with:

$$b(i, j) = G_R(i, j)L_a(i, j) + (1 - G_R(i, j))L_b(i, j). \tag{6}$$

**GAN blending.** As a third blending option, we employed GP-GAN [34]. The generator of GP-GAN converts a naive copy-paste blended image to a realistic well-blended image. Besides conditional GAN loss, GP-GAN employs an auxiliary $l_2$ loss to sharpen the image.

$$\mathcal{L}(x, x_g) = \lambda \mathcal{L}_{l_2}(x, x_g) + (1 - \lambda)\mathcal{L}_{\text{adv}}(x, x_g) \tag{7}$$

Where $\mathcal{L}(x, x_g)$ is the final loss, $\mathcal{L}_{l_2}$ is the $l_2$ loss and $\mathcal{L}_{\text{adv}}$ is the adverserial loss. $\lambda$ is a hyperparameter and set to 0.999. Network details are given in Appendix B.

## 4 Experiments

### 4.1 Implementation details

We used a pre-trained SPADE network provided by the original authors [12]. The network was trained on Cityscapes [32], an urban driving dataset with paired semantic mask and image data. CARLA

(a) Semantic image

(c) **Proposed:** (cGAN, partial render) blend

(b) Fully rendered

(d) only-cGAN

Figure 3: The proposed framework converts semantic images (a) to photo-realistic partially rendered images (c) for an interactive virtual environment frame-by-frame. Obtaining (a) is relatively easy in a virtual environment, whereas (b) is more expensive and can yield unrealistic results. For important objects of interest such as lane markings and cars, a pure generative adversarial image-synthesizer (d) cannot be trusted completely. (c) alleviates the shortcomings of (b) and (d).

[35], an open-source driving simulator built upon Unreal Engine 4 was utilized to obtain the semantic layout and partially rendered images. We used the shading and lighting engine [17] of Unreal Engine 4 in our experiments. Only vehicles and lane markings were considered as objects of interest. For blending, we used a GP-GAN [34] trained on the Transient Attributes Database [36].

## 4.2 Evaluation

### 4.2.1 Semantic retention

A common [11, 14, 12] evaluation method for realistic image synthesis is semantic retention analysis. Semantic retention measures the semantic correspondence between the conditional semantic mask and the synthesized image. In summary, an external semantic segmentation network is used to segment the synthesized image. Then, the discrepancy between the conditional semantic layout (input of the synthesizer) and the semantic mask obtained from the generated image (output of the pre-trained external segmentation network) is calculated with top-1 accuracy. A good synthesizer should produce photo-realistic images while retaining the initial conditional semantic layout. In other words, the initial semantic layout is accepted as the ground truth, and the image synthesizer's mask accuracy is calculated to obtain the retention score. Figure 4 illustrates the semantic retention analysis. A higher retention score is favorable.

In this study, we employed DeepLabV3 [37], a state-of-the-art semantic segmentation network, to measure semantic retention. DeepLabV3 was trained on Cityscapes, an urban driving dataset [32].

### 4.2.2 FID

Fréchet Inception Distance (FID) [38] is a commonly used [12, 15] performance metric for measuring visual fidelity. In summary, a deep neural network is employed to extract features of all images in a dataset. Then, the covariance and mean of features obtained from synthesized and real images are compared to get a score. We do not have any real-data corresponding to our virtual 3D scene, but FID can still be used with unpaired data. As such, two different real-world datasets [32, 39] were utilized as the ground truth. The synthesized images were then compared with FID scores. A lower FID indicates high visual fidelity.

Figure 4: An illustration of semantic retention analysis. (a) Full-render yields unrealistic shadows. On the bottom right-hand side of the image, shadows of trees cast on the sidewalk were misclassified as a road by DeepLabV3. (b) Vanilla cGAN vehicles do not retain their shapes perfectly. (c) Blending retains the semantic relationship with the source layout. This figure employs different color codes to distinguish the semantic image formation and semantic segmentation processes for illustration purposes. In practice, the same IDs were used for corresponding classes.

We used two datasets for FID calculations: Cityscape [32] and ADE20K [39]. Details of network architectures, semantic retention, and FID networks are given in Appendix C.

### 4.2.3 Comparisions and ablations

Ablation studies were conducted to demonstrate the effect of each component of the proposed method. Ablation list:

1. No partial-render (only vanilla cGAN or Cy-GAN)
2. No cGAN or Cy-GAN (only full render)
3. Alpha blend (cGAN or Cy-GAN + partial render)
4. Pyramid blend (cGAN or Cy-GAN + partial render)
5. GAN blend (cGAN or Cy-GAN + partial render)

We also compared the performance of 2 alternative generative neural image synthesizers: SPADE [12] and Cycle-GAN (Cy-GAN)[10].

We used vanilla CARLA [35] to obtain fully rendered images of an urban scene. The semantic layout of the scene was also imported from CARLA and used as the conditional input for the generative adversarial image synthesizers. Only vehicles and lane markings were considered by the partial-renders.

In this work, the image synthesis was done frame-by-frame.

### 4.3 Results

#### 4.3.1 Qualitative results

The qualitative results are shown in Figures 3, and 4. These figures illustrate fully rendered, hybrid, and only-cGAN images. As can be seen in Figure 4, rendered shadows are unrealistic, while only-cGAN generated vehicles cannot retain their shapes. These results underline the importance of partial rendering of objects of interest such as cars, vans, and lane markings. The hybrid approach combines the accuracy of a full-render with the realism of a generative model. More qualitative results are given in Appendix D.

Table 1: Performance. Our methods outperform the physics-based computer graphics pipeline.

| Method | Cityscapes [32] | | ADE20K [39] | |
| --- | --- | --- | --- | --- |
| | Sem. ret. ↑ | FID ↓ | Sem. ret ↑ | FID ↓ |
| only render [35] | 0.819 | 231.768 | 0.804 | 361.496 |
| **ours** | | | | |
| only Cycle GAN [10] | 0.343 | 221.609 | 0.430 | 246.597 |
| only cGAN [12] | 0.879 | 208.139 | 0.809 | **240.333** |
| cy-GAN alpha blend | 0.362 | **175.832** | 0.467 | 272.069 |
| cy-GAN pyramid blend | 0.353 | 196.911 | 0.457 | 279.704 |
| cy-GAN GAN blend | 0.318 | 194.191 | 0.454 | 266.615 |
| cGAN alpha blend | **0.879** | 188.809 | 0.801 | 272.877 |
| cGAN pyramid blend | 0.868 | 202.120 | 0.800 | 265.603 |
| cGAN GAN blend | 0.846 | 194.898 | **0.809** | 260.404 |

### 4.3.2 Quantitative results

FID and semantic retention scores are given in Table 1.

- The proposed hybrid blending approach outperformed conventional rendering and pure generative adversarial image synthesis.

- Cityscapes dataset contains only urban driving scenes, while ADE20K has miscellaneous scenes also. All of our virtual 3D scenes were in an urban environment. As such, most of the methods got better FID scores for the Cityscapes dataset.

- GAN blend and Alpha blend showed similar performances. However, it should be noted that the blending GAN was not trained on an urban driving dataset. The blending performance can be possibly increased with a better blending dataset for training the blending GAN.

- The cGAN variants performed better on average as expected, as shown in Table 1. The synthesized images were both realistic and loyal to the initial semantic layout. However, cGAN requires a paired dataset for training.

- The full render was better at semantic retention than Cy-GAN variants. But Cy-GAN variants had a higher FID score than rendering. This means that Cy-GAN can generate realistic images, but fails to retain the semantic constraints.

## 5 Conclusions

This work introduced and investigated the feasibility of Hybrid Neural Computer Graphics (HNCG). Preliminary results indicate that conventional computer graphics pipelines now have a strong alternative.

However, without a paired-dataset to train the cGAN, the proposed system cannot outperform the conventional pipelines. Cycle consistency is not enough by itself, but potential future developments in domain adaptation can benefit HNCG a great deal.

This work focused on frame-by-frame image formation. However, computer graphics applications such as video games and simulations may require temporally more consistent approaches. Future work can focus on video-to-video synthesis to this end.

## 6 Broader Impacts

This work does not reinforce an unfair bias nor has any purpose of harm or injury. On the contrary, it has the potential for increasing the efficiency of artists, video game developers, simulation makers, and other visual content creators.

# References

[1] James T Kajiya. The rendering equation. In *Proceedings of the 13th annual conference on Computer graphics and interactive techniques*, pages 143–150, 1986.

[2] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2018.

[4] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. Self-attention generative adversarial networks. *stat*, 1050:21, 2018.

[5] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 5907–5915, 2017.

[6] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1947–1962, 2018.

[7] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Bolei Zhou, Joshua B. Tenenbaum, William T. Freeman, and Antonio Torralba. Visualizing and understanding generative adversarial networks. In *International Conference on Learning Representations*, 2019.

[8] Ming-Yu Liu, Thomas Breuel, and Jan Kautz. Unsupervised image-to-image translation networks. In *Advances in neural information processing systems*, pages 700–708, 2017.

[9] Takeru Miyato and Masanori Koyama. cGANs with projection discriminator. In *International Conference on Learning Representations*, 2018.

[10] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017.

[11] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1511–1520, 2017.

[12] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2337–2346, 2019.

[13] Xiaojuan Qi, Qifeng Chen, Jiaya Jia, and Vladlen Koltun. Semi-parametric image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8808–8816, 2018.

[14] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Andrew Tao, Jan Kautz, and Bryan 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.

[15] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. Video-to-video synthesis. In *Advances in Neural Information Processing Systems*, pages 1144–1156, 2018.

[16] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[17] Brian Karis and Epic Games. Real shading in unreal engine 4. *Proc. Physically Based Shading Theory Practice*, 4, 2013.

[18] Thu H Nguyen-Phuoc, Chuan Li, Stephen Balaban, and Yongliang Yang. Rendernet: A deep convolutional network for differentiable rendering from 3d shapes. In *Advances in Neural Information Processing Systems*, pages 7891–7901, 2018.

[19] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7708–7717, 2019.

[20] Ioannis Gkioulekas, Anat Levin, and Todd Zickler. An evaluation of computational imaging techniques for heterogeneous inverse scattering. In *European Conference on Computer Vision*, pages 685–701. Springer, 2016.

[21] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3907–3916, 2018.

[22] Matthew M Loper and Michael J Black. Opendr: An approximate differentiable renderer. In *European Conference on Computer Vision*, pages 154–169. Springer, 2014.

[23] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. Generative adversarial text to image synthesis. In *33rd International Conference on Machine Learning*, pages 1060–1069, 2016.

[24] Weijian Deng, Liang Zheng, Qixiang Ye, Guoliang Kang, Yi Yang, and Jianbin Jiao. Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 994–1003, 2018.

[25] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.

[26] Konstantinos Bousmalis, Nathan Silberman, David Dohan, Dumitru Erhan, and Dilip Krishnan. Unsupervised pixel-level domain adaptation with generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3722–3731, 2017.

[27] Ming-Yu Liu and Oncel Tuzel. Coupled generative adversarial networks. In *Advances in neural information processing systems*, pages 469–477, 2016.

[28] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

[29] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7167–7176, 2017.

[30] Mario Botsch, Leif Kobbelt, Mark Pauly, Pierre Alliez, and Bruno Lévy. *Polygon mesh processing*. CRC press, 2010.

[31] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.

[32] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.

[33] Tom Mertens, Jan Kautz, and Frank Van Reeth. Exposure fusion: A simple and practical alternative to high dynamic range photography. In *Computer graphics forum*, volume 28, pages 161–171. Wiley Online Library, 2009.

[34] Huikai Wu, Shuai Zheng, Junge Zhang, and Kaiqi Huang. Gp-gan: Towards realistic high-resolution image blending. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 2487–2495, 2019.

[35] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. *arXiv preprint arXiv:1711.03938*, 2017.

[36] Pierre-Yves Laffont, Zhile Ren, Xiaofeng Tao, Chao Qian, and James Hays. Transient attributes for high-level understanding and editing of outdoor scenes. *ACM Transactions on graphics (TOG)*, 33(4):1–11, 2014.

[37] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.

[38] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.

[39] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017.

[40] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.

[41] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

## Appendix A   Conditional GANs & Cycle GANs

### A.1   Vanilla GAN

Generative Adverserial Networks (GAN) use a generator $G$ model and a discriminator $D$ model in a simultaneous adverserial training strategy. The ultimate goal of $G$ is realistic fake data $\hat{x}$ generation that is indistinguishable from real data $x \in X$. During training, $G$ captures the data distrubition $p_g$. This is achieved with training a generative mapping function $G(z)$ that maps an a priori noise distrubition $p_z(z)$ to the data domain $X$. While $G$ tries to genereate the most realistic $\hat{x}$, the discriminator $D$ tries to discriminate fake data $\hat{x}$ from real data $x$. The output of $D(x)$ is a scalar indicating $x$ being fake or real with a probability. $G(z)$ and $D(x)$, both can be neural networks, are then trained simultaneously with the following min max game:

$$\min_{G}\max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} \left[ \log D(\boldsymbol{x}) \right] + \mathbb{E}_{\boldsymbol{z} \sim p_z(\boldsymbol{z})} \left[ \log(1 - D(G(\boldsymbol{z}))) \right]. \tag{8}$$

### A.2   cGAN

Vanilla GAN learns to generate realistic data, but cannot impose a condition on the fake data. Conditional Generative Adversarial Networks(cGANs) extended the original GAN and can generate realistic fake data while retaining a conditional constraint. This is achieved by pairing the conditional constraint $y$ with the data $x$ and creating a new paired dataset $(x, y)$. This pair can be an (RGB image, pixel-wise semantic layout image), (image, text), and so on. $x$ and $y$ do not have to share the modality. Figure 5 shows an overview of cGAN and details can be found in [16].

$$\min_{G}\max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} \left[ \log D(\boldsymbol{x}|\boldsymbol{y}) \right] + \mathbb{E}_{\boldsymbol{z} \sim p_z(\boldsymbol{z})} \left[ \log(1 - D(G(\boldsymbol{z}|\boldsymbol{y}))) \right]. \tag{9}$$



Figure 5: Overview of cGAN.

### A.3   Cy-GAN

cGAN can successfully generate photo-realistic fake data with a conditional constraint. However, the paired dataset requirement increases the cost of this approach. In comparison, building an unpaired $X$ and $Y$ is relatively easy. Cycle GAN (Cy-GAN) enabled photo-realistic image synthesis with unpaired data. In summary, Cy-GAN contains two generators, $G(x)$ and $F(y)$, which maps $X \rightarrow Y$ and $Y \rightarrow X$ respectively. Also, two discriminators, $D_X$ and $D_Y$, try to distinguish fake data. The adversarial losses are similar to the original GAN; the addition is the novel cycle consistency loss. This loss prevents the mappings of $G$ and $F$ from straying away from each other.

$$\mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = \mathbb{E}_{\boldsymbol{y} \sim p_{\text{data}}(\boldsymbol{y})} \left[ \log D_Y(\boldsymbol{y}) \right] + \mathbb{E}_{\boldsymbol{x} \sim p_x(\boldsymbol{z})} \left[ \log(1 - D_Y(G(\boldsymbol{x}))) \right]. \qquad (10)$$

The key idea of cycle GAN is using the two generators to create a cycle. First, $G(x)$ generates fake $\hat{y}$, then $F(G(x))$ translates the fake $\hat{y}$ back to $\hat{x}$. If the cycle is consistent, then $x \approx \hat{x}$. $x \rightarrow G(x) \rightarrow F(G(x)) \approx x$.

$$\mathcal{L}_{\text{cyc}}(G, F) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})} \left[ \| F(G(x)) - x \|_1 \right] + \mathbb{E}_{\boldsymbol{y} \sim p_{\text{data}}(\boldsymbol{y})} \left[ \| G(F(y)) - y \|_1 \right] \qquad (11)$$

The final loss is:

$$\mathcal{L}(G, F, D_X, D_Y) = \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) + \lambda \mathcal{L}_{\text{cyc}}(G, F) \qquad (12)$$



Figure 6: Cy-GAN

An illustration of Cy-GAN is shown in Figure 6, and details can be found in [28].

## A.4 Hyperparameters and implementation details

We used an Nvidia RTX 2080 to conduct all our computational experiments.

### A.4.1 SPADE architecture and hyperparameters

We used a pre-trained SPADE network provided by the original authors [12]. The network was trained on Cityscapes [32], an urban driving dataset with paired semantic mask and image data. Architecture details are shown in Figure 7. We refer the readers to the original paper [12] for more details.

### A.4.2 Cycle-GAN architecture and hyperparameters

We employed a pre-trained Cycle-GAN provided by the original authors [28]. The network was trained on Cityscapes [32]. The architecture details are given below.

**Generator.** Generator consists of 6 residual blocks. The coded architecture is as follows: c7s1-64, d128, d256, R256, R256, R256, R256, R256, R256, u128, u64, c7s1-3.

Where c7s1-64 is a 7x7 Conv instance normalization ReLU (Conv) layer with 64 filters and stride of 1, d128 is a 3x3 Conv layer with 128 filters, R256 is a residual block with two 3x3 conv layers, u128 is a 3x3 fractional stridden Conv layer with 128 filters and a stride of 0.5.

13

## SPADE Generator



| |
|---|
| Linear(256, 16384) |
| Reshape(1024, 4, 4) |
| SPADE ResBlk(1024), Upsample(2) |
| SPADE ResBlk(1024), Upsample(2) |
| SPADE ResBlk(1024), Upsample(2) |
| SPADE ResBlk(512), Upsample(2) |
| SPADE ResBlk(256), Upsample(2) |
| SPADE ResBlk(128), Upsample(2) |
| SPADE ResBlk(64), Upsample(2) |
| 3x3 Conv-3, Tanh |

## SPADE Discriminator



| |
|---|
| Concat |
| 4x4 ↓ 2-Conv-64 LReLU |
| 4x4 ↓ 2-Conv-128 IN LReLU |
| 4x4 ↓ 2-Conv-256 IN LReLU |
| 4x4 Conv-512 IN LReLU |
| 4x4 Conv-1 |

Fake/real

Figure 7: SPADE [12] network architecture

**Discriminator.** The coded architecture of the discriminator is as follows: C64, C128, C256, C512, where C64 denotes a Conv instance norm. Leaky ReLU layer with 64 filters. The final conv layer outputs a 1-dimensional output.

# Appendix B  Blending

GP-GAN Generator

Foreground (fg)  Background (bg)

Copy and paste fg over bg

encoder

Conv 4x4 64filters

Conv 4x4 128 filters

Conv 4x4 256 filters

Conv 4x4 512 filters

FC(4000)

decoder

Deconv 4x4 512 filters

Deconv 4x4 256 filters

Deconv 4x4 128 filters

Deconv 4x4 64filters

Blended image

GP-GAN Discriminator

Conv 4x4 64filters

Conv 4x4 128 filters

Conv 4x4 256 filters

Conv 4x4 512 filters

Real/fake

Figure 8: Overview of GP-GAN [34] blend.

GP-GAN was trained on the Transient Attributes Database [36]. We used a pre-trained GP-GAN provided by the authors [34]. Network details are shown in Figure 8.

# Appendix C   Semantic retention & FID

## C.1   Semantic retention

In this study, we employed DeepLabV3 [37] to measure semantic retention. Two different pre-trained DeepLabV3 were utilized. The first one was trained on Cityscapes, an urban driving dataset [32], and the second was trained on ADE20K [39]. ADE20K is a more diverse dataset and includes indoor scenes also. Overview of semantic retention is shown in Figure 9.



Figure 9: Overview of semantic retention with Deeplabv3 [37].

## C.2   Fréchet Inception Distance (FID)

FID[38] is commonly used to measure visual fidelity. In summary, a deep neural network is employed to extract features of all images in a dataset. Then, the covariance and mean of features obtained from synthesized and real images are compared to get a score.

An InceptionV3 [40] that was trained on ImageNet [41] was employed as the feature extractor. After features were extracted from the generated images and real-world images from Cityscapes [32] and ADE20K [39], the FID is calculated as follows:

$$d^2 = ||\mu_1 - \mu_2||^2 + Tr(C_1 + C_2 - 2\sqrt{(C_1 C_2)}) \tag{13}$$

Where $\mu_1$ is the mean of features obtained from dataset 1, and $C_1$ is the covariance. The smaller the distance $d^2$ is, the more similar are the two datasets. In other words, a small FID indicates that fake data is close to real-world data.

# Appendix D  More qualitative results



Conventional pipe: render only

Semantic map

cGAN only

cGAN alpha blend

cGAN GP-GAN blend

cGAN Pyramid blend

Cy-GAN only

Cy-GAN alpha blend

Cy-GAN GP-GAN blend

Cy-GAN Pyramid blend

Figure 10: Qualitative results.