

# Do 2D GANs Know 3D Shape? Unsupervised 3D Shape Reconstruction from 2D Image GANs

ICLR2021 Oral paper

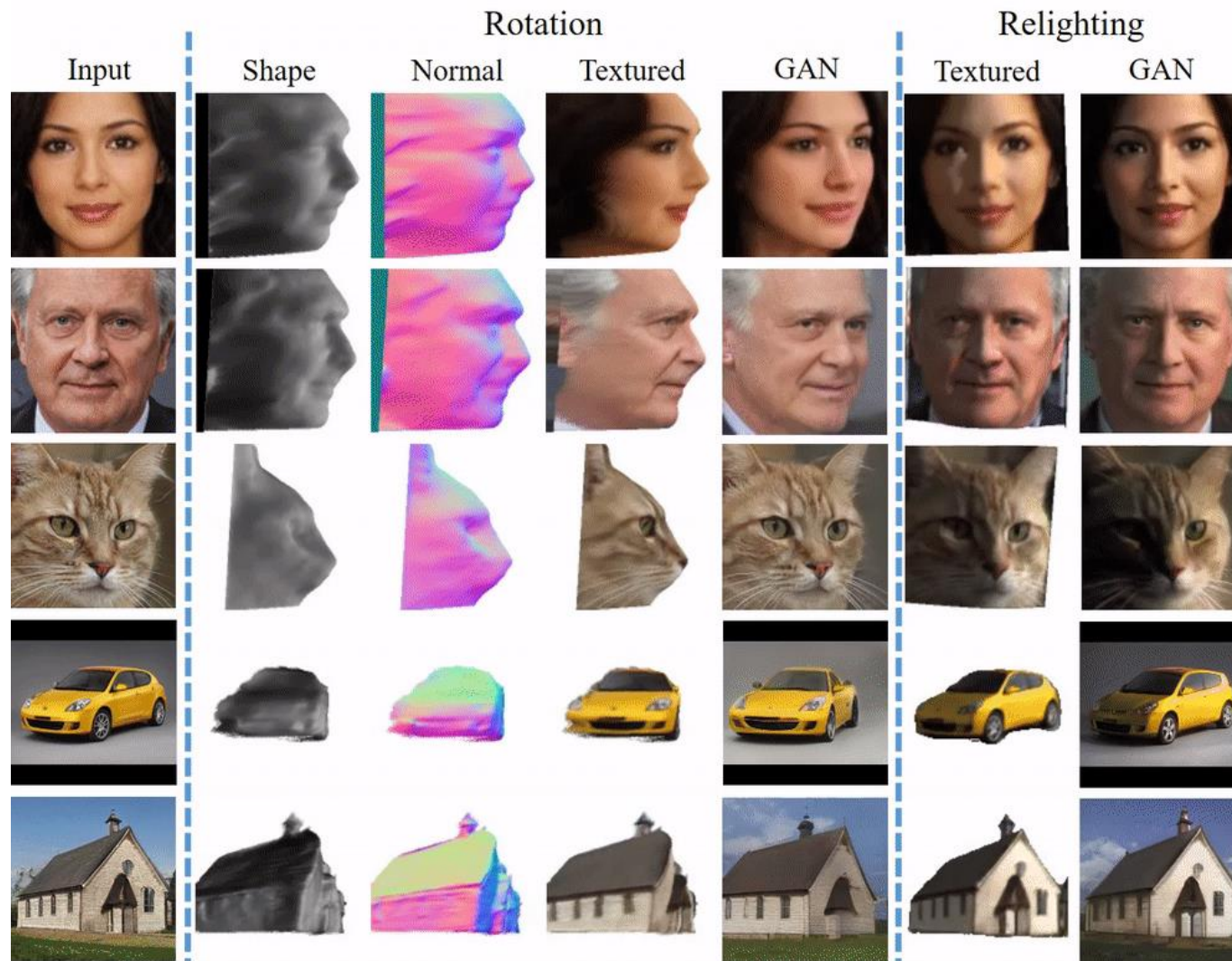
21.04.12 Leeminsoo



**DAVIAN**  
Data and Visual Analytics Lab

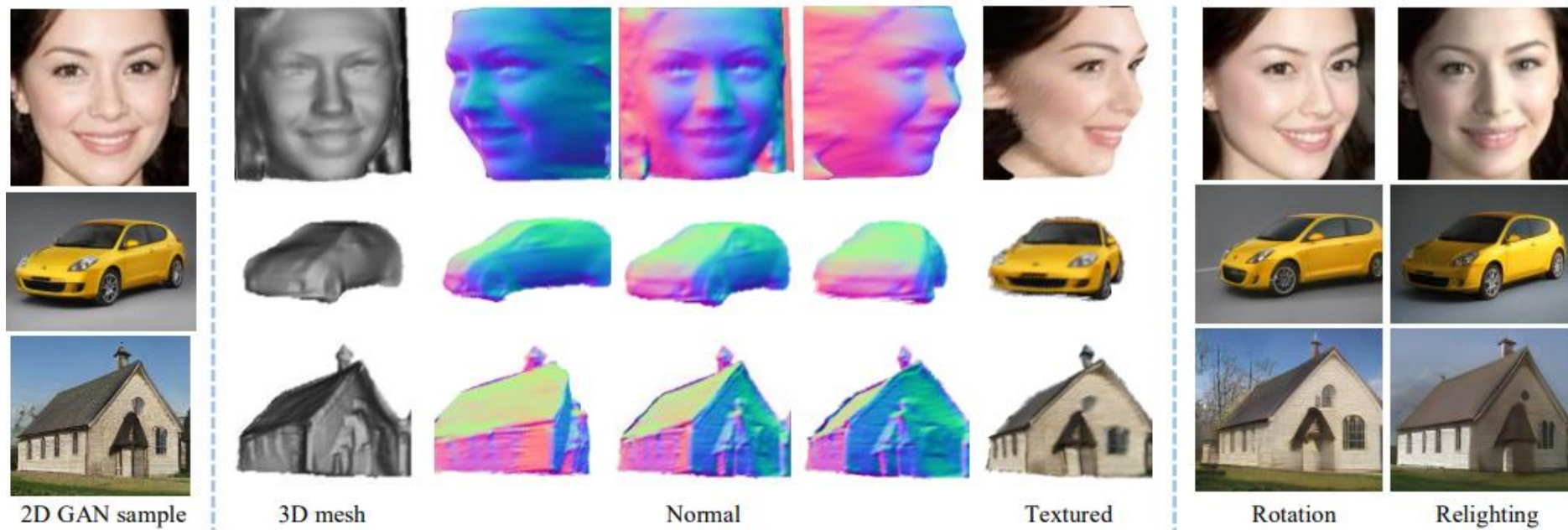
<https://arxiv.org/pdf/2011.00844.pdf>  
<https://github.com/XingangPan/GAN2Shape>  
<https://www.youtube.com/watch?v=Ea4eNXf3s24>

# Demo



# Introduction

## 2D GANs Know 3D Shape? Unsupervised 3D Shape Reconstruction from 2D Image GANs



- The framework is an iterative strategy that explores and exploits diverse viewpoint and lighting variations in the GAN image manifold.
- The framework does not require 2D keypoint or 3D annotations, or strong assumptions on object shapes, yet it successfully recovers 3D shapes.

# Background Knowledge

*Depth Map*



*Normal Map*



*Shading*



*Albedo*



*Original Img*



*Other Viewpoint*



# Background Knowledge

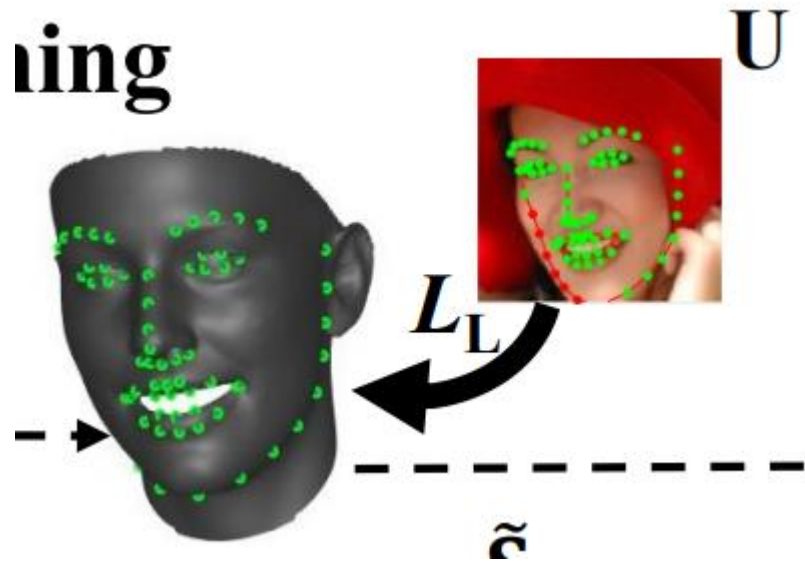
## Renderer

$$\begin{aligned} \hat{\mathbf{I}} &= \Phi(\mathbf{d}, \mathbf{a}, \mathbf{v}, \mathbf{l}) \\ \text{rendered Image} &= \Pi(\underbrace{\Lambda(\mathbf{d}, \mathbf{a}, \mathbf{l})}_{\text{reprojection}}, \underbrace{\mathbf{d}, \mathbf{v}}_{\text{lighting}}) \end{aligned}$$

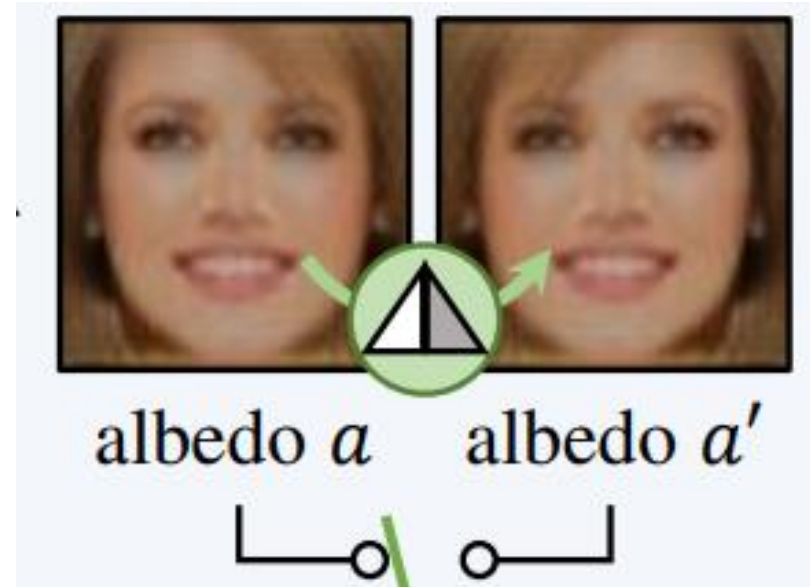
$$\begin{aligned} \mathbf{J}_{uv} &= (k_s + k_d \max\{0, \langle \mathbf{l}, \mathbf{n}_{uv} \rangle\}) \cdot a_{uv} \\ \mathbf{p} &\propto \mathbf{K} \mathbf{P}, \quad \mathbf{K} = \begin{bmatrix} f & 0 & c_i \\ 0 & f & c_j \\ 0 & 0 & 1 \end{bmatrix}, \quad \begin{cases} c_i = \frac{W-1}{2}, \\ c_j = \frac{H-1}{2}, \\ f = \frac{W-1}{2 \tan \frac{\theta_{FOV}}{2}} \end{cases} \\ \mathbf{p}' &\propto \mathbf{K} (d_{ij} \cdot \mathbf{R} \mathbf{K}^{-1} \mathbf{p} + \mathbf{T}) \end{aligned}$$



## Previous Work

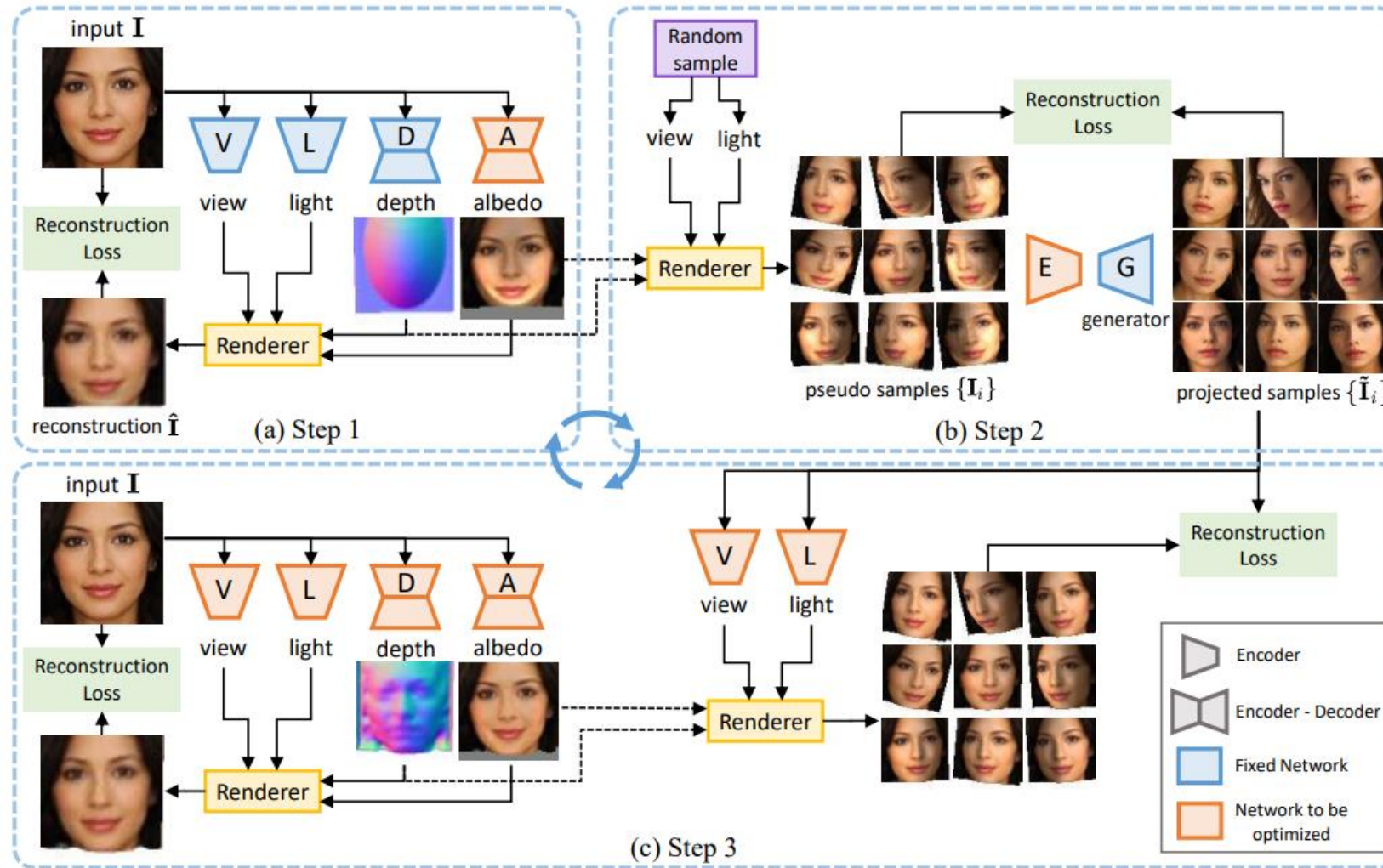


[3D model or 2D keypoints supervision]



[Symmetry Assumption]

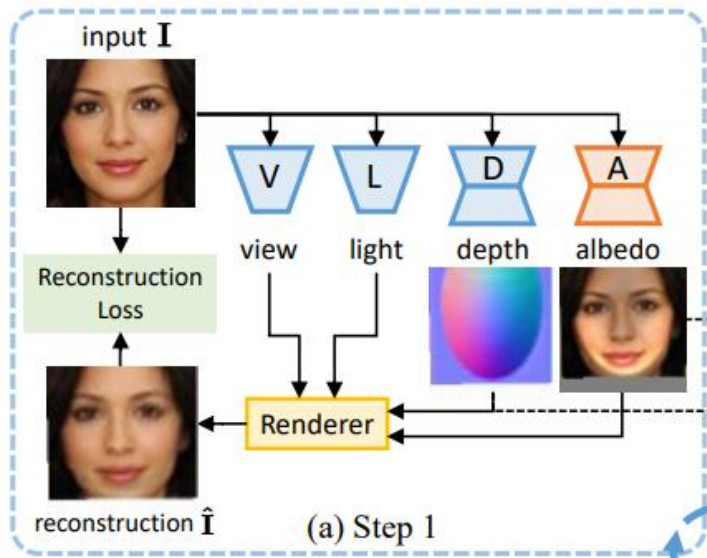
# Method



- Creating pseudo samples with different viewpoints and light via GAN to recover 3D shape.

# Method

## Step 1: Using a Weak Shape Prior.



- Initialize
  - view = 0, light = front, depth = ellipsoid
- Renderer

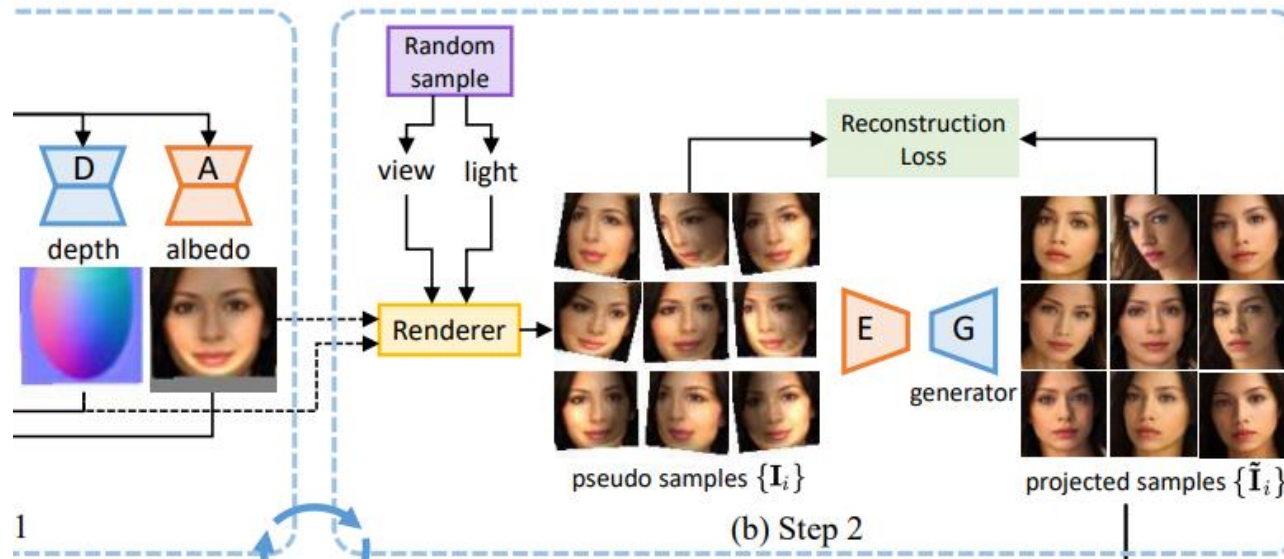
$$\begin{aligned}\hat{\mathbf{I}} &= \Phi(\mathbf{d}, \mathbf{a}, \mathbf{v}, \mathbf{l}) \\ &= \Pi(\Lambda(\mathbf{d}, \mathbf{a}, \mathbf{l}), \mathbf{d}, \mathbf{v})\end{aligned}$$

- Reconstruction loss
  - Weighted sum of L1 and Perceptual Loss



# Method

## Step2: Sampling and Projecting to the GAN Image Manifold.



$$\tilde{\mathbf{I}}_i = G(E(\mathbf{I}_i) + \mathbf{w})$$

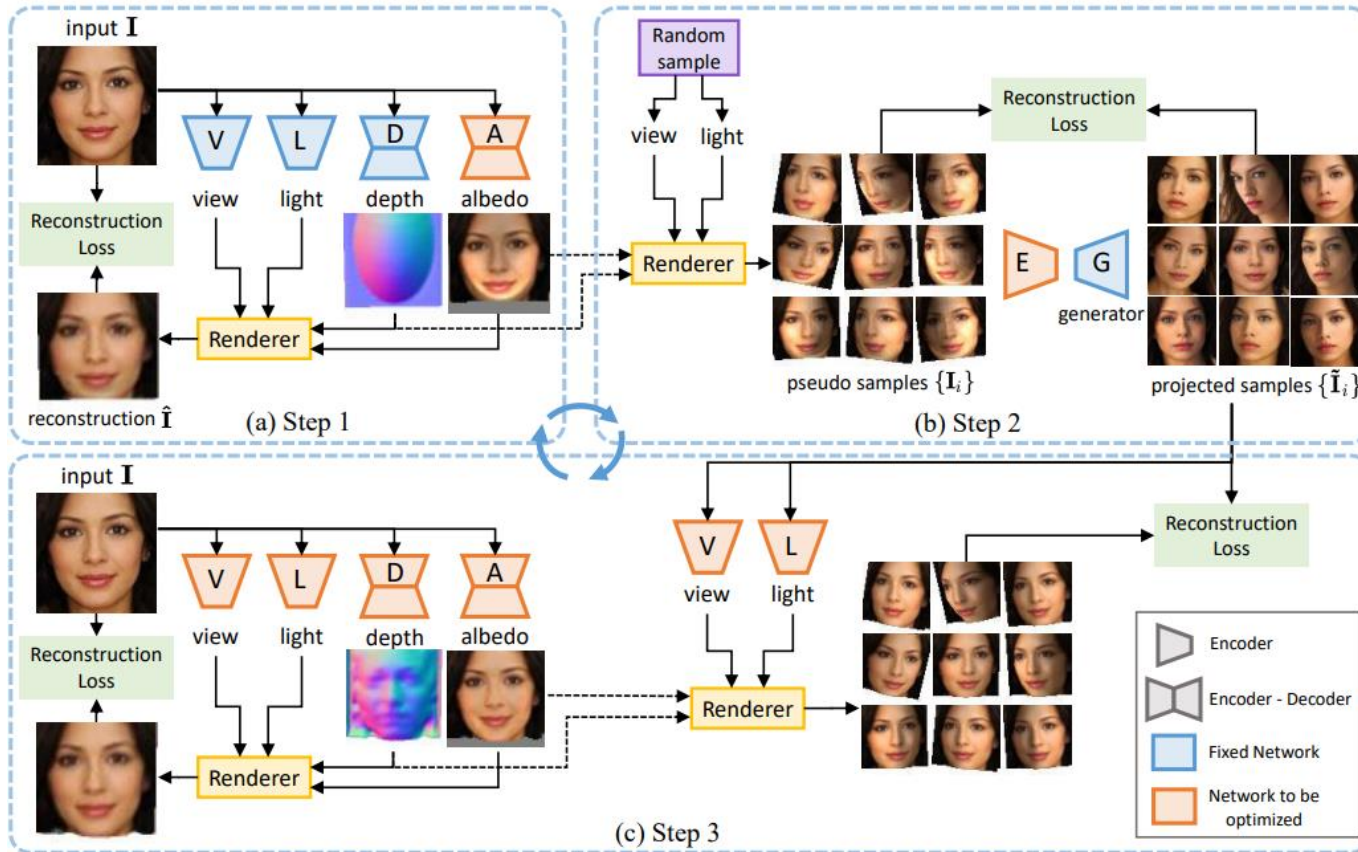
$$\theta_E = \arg \min_{\theta_E} \frac{1}{m} \sum_{i=0}^m \mathcal{L}'(\mathbf{I}_i, G(E(\mathbf{I}_i) + \mathbf{w})) + \lambda_1 \|E(\mathbf{I}_i)\|_2$$

$\mathcal{L}' = L1$  of the discriminator features

- With the shape prior as an initialization, we are able to create ‘pseudo samples’ by sampling a number of random viewpoints and lighting directions
- In order to leverage such cues to guide novel viewpoint and lighting direction exploration in the GAN image manifold, we perform GAN inversion to these pseudo samples (reconstruct them with the GAN generator  $G$ ).
- The generator  $G$  could regularize the projected samples to lie in the natural image manifold.

# Method

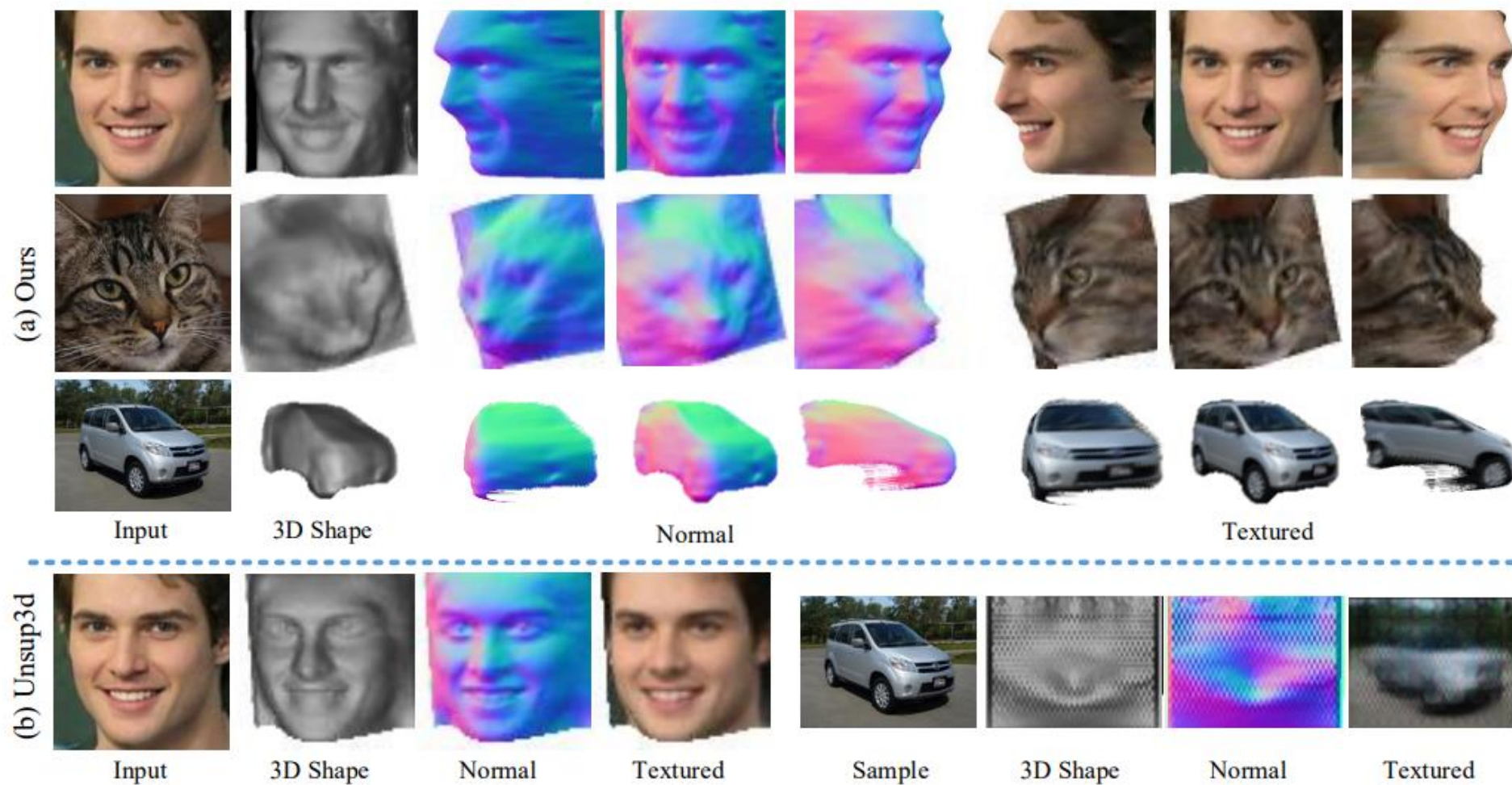
## Step3: Learning the 3D Shape.



- The projected samples provide images of multiple viewpoint and lighting conditions of nearly the same object content.
- Such single instance multiple view and lighting wetting makes it possible to learn the underlying 3D shape.

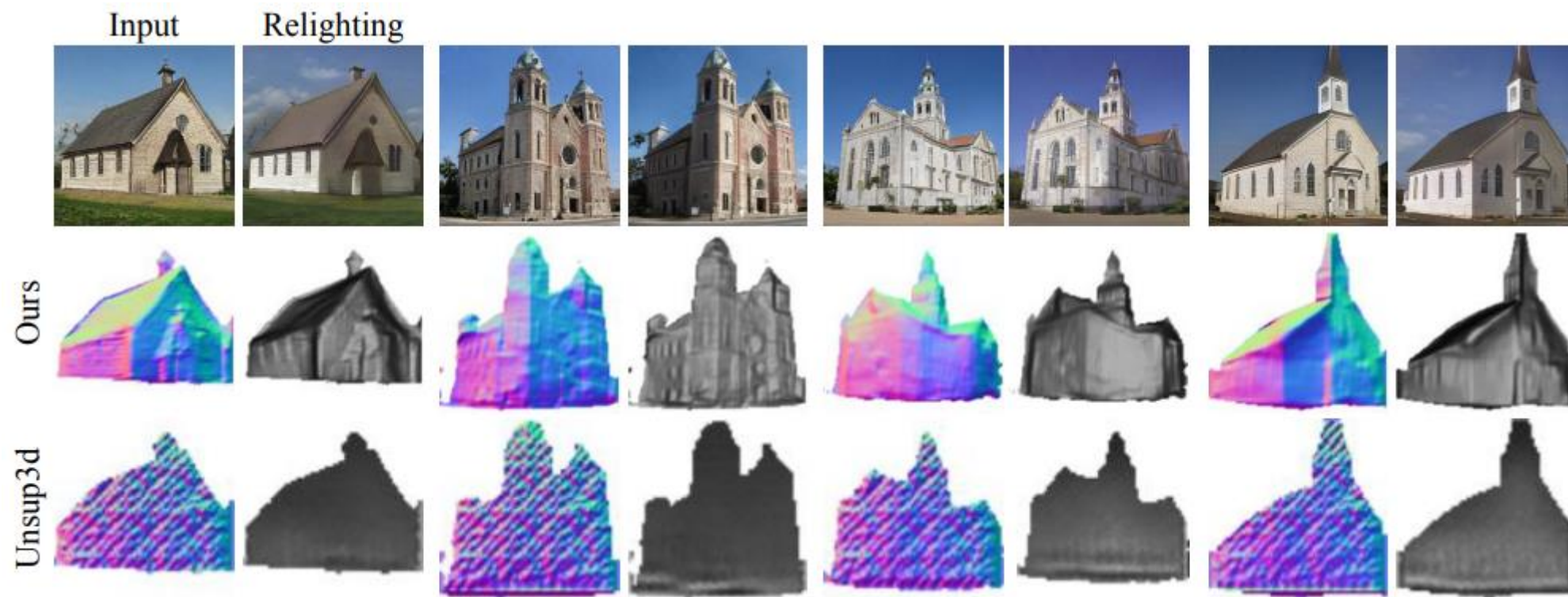
$$\theta_D, \theta_A, \theta_V, \theta_L = \arg \min_{\theta_D, \theta_A, \theta_V, \theta_L} \frac{1}{m} \sum_{i=0}^m \mathcal{L}(\tilde{I}_i, \Phi(D(I), A(I), V(\tilde{I}_i), L(\tilde{I}_i))) + \lambda_2 \mathcal{L}_{smooth}(D(I))$$

# Experiments





# Experiments



# Experiments

Table 1: **Comparisons on the BFM dataset.** We report SIDE and MAD errors. ‘Symmetry’ indicates whether the symmetry assumption on object shape is used. We outperform others on both metrics.

No.	Method	Symmetry	SIDE ( $\times 10^{-2}$ ) $\downarrow$	MAD (deg.) $\downarrow$
(1)	Supervised	N	0.419	10.83
(2)	Const. null depth	/	2.723	43.22
(3)	Average g.t. depth	/	1.978	22.99
(4)	Unsup3d (Wu et al., 2020)	Y	0.807	16.34
(5)	Ours (w/o regularize)	Y	0.925	16.42
(6)	Ours	Y	<b>0.756</b>	<b>14.81</b>
(7)	Unsup3d (Wu et al., 2020)	N	1.334	33.79
(8)	Ours	N	<b>1.023</b>	<b>17.09</b>

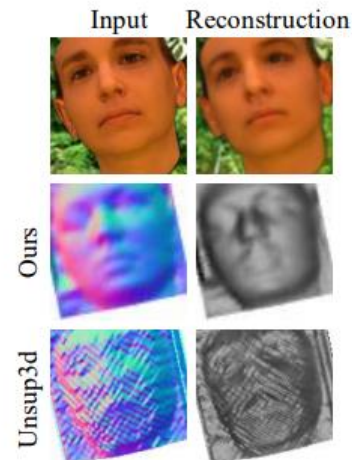
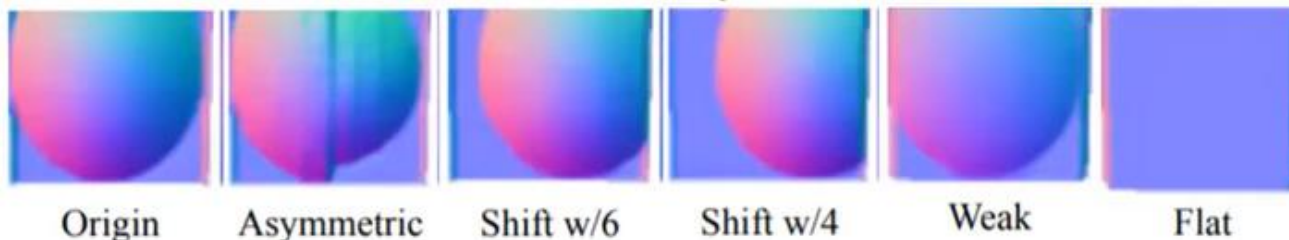


Figure 6: Results without symmetry assumption.

Table 2: **Effects of different shape prior.** We report results of the original ellipsoid shape, asymmetric shape with ellipsoid for the left half and sphere for the right half, shape with its position shifted by 1/6 and 1/4 image width, weaker shape prior whose height is half of the original one, and no shape prior. Qualitative results can be found in Fig.14 in the Appendix.

Shape prior	Origin	Asymmetric	Shift w/6	Shift w/4	Weak	Flat
SIDE ( $\times 10^{-2}$ ) $\downarrow$	0.756	0.769	0.767	0.775	0.764	1.021
MAD (deg.) $\downarrow$	14.81	14.95	14.93	15.07	14.97	20.46





# Experiments

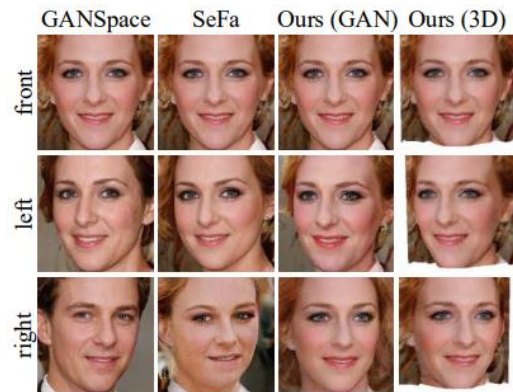


Figure 8: **Qualitative comparison on face rotation.** "Ours (GAN)" and "Ours (3D)" indicate results generated by GAN and rendered from 3D mesh respectively. The face identities in the baseline methods tend to drift during rotation.

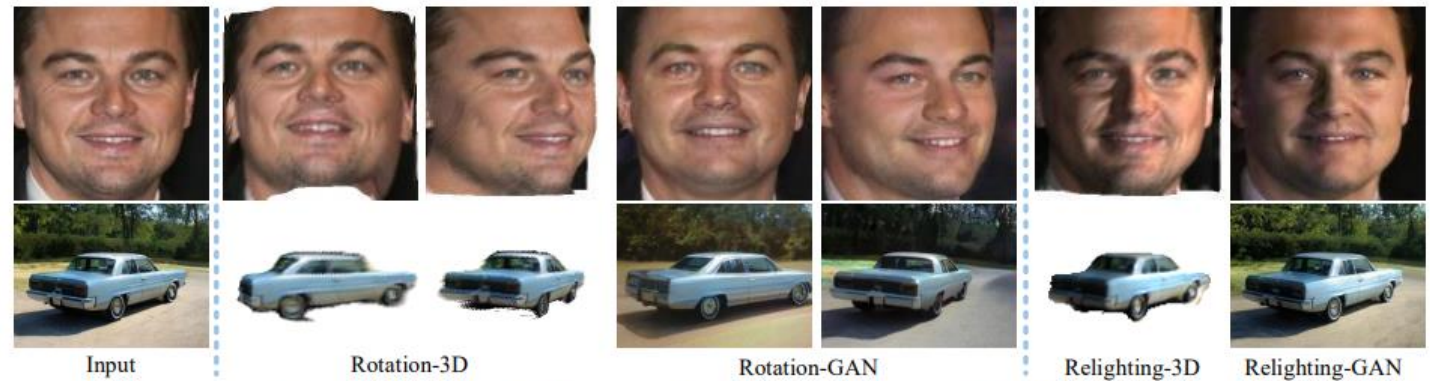


Figure 7: **3D-aware image manipulation**, including rotation and relighting. We show results obtained via both 3D mesh and GANs. The input of the first row is a real natural image. Our method achieves photo-realistic manipulation effects obeying the objects' underlying 3D structures.

## Limitation

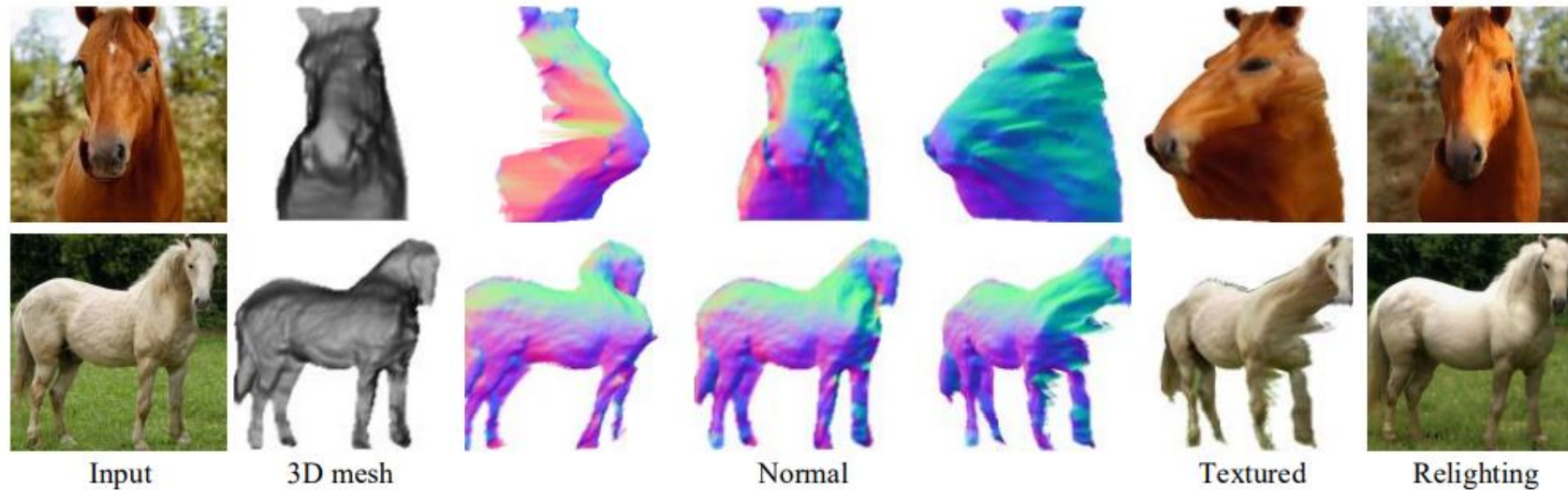


Figure 9: **Qualitative results on the LSUN Horse dataset** (Yu et al., 2015).

- For objects with more sophisticated shapes like horses, a simple convex shape prior may not well reflect the viewpoint and lighting variations.
- The model could not predict the back-side shape of objects.

EOD