

Visual Object Networks: Image Generation with Disentangled 3D Representations

Jun-Yan Zhu¹

Motivation

Problem: Existing 2D generative models

Do not capture the 3D nature of the world. **Do not allow 3D-aware image manipulations.**

Given an image of a car:



- What if we apply its texture to a van?
- Can we mix different 3D designs?

Goal: Joint 3D & 2D generation with disentangled representation.





VON (ours)

Chengkai Zhang¹ Jiajun Wu¹ Joshua B. Tenenbaum¹ Antonio Torralba¹ 2 Google Research 1 Massachusetts Institute of Technology

Visual Object Networks (VON)

shape network G_{shape}



Training data: Unpaired 3D shapes & Image datasets



Zhoutong Zhang¹







Total Loss:

3D Shape Loss:

Texture Loss:

Formulation

Image Formation: $x = G_{\text{texture}}(\mathcal{P}(G_{\text{shape}}(\mathbf{z}_{\text{shape}}), \mathbf{z}_{\text{view}}), \mathbf{z}_{\text{texture}})$

 $\mathcal{L} = \lambda_{\text{shape}} \mathcal{L}_{\text{shape}} + \mathcal{L}_{\text{texture}}$

Adversarial losses

Image Generation (2D GANs vs. VON)

VON (ours)





William T. Freeman^{1,2}



Shape Generation

| | 3D-GAN (voxels) | VON (voxels) |
|--------|-----------------|--------------|
| Cars | 3.021 | 0.021 |
| Chairs | 2.598 | 0.082 |
| | 3D-GAN (DF) | VON (DF) |
| Cars | 3.896 | 0.002 |
| Chairs | 1.790 | 0.006 |