

3D-aware Image Generation using 2D Diffusion Models

Jianfeng Xiang^{1,2} Jiaolong Yang² Binbin Huang³ Xin Tong² ¹Tsinghua University ²Microsoft Research Asia ³ShanghaiTech University

Overview

Background

Learning to generate 3D contents has become an increasingly prominent task due to its numerous applications.

Current Limitations

Predominant method rely on GANs for generative modeling.

- Limited generative power
- \circ Hard to scale up

Our Solution

2D diffusion-based paradigm with sequential unconditionalconditional multi-view image generation formulation.



NeRF-GAN paradigm



IVID paradigm

- Superior generative power
- Scaling up to ImageNet
- o Only using still images with depth from depth estimators to construct the training data

We undertake 3D-aware generation on ImageNet. We show the capability of our method for large-angle generation from unaligned data (up to 360 degrees).





Problem Formulation

- o Factorize it into the production of an unconditional distribution and a series of conditional distributions.
- o Approximate the conditional distributions by changing conditions to **warped images with holes** with **depth** information.





Training

Erosion

Inference

Orig

By sampling views iteratively following the unconditional-conditional image generation process, the final 3D assets are obtained. Aggregated conditioning is applied to ensure the novel view is conditioned on all previous sampled views.

Fast free-view synthesis

Warp

Once a set of views uniformly covering the desired viewing range are generated, **fusion**based free-view synthesis can be applied to interpolate viewpoint in real-time.



Our diffusion-based 3D-aware image generation trained on ImageNet

Methodology

o The distribution of 3D assets is equivalent to the joint distribution of its multi-view images. $= q_i(\Gamma(\mathbf{x}, \boldsymbol{\pi}_0)) \cdot$ $q_i(\Gamma(\mathbf{x}, \pi_1) | \Gamma(\mathbf{x}, \pi_0))$.

 $q_i(\Gamma(\mathbf{x}, \boldsymbol{\pi}_N) | \Gamma(\mathbf{x}, \boldsymbol{\pi}_0), \cdots, \Gamma(\mathbf{x}, \boldsymbol{\pi}_{N-1}))$ $q_a(\mathbf{x}) \approx q_i(\Gamma(\mathbf{x}, \boldsymbol{\pi}_0))$ $q_i(\Gamma(\mathbf{x}, \boldsymbol{\pi}_1) | \Pi(\Gamma(\mathbf{x}, \boldsymbol{\pi}_0), \boldsymbol{\pi}_1))$

 $q_i(\Gamma(\mathbf{x}, \boldsymbol{\pi}_N) | \Pi(\Gamma(\mathbf{x}, \boldsymbol{\pi}_0), \boldsymbol{\pi}_N), \cdots)$

Data Preparation

We construct training data using depth-based image warping. By forward-backward warping the images with estimated depth, we can construct training data without actual multi-view pairs.

Two DMs are trained to fit the unconditional and conditional distributions, respectively.

Blur and texture erosion augmentations are applied to the warped images to boost performance.

9.45

40.4

67.3

Large View Synthesis (Up to 360°)





Ablation Study on Augmentations

Geometry Visua



Project page: https://jeffreyxiang.github.io/ivid/













