



# Shape, Pose, and Appearance from a Single Image via Bootstrapped Radiance Field Inversion

Dario Pavlo  
David Joseph Tan  
Marie-Julie Rakotosaona  
Federico Tombari

CVPR 2023

**TUE-PM-025**

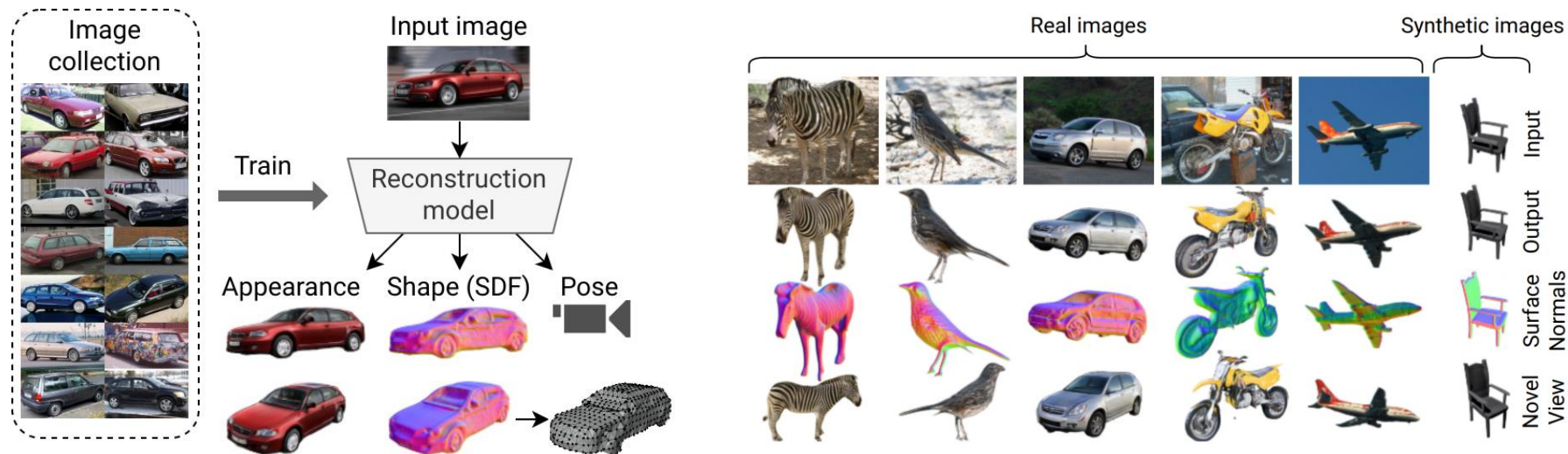


**ETH** zürich



# NeRF from a single image

- **Goal:** learn a model to reconstruct 3D shape, pose, and appearance from a single view of an object
  - NeRF with SDF shape parameterization
- Training **without** multiple views
- Focus on real datasets as opposed to synthetic datasets
  - Poses may be inaccurate



# Reconstruction demo



# Reconstruction demo

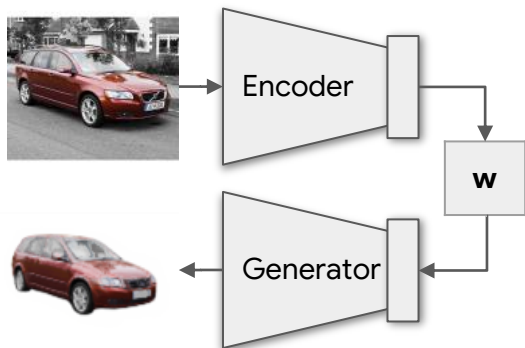




# Reconstruction demo

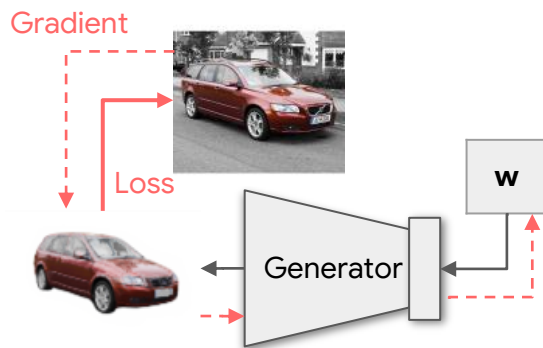


# Reconstruction frameworks & motivation



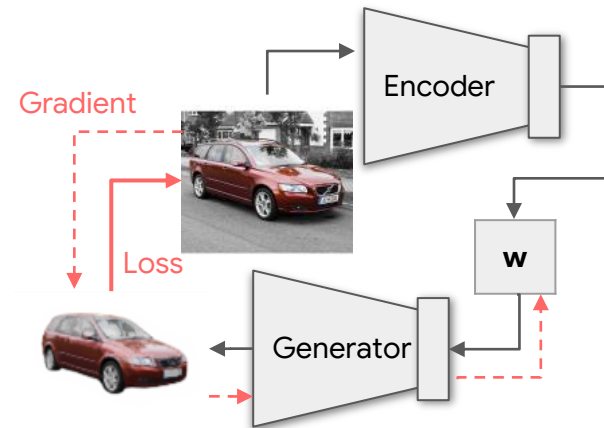
## Encoder-based

- E.g. CMR, PixelNeRF, Pix2NeRF
- Autoencoder setup: a ConvNet encoder predicts a latent code  $w$  which is decoded into a 3D scene
- **Fast but relies on accurate poses**, which are typically available only on synthetic datasets



## GAN Inversion via Optimization

- Leverages a pretrained unconditional GAN (e.g. pi-GAN, EG3D)
- Gradient-based optimization w.r.t. pose and latent code  $w$
- **Better results, robust to inaccurate poses, but very slow**

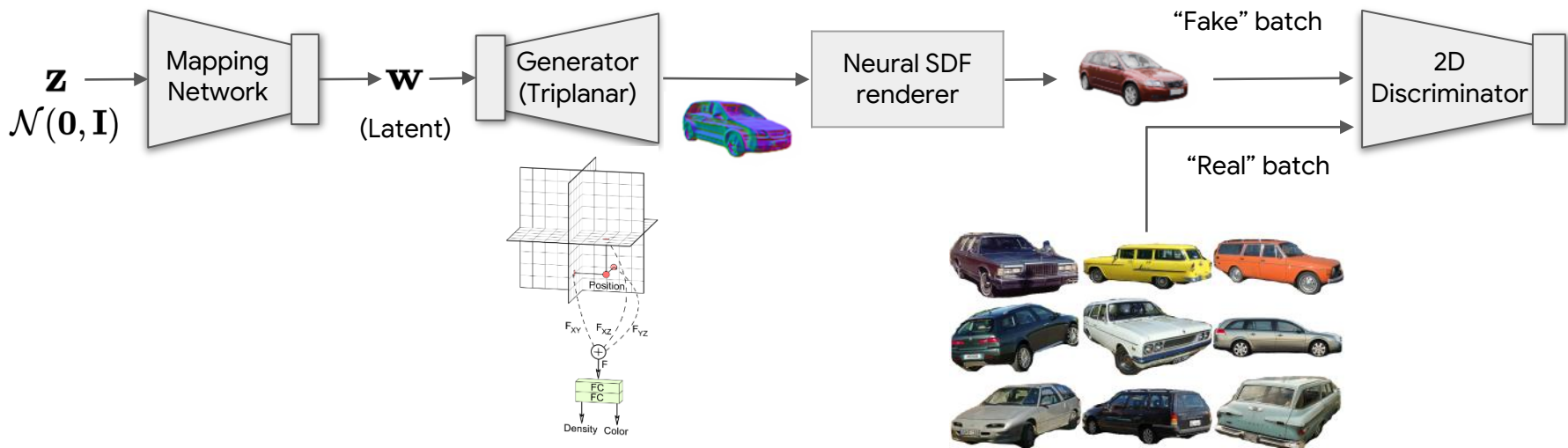


## Bootstrapped Inversion (Ours)

- **Not explored for NeRFs**
- A ConvNet encoder produces a first guess of the pose and latent code
- These are then refined via optimization for a small number of steps

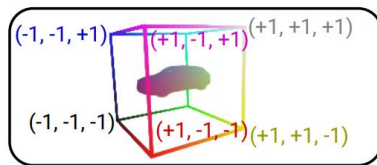
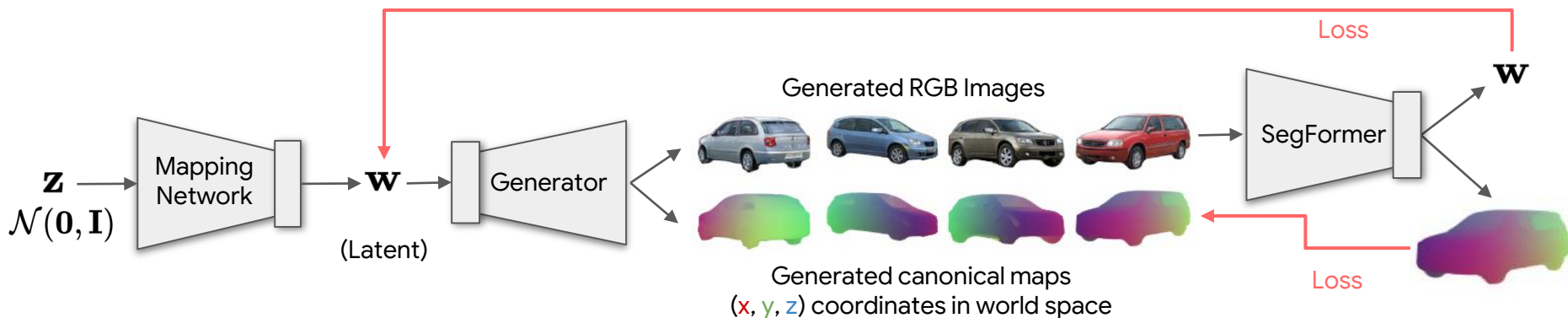
## Method (1/3): unconditional GAN training

- We train an unconditional 3D GAN on the collection of images (foundation model)
- Backbone inspired by EG3D (triplanar NeRF representation), with some improvements
  - **SDF representation (VolSDF)** → better surface reconstruction
  - **Color mapping:** disentanglement between color and semantics → facilitates inversion & manipulation
  - **Path length regularization** → facilitates inversion



## Method (2/3): bootstrapping & pose estimation

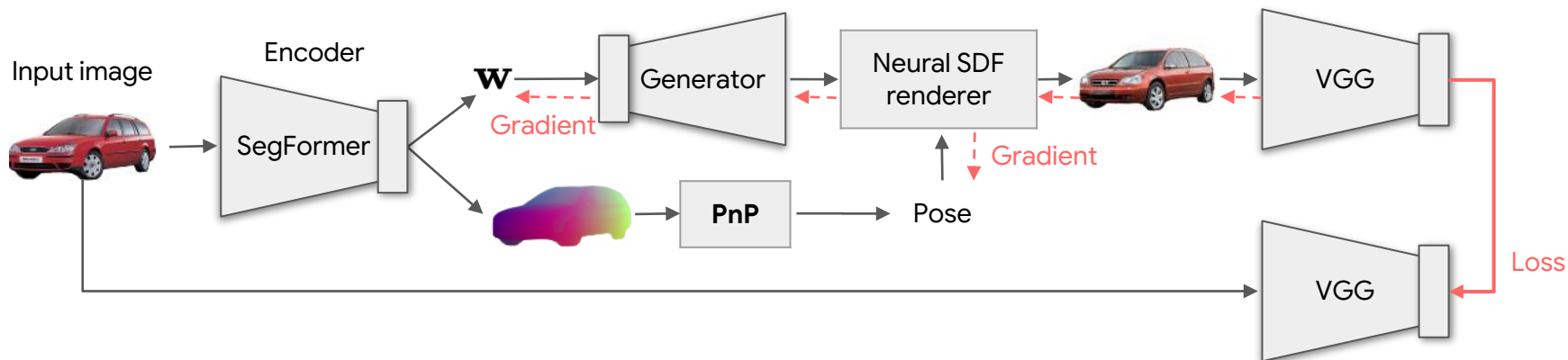
- Using generated data, we learn a model that jointly predicts the pose and the latent code  $\mathbf{w}$
- **NOCS approach:** predict a canonical map, convert to point cloud, and recover the pose using a PnP solver
  - More robust than directly regressing the pose parameters
  - Pseudo-ground-truth canonical maps generated using the GAN itself (rasterize  $xyz$  instead of  $rgb$ )





## Method (3/3): hybrid inversion

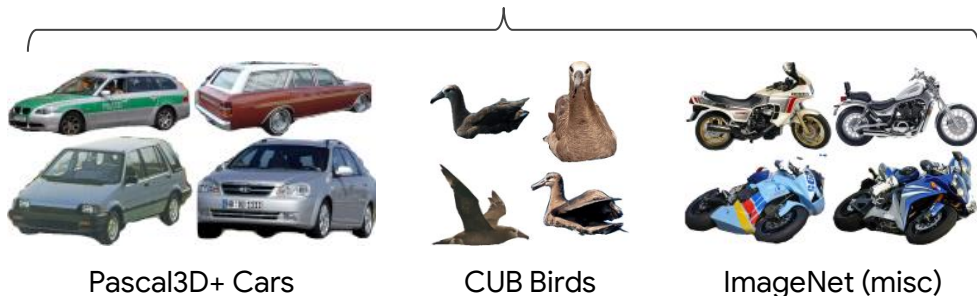
- At inference, we initially estimate the latent code  $\mathbf{w}$  and pose using the previous model
- These are then refined via optimization using a VGG loss
  - In practice, we use multiple crops to reduce the variance of the gradient
- We further investigated strategies to achieve maximum speed
  - We can invert an image in as few as **10 steps** (vs 100s of related work)



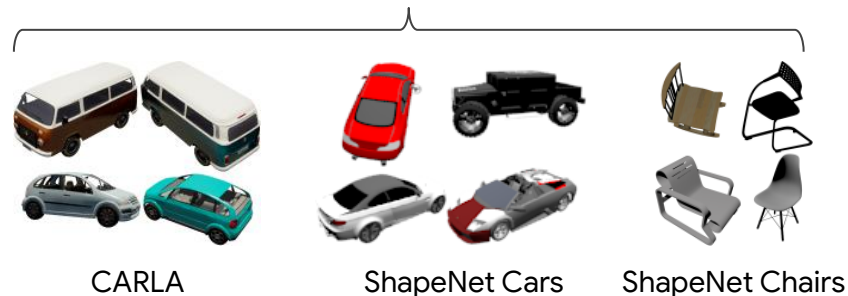
# Datasets & Evaluation

- Approach evaluated on a mix of synthetic and real datasets
- For real datasets, we compare to [CMR \(Kanazawa et al. 2018\)](#) & follow-up papers ([U-CMR](#), [UMR](#), ...)
  - Reconstruction evaluated using IoU against input view → easy to overfit
  - No ground-truth novel views → quality evaluated using FID on renderings from random views
- For ShapeNet, in addition to the FID, we evaluate the PSNR on novel views from the test set
  - Comparison against [Pix2NeRF \(Cai et al. 2022\)](#)

Real datasets



Synthetic datasets



# Qualitative results

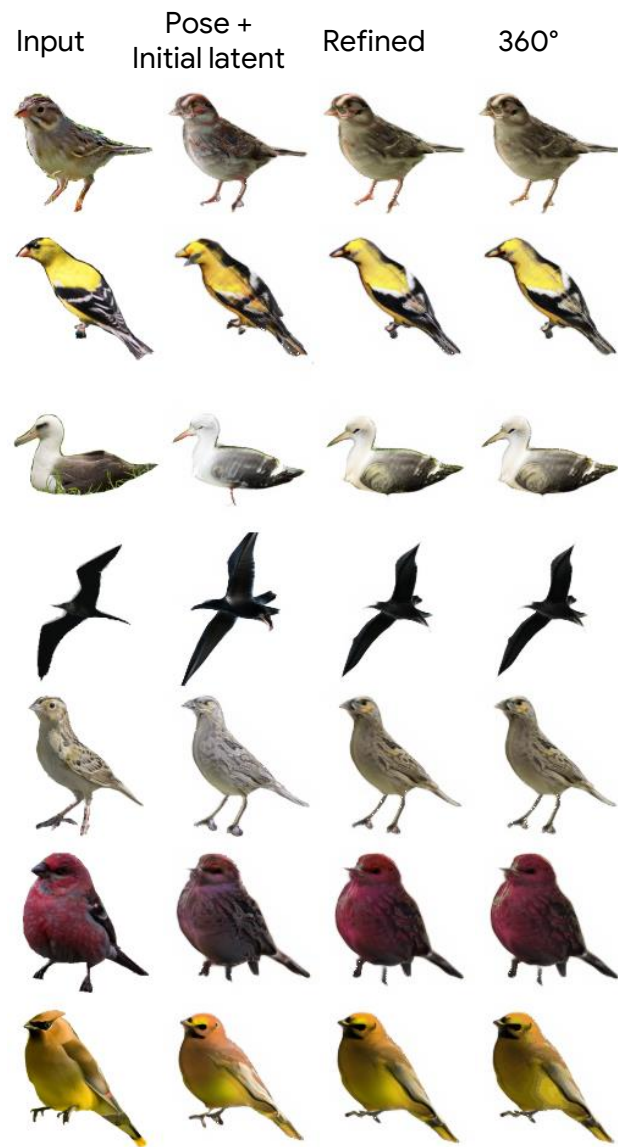
Real images

**Left:**

Pascal3D+ Cars  
(test set)

**Right:**

CUB Birds  
(test set)



# Qualitative results

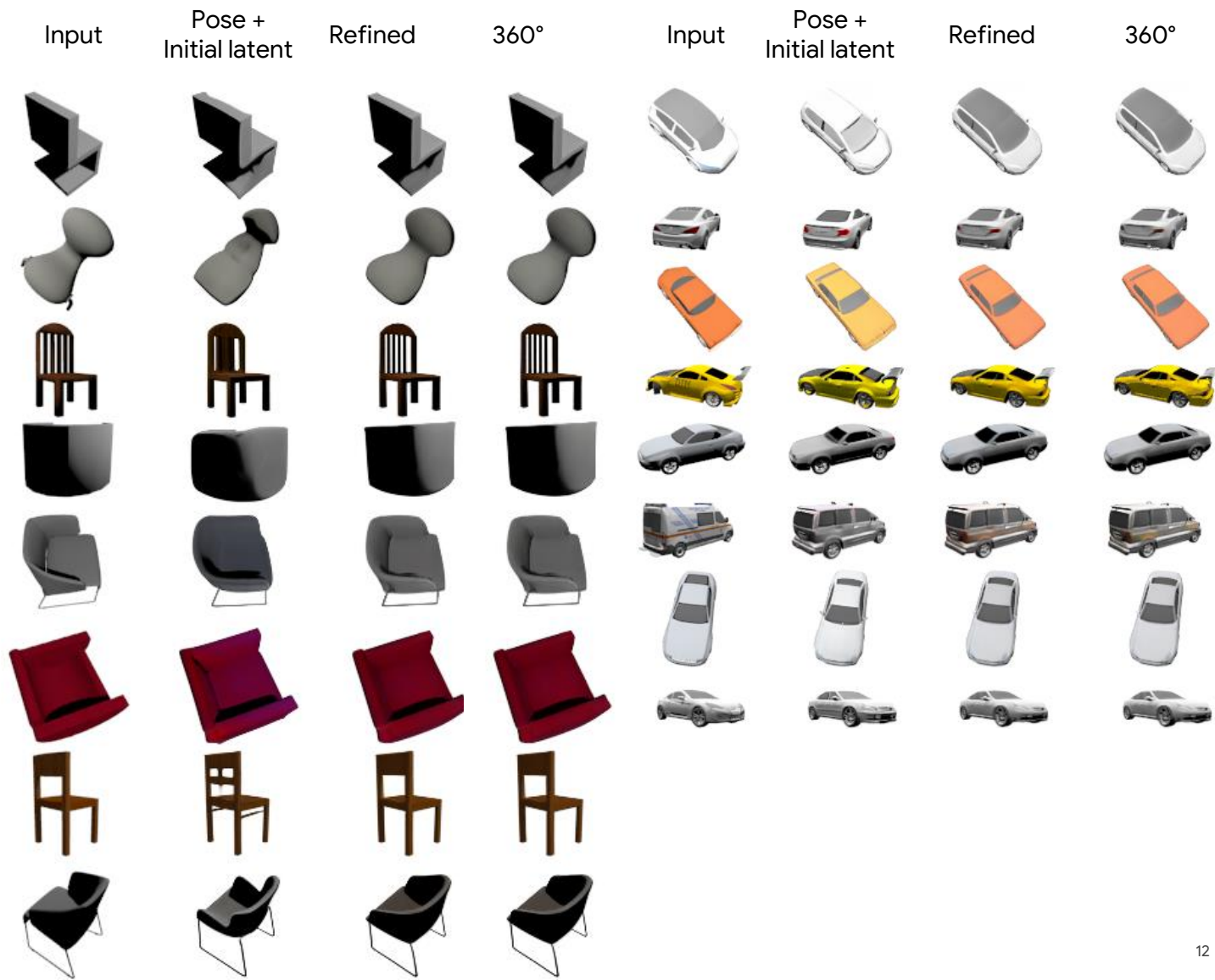
## Synthetic images

### Left:

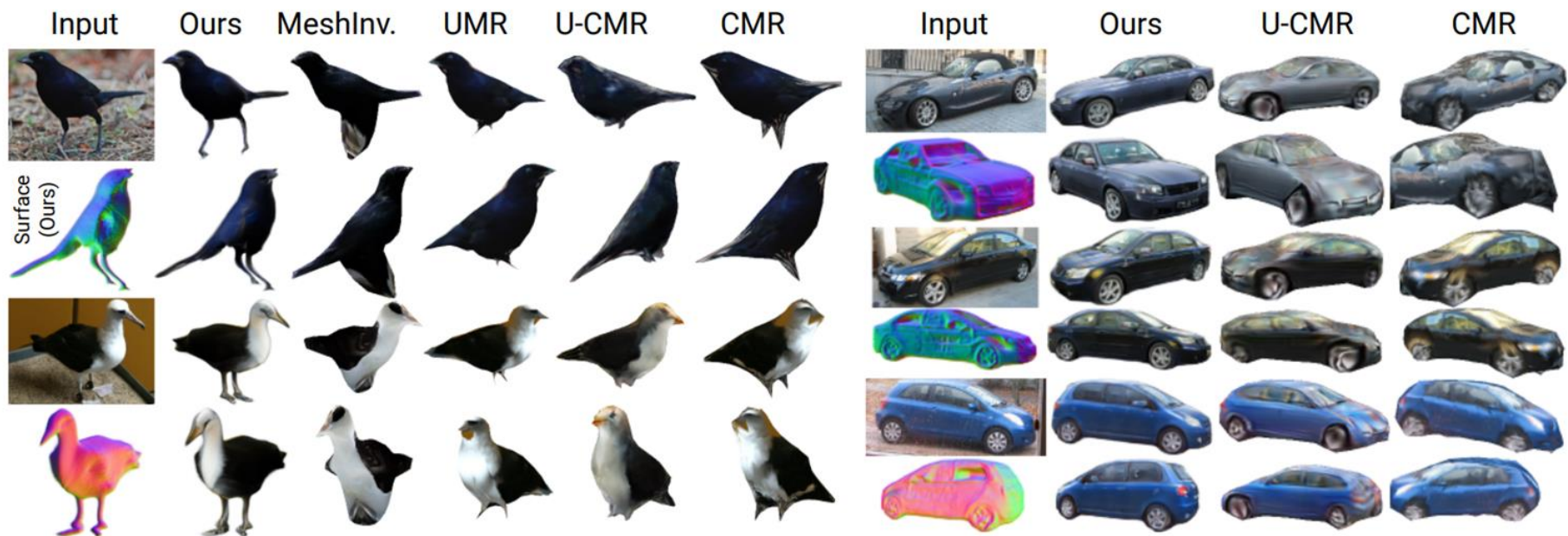
ShapeNet Chairs  
(test set)

### Right:

ShapeNet Cars  
(test set)



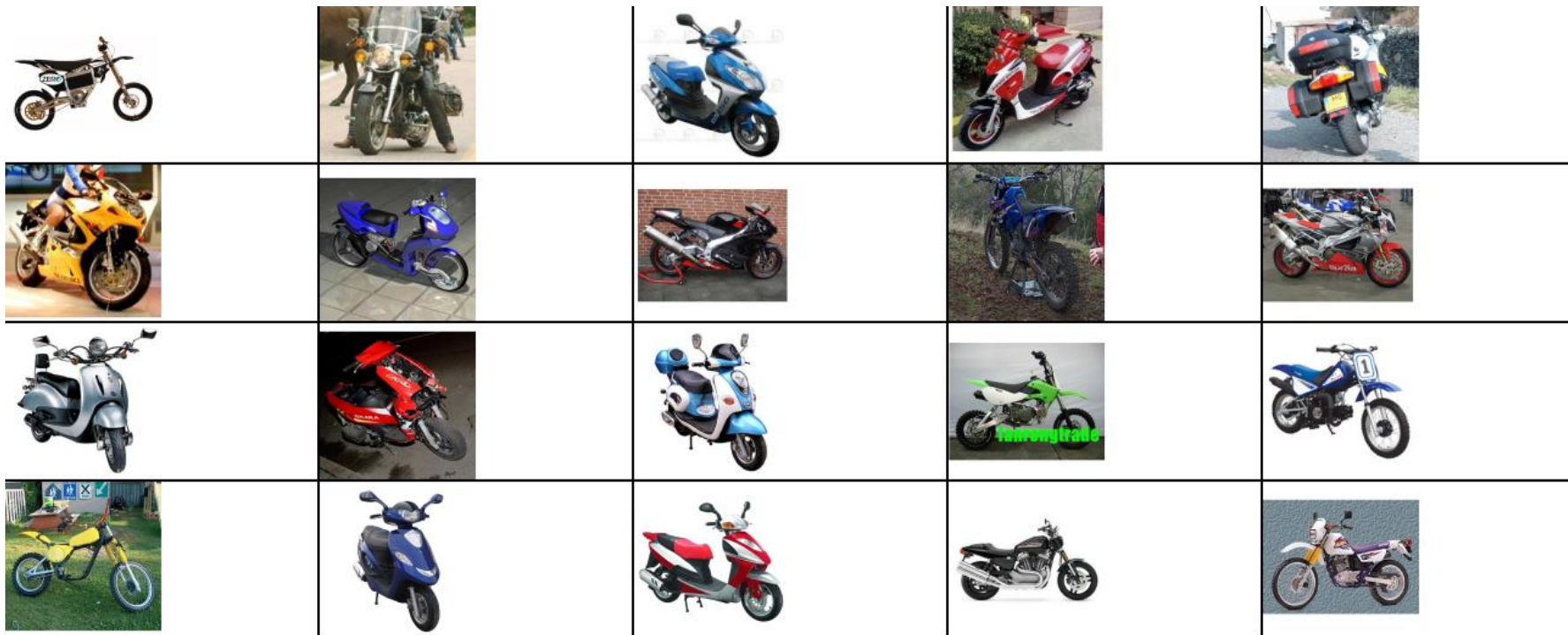
# Side-by-side comparison





# Qualitative results

End-to-end reconstruction pipeline on real images in the wild (ImageNet)



## End-to-end reconstruction: quantitative results

- Even without inversion ( $N = 0$  steps), we show an improvement over existing work
  - 36% decrease in FID on CUB over SOTA; 9% increase in IoU on P3D Cars
  - 68% decrease in FID on ShapeNet Chairs; 83% decrease in FID on CARLA
- Applying our hybrid inversion approach further widens the gap

Method	Pascal3D+ Cars		CUB Birds	
	IoU $\uparrow$	FID $\downarrow$	IoU $\uparrow$	FID $\downarrow$
CMR [27]	0.64	273.28	0.706	105.04
U-CMR [17]	0.646	<u>223.12</u>	0.644	69.42
UMR [33]	-	-	<u>0.734</u>	<u>43.83</u>
SDF-SRN [34]	<u>0.81</u>	254.90	-	-
ViewGeneralization [2]	0.78	-	0.629	-
StyleGANRender [71]	0.80	-	-	-
Ours Init. ( $N=0$ )	<b>0.883</b>	<b>75.90</b> (15.08)	<b>0.739</b>	<b>28.15</b>
MeshInv. ( $N=200$ ) (*) ( $\dagger$ ) [69]	-	-	0.752	31.60
Ours Hybrid Slow ( $N=30$ ) ( $\dagger$ )	<b>0.920</b>	<u>73.53</u> (14.36)	<b>0.844</b>	<b>24.70</b>
Ours Hybrid Fast ( $N=10$ ) ( $\dagger$ )	<u>0.917</u>	<b>73.12</b> (14.36)	<u>0.835</u>	<u>25.65</u>

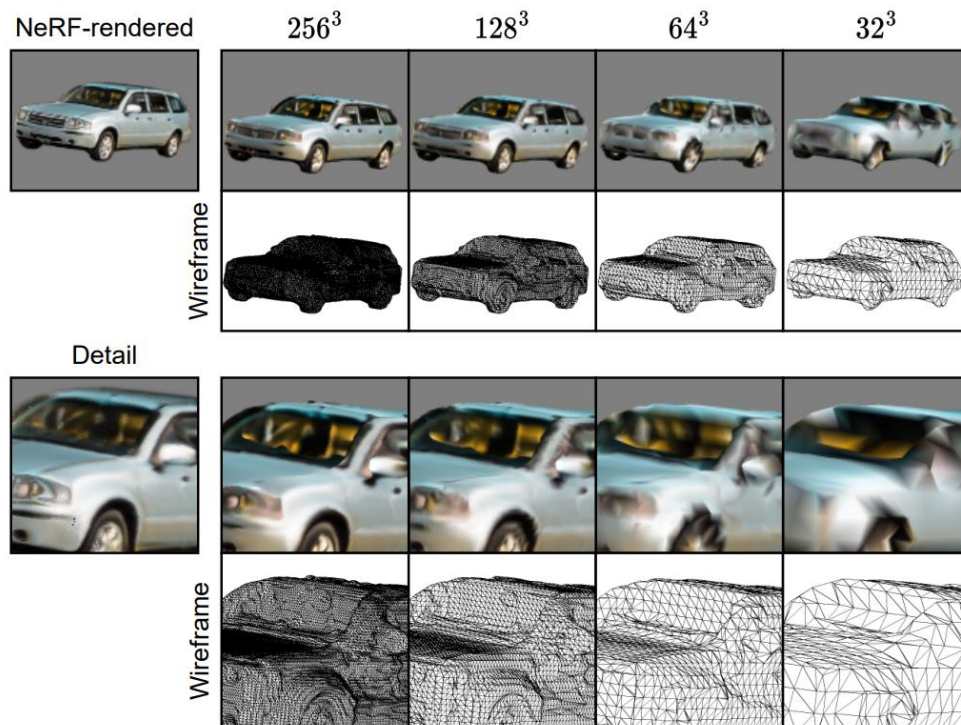
Evaluation on real datasets  
(Novel views **not** available)

Method	SRN Cars			SRN Chairs			CARLA
	PSNR $\uparrow$	SSIM $\uparrow$	FID $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	FID $\downarrow$	FID $\downarrow$
Pix2NeRF [4]	-	-	-	18.14	0.84	26.81	38.51
Ours Init. ( $N=0$ )	18.54	0.848	12.39	18.26	0.857	8.64	6.49
Ours Hybrid Slow ( $N=30$ )	<b>19.55</b>	<b>0.864</b>	<b>11.37</b>	<b>19.36</b>	<b>0.875</b>	<b>7.44</b>	<b>5.97</b>
Ours Hybrid Fast ( $N=10$ )	<u>19.24</u>	<u>0.861</u>	<u>12.26</u>	<u>19.02</u>	<u>0.871</u>	<u>7.62</u>	<u>6.18</u>

Evaluation on synthetic datasets  
(Novel views **available** on ShapeNet)

## Bonus: extraction of a triangle mesh from the SDF

- The adoption of an SDF representation allows us extract its 0-level set and obtain a colored triangle mesh
- Results at different step sizes for the marching cubes algorithm





Thank you for your attention!

Feel free to visit our poster:

**TUE-PM-025**

(Tuesday afternoon session, stand 25)