Generative Neural Scene Representations for 3D-Aware Image Synthesis

Andreas Geiger

Autonomous Vision Group University of Tübingen and MPI for Intelligent Systems











Covered Papers

GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis

Katja Schwarz and Yiyi Liao and Michael Niemeyer and Andreas Geiger NeurIPS 2020

GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields

Michael Niemeyer, Andreas Geiger CVPR 2021

CAMPARI: Camera-Aware Decomposed Generative Neural Radiance Fields

Michael Niemeyer, Andreas Geiger 3DV 2021

Collaborators



Michael Niemeyer

Katja Schwarz



Yiyi Liao

Generative models are fantastic!

Sample a latent code from the prior distribution.



Latent Code

Pass latent code to trained generator G_{θ} .



Latent Code

Generator G_{θ}

The generator outputs a synthesized image.





Latent Code

Generator G_{θ}

Generated Image*

* The generated images are samples from StyleGAN2.

Sample more latent codes to get different generated images.





Generated Image*

* The generated images are samples from StyleGAN2.

Sample more latent codes to get different generated images.





Generated Image*

* The generated images are samples from StyleGAN2.

Is generating photorealistic images all we need?

For many applications, we require **control over the image generation process:**

For many applications, we require **control over the image generation process:**



Animation Movies

Video Source: Disney's Toy Story 4 Trailer

For many applications, we require **control over the image generation process:**



Video Games

Video Source: Gran Turismo 7 Trailer

For many applications, we require **control over the image generation process:**

Virtual Reality



Video Source: Oculus Rift Trailer

Goal: A generative model for 3D-aware image synthesis which allows us to:

► Generate photorealistic images

- ► Generate photorealistic images
- Control individual objects wrt. their appearance, pose and size in 3D

- ► Generate photorealistic images
- ► Control individual objects wrt. their appearance, pose and size in 3D
- ► Control the camera viewpoint in 3D

- ► Generate photorealistic images
- ► Control individual objects wrt. their appearance, pose and size in 3D
- ► Control the camera viewpoint in 3D
- ► Train from unstructured and unposed images

What representation should we use for 3D-aware image synthesis?

Voxel-based 3D Shape with Volumetric Rendering



PlatonicGAN [Henzler et al., ICCV 2019]

Voxel-based 3D Shape with Volumetric Rendering



PlatonicGAN [Henzler et al., ICCV 2019]

+ Multi-view consistent

Schwarz, Liao, Niemeyer, Geiger: GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis. NeurIPS, 2020.

Voxel-based 3D Shape with Volumetric Rendering



PlatonicGAN [Henzler et al., ICCV 2019]

- Multi-view consistent
- Low image fidelity, high memory consumption

Voxel-based 3D Latent Feature with Learnable Projection



HoloGAN [Nguyen-Phuoc et al., ICCV 2019]

Voxel-based 3D Latent Feature with Learnable Projection



HoloGAN [Nguyen-Phuoc et al., ICCV 2019]

+ High image fidelity

Schwarz, Liao, Niemeyer, Geiger: GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis. NeurIPS, 2020.

Voxel-based 3D Latent Feature with Learnable Projection



HoloGAN [Nguyen-Phuoc et al., ICCV 2019]

- ✤ High image fidelity
- Object shape and identity vary with viewpoint due to learnable projection

Generative Radiance Fields



Generative Radiance Fields



+ Continuous representation, multi-view consistent

Generative Radiance Fields



- + Continuous representation, multi-view consistent
- + High image fidelity, low memory consumption

Sample camera matrix **K**, camera pose $\boldsymbol{\xi} \sim p_{\boldsymbol{\xi}}$, and patch sampling pattern $\boldsymbol{\nu} \sim p_{\boldsymbol{\nu}}$.

 \mathbf{K} $\boldsymbol{\xi} \sim p_{\boldsymbol{\xi}}$

 $\boldsymbol{\nu} \sim p_{\nu}$

Pass **K**, $\boldsymbol{\xi}$, and $\boldsymbol{\nu}$ to generator G_{θ} which samples corresponding pixels / rays.



For each ray, get viewing direction \mathbf{d}_r and sample 3D points \mathbf{x}_r^i along ray.



Schwarz, Liao, Niemeyer, Geiger: GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis. NeurIPS, 2020.

Push \mathbf{d}_r and \mathbf{x}_r^i through positional encoding γ as input to conditional radiance field g_{θ} .



Sample latent shape and appearance codes $\mathbf{z}_s, \mathbf{z}_a$ and pass them as condition to g_{θ} .



Perform volume-rendering for each ray, resulting in predicted patch P'.



Sample patch **P** from real image **I** drawn from the data distribution $p_{\mathcal{D}}$ with pattern $\boldsymbol{\nu}$.


GRAF: Generative Radiance Fields

Pass fake and real patch \mathbf{P}', \mathbf{P} to discriminator D_{ϕ} and train with adversarial loss.



GRAF: Generative Radiance Fields



- ► Key idea: Generator/discriminator for image patches of size 32 × 32 pixels
- ► Patches sampled at **random locations** and with **random scales** (using dilation)

Conditional Radiance Fields



Conditional Radiance Field Architecture:

- Fully-connected MLP g_{θ} with ReLU activations
- Appearance code \mathbf{z}_a and view direction influence only the output color
- ► Inductive biases allow the model to disentangle shape and appearance

Results

Generative Radiance Fields

Results on synthetic Carla dataset at 256^2 pixels:





Schwarz, Liao, Niemeyer, Geiger: GRAF: Generative Radiance Fields for 3D-Aware Image Synthesis. NeurIPS, 2020.

Generative Radiance Fields

Results on real CelebA-HQ dataset at 256^2 pixels:



How can we scale this idea to more complex, multi-object scenes?

GIRAFFE: Compositional Generative Neural Feature Fields

GRAF:

► Incorporate a **3D representation** into the generative model

GIRAFFE: Compositional Generative Neural Feature Fields

GRAF:

► Incorporate a **3D representation** into the generative model

GIRAFFE:

► Incorporate a **compositional 3D scene representation** into the generative model

GIRAFFE: Compositional Generative Neural Feature Fields

GRAF:

► Incorporate a **3D representation** into the generative model

GIRAFFE:

- ► Incorporate a **compositional 3D scene representation** into the generative model
- ► Incorporate a **neural renderer** to yield fast and high-quality inference

Sample N shape and appearance codes.





Get N feature fields. Note: We show features in RGB color for clarity.



Sample size and pose for each feature field.



Get posed feature fields.



Composite all feature feature fields into a single 3D scene representation.



Sample a camera pose.



Niemeyer, Geiger: GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR, 2021.

Perform volume rendering and get low-resolution feature image.



Niemeyer, Geiger: GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR, 2021.

Pass feature image to neural renderer to obtain final high-resolution image.



At test time, we sample individual codes to control shapes, appearances and poses.



Niemeyer, Geiger: GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR, 2021.



- ▶ We volume-render the feature image at 16² pixels and upsample via the decoder
- ► Unlike GRAF, we can directly train with an adversarial loss at full resolution

We use GRAF's architecture, but replace the color head with a feature head:



We use GRAF's architecture, but replace the color head with a feature head:





Scene Composition

We have N feature fields

$$h_i(\mathbf{x}, \mathbf{d}) = (\sigma_i, \mathbf{f}_i)$$

which predict a density σ_i and a feature vector \mathbf{f}_i for any (\mathbf{x}, \mathbf{d}) .

We combine densities by summation and features by density weighted averaging:

$$\sigma = \sum_{i=1}^{N} \sigma_i$$
 $\mathbf{f} = \frac{1}{\sigma} \sum_{i=1}^{N} \sigma_i \mathbf{f}_i$

Niemeyer, Geiger: GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR, 2021.

Results

We compare object translation for a 2D-based GAN (left) and our method (right):



We can perform more complex operations like circular translations



We can add more objects at test time (trained on two-object)



We can rotate the object



Niemeyer, Geiger: GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. CVPR, 2021.

We can translate the object



We can change the object shape



We can change the object appearance



We can generate out-of-distribution samples



(a) Increase Depth Translation.



(b) Increase Horizontal Translation.



(c) Add Additional Objects (Trained on Two-Object Scenes).



(d) Add Additional Objects (Trained on Single-Object Scenes).



Rendering Time

	64×64	256×256
GRAF	110.1ms	1595.0ms
GIRAFFE	4.8ms	5.9ms

- We volume-render the feature image at 16×16 pixels.
- CNN-based neural rendering and upsampling yields fast inference.

Can we learn the camera distribution?

CAMPARI: Camera-Aware Generative Neural Radiance Fields

GRAF, GIRAFFE:

- ► Learn a 3D-aware image generator from unposed image collections
- ► Requires careful tuning of camera pose distributions

CAMPARI: Camera-Aware Generative Neural Radiance Fields

GRAF, GIRAFFE:

- ► Learn a 3D-aware image generator from unposed image collections
- ► Requires careful tuning of camera pose distributions

CAMPARI:

► Learn a 3D aware image generator and a **camera generator** jointly.
Sample prior camera $\boldsymbol{\xi}^{\text{prior}} \sim p_{\boldsymbol{\xi}}$. Parameters: intrinsics f_x, f_y and extrinsics r, α_r, α_e .



Pass $\boldsymbol{\xi}^{\text{prior}}$ to camera generator G_{θ}^{C} and obtain predicted camera $\boldsymbol{\xi}^{\text{pred}}$.



Camera Generator

Niemeyer, Geiger: CAMPARI: Camera-Aware Decomposed Generative Neural Radiance Fields. 3DV, 2021.

Pass $\pmb{\xi}^{\text{pred}}$ and sampled FG / BG latent codes to 3D-aware image generator.



Train entire method with adversarial objective similar to GRAF, GIRAFFE.



Results

CAMPARI learns to match the GT distribution for synthetic datasets.



CAMPARI learns to match the GT distribution for synthetic datasets.



This leads to significantly improved performance compared to untuned models.



This leads to significantly improved performance compared to untuned models.



Camera Rotation

Niemeyer, Geiger: CAMPARI: Camera-Aware Decomposed Generative Neural Radiance Fields. 3DV, 2021.

Generative Neural Scene Representations

Summary

- ► We propose novel methods for **3D controllable image synthesis**
- ► We demonstrate training from **raw**, **unposed image collections**
- ► We incorporate **compositional 3D scene structure** into the generative model
- ► We gain explicit control over individual objects during synthesis

Future Research

- Scale to larger and more complex multi-object scenes
- Enhance **photorealism** to match 2D GANs (e.g., StyleGAN2)
- ► Disentangle lighting, materials, etc.

Thank you!

http://autonomousvision.github.io

