3D Controllable Image Synthesis

Yiyi Liao

December 24, 2020



University of Tübingen MPI for Intelligent Systems

Autonomous Vision Group



Collaborators



Katja Schwarz



Lars Mescheder



Michael Niemeyer



Jun Xie



Andreas Geiger





Introducing Town10 A vibrant city with new assets to improve visual quality H

Can we learn a simulation tool from 2D images?

Liao*, Schwarz*, Mescheder, Geiger: Towards Unsupervised Learning of Generative Models for 3D Controllable Image Synthesis. CVPR, 2020

Goal

3D Controllable Image Synthesis

- Learning a generative image model
- With controllable 3D factors:
 - Object shape
 - Object appearance
 - Object pose
 - Camera viewpoint
- Learn from 2D observations (unposed images)
- No supervision (segmentation, bounding box, depth)









Classical Rendering Pipeline



3D factors can be controlled Expensive 3D content creation

2D Generative Models



Unsupervised learning from 2D images 3D factors cannot be controlled

Our approach



3D factors can be controlled Unsupervised learning from 2D images



Our approach



Idea: Learning the image generation process jointly in 3D and 2D space





- Foreground/background primitives: $\{\mathbf{o}_1, \dots, \mathbf{o}_N, \mathbf{o}_{bg}\}$, $\mathbf{o}_i = (\mathbf{R}_i, \mathbf{t}_i, \mathbf{s}_i, \phi_i)$
- Primitive type: point cloud, sphere, cuboid



- Sample camera viewpoint, render each primitive individually
- Obtain feature map X, alpha map A and depth map D



- Convert features to RGB pixel values
- Render to image via alpha composition (based on depth ordering)

Loss Functions



 $\mathcal{L}_{adversarial}(\theta, \psi, c) = \mathbb{E}_{p(\mathbf{z})}[f(d_{\psi}(g_{\theta}(\mathbf{z}, c), c))] + \mathbb{E}_{p_{\mathcal{D}}(\mathbf{I}|c)}[f(-d_{\psi}(\mathbf{I}, c))]$

Loss Functions



$$\mathcal{L}_{compactness}(\theta) = \mathbb{E}_{p(\mathbf{z})} \left[\sum_{i=1}^{N} \max\left(\tau, \frac{1}{P}\right) \right]$$



Loss Functions



$$\mathcal{L}_{geometric}(\theta) = \mathbb{E}_{p(\mathbf{z})} \left[\sum_{i=1}^{N} \|\mathbf{A}'_{i} \odot (\mathbf{X}'_{i} - \tilde{\mathbf{X}}'_{i})\|_{1} \right] + \mathbb{E}_{p(\mathbf{z})}$$



Datasets









Baselines

Layout2lm [Zhao et al., CVPR 2019]

- Only 2D translation control, requires
 2D bounding box supervision
- Fails to disentangle object identity and pose











Baselines

Ours 2D Replace 3D primitives with 2D primitives











Baselines

Ours 2D Replace 3D primitives with 2D primitives

• Only 2D translation control









Baselines

Ours 2D

Replace 3D primitives with 2D primitives

- Only 2D translation control
- Fails to disentangle rotation and translation









Ours: Object Translation









Ours: Object Rotation









Ours: Camera Rotation











Ours: Camera Rotation











Ours: Camera Rotation











Ours: Object Translation











Is there a better 3D representation for 3D controllable image synthesis?

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

3D Representations

3D Latent Feature with Learnable Projection



HoloGAN [Nguyen et al., ICCV 2019]

+ High image fidelity

Object identity may vary with viewpoint due to learnable projection

3D Representations

3D Shape with Volumetric Rendering



PlatonicGAN [Henzler et al., ICCV 2019]

Multi-view consistent

- Low image fidelity, high memory consumption

3D Representations

Generative Radiance Fields



Radiance Field

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

Comparison to Baselines

PlatonicGAN 🛎 🔋 🛸 🍉 HoloGAN 🚔 🌰 🚔 🍩 Ours 斗 📥 🦛 📥







Scalability to High Resolution











Disentangling Shape & Appearance







Real-World Datasets





Training Images







Training Images





 256×256

Failure Cases

White artifacts



Inward facing depth





More complex real world

• Incorporate more supervision • Object disentanglement

What's next?





RGB



Semantic





Bounding Box

Instance



Semantic





Confidence



Instance



Bounding Box



• Simulation

- 3D Controllable image synthesis
- Novel view & semantic synthesis

• Perception

- Semantic & instance segmentation in 2D & 3D
- Holistic scene understanding

• Robotics

Semantic SLAM

Thank you!

Live slides with videos: yiyiliao.github.io/20201224_GAMES

