

Share With Thy Neighbors: Single-View Reconstruction by Cross-Instance Consistency

Tom Monnier¹ Matthew Fisher² Alexei A. Efros³ Mathieu Aubry¹

¹LIGM, Ecole des Ponts, Univ Gustave Eiffel ²Adobe Research ³UC Berkeley

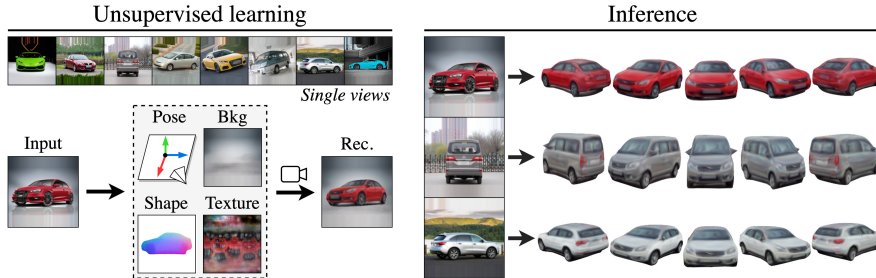


Fig. 1: **Single-View Reconstruction by Cross-Instance Consistency.** (left) Given a collection of single-view images from an object category, we learn without additional supervision an autoencoder that explicitly generates shape, texture, pose and background. (right) At inference time, our approach reconstructs high-quality textured meshes from raw single-view images.

Abstract. Approaches for single-view reconstruction typically rely on viewpoint annotations, silhouettes, the absence of background, multiple views of the same instance, a template shape, or symmetry. We avoid all such supervision and assumptions by explicitly leveraging the consistency between images of different object instances. As a result, our method can learn from large collections of unlabelled images depicting the same object category. Our main contributions are two ways for leveraging cross-instance consistency: (i) *progressive conditioning*, a training strategy to gradually specialize the model from category to instances in a curriculum learning fashion; and (ii) *neighbor reconstruction*, a loss enforcing consistency between instances having similar shape or texture. Also critical to the success of our method are: our structured autoencoding architecture decomposing an image into explicit shape, texture, pose, and background; an adapted formulation of differential rendering; and a new optimization scheme alternating between 3D and pose learning. We compare our approach, UNICORN, both on the diverse synthetic ShapeNet dataset — the classical benchmark for methods requiring multiple views as supervision — and on standard real-image benchmarks (Pascal3D+ Car, CUB) for which most methods require known templates and silhouette annotations. We also showcase applicability to more challenging real-world collections (CompCars, LSUN), where silhouettes are not available and images are not cropped around the object.

Keywords: single-view reconstruction, unsupervised learning

1 Introduction

One of the most magical human perceptual abilities is being able to see the 3D world behind a 2D image – a mathematically impossible task! Indeed, the ancient Greeks were so incredulous at the possibility that humans could be “hallucinating” the third dimension, that they proposed the utterly implausible Emission Theory of Vision [10] (eye emitting light to “sense” the world) to explain it to themselves. In the history of computer vision, single-view reconstruction (SVR) has had an almost cult status as one of the holy grail problems [19, 20, 46]. Recent advancements in deep learning methods have dramatically improved results in this area [6, 36]. However, the best methods still require costly supervision at training time, such as multiple views [34, 43]. Despite efforts to remove such requirements, the works with the least supervision still rely on two signals limiting their applicability: (i) silhouettes and (ii) strong assumptions such as symmetries [21, 25], known template shapes [12, 50], or the absence of background [56]. Although crucial to achieve reasonable results, priors like silhouettes and symmetry can also harm the reconstruction quality: silhouette annotations are often coarse [5] and small symmetry prediction errors can yield unrealistic reconstructions [12, 56].

In this paper, we propose the most unsupervised approach to single-view reconstruction to date, which we demonstrate to be competitive for diverse datasets. Table 1 summarizes the differences between our approach and representative prior works. More precisely, we learn in an analysis-by-synthesis fashion a network that predicts for each input image: 1) a 3D shape parametrized as a deformation of an ellipsoid, 2) a texture map, 3) a camera viewpoint, and 4) a background image (Fig. 1). Our main insight to remove the supervision and assumptions required by other methods is to leverage the consistency across different instances. First, we design a training procedure, *progressive conditioning*, which encourages the model to share elements between images by strongly constraining the variability of shape,

Method	Supervision	Data	Output
Pix2Mesh [53], AtlasNet [14], OccNet [36]	3D	ShapeNet	3D
PTN [59], NMR [30]	MV , CK , S	ShapeNet	3D
DRC [51], SoftRas [34], DVR [43]	MV , CK , S	ShapeNet	3D, T
GANverse3D [66]	MV , CK , S	Bird, Car, Horse	3D, T
DPC [23], MVC [49]	MV , S	ShapeNet	3D, P
Vicente <i>et al.</i> [52], CSDM [26], DRC [51]	CK , S	Pascal3D	3D
CMR [25]	CK , S , A (†)	Bird, Car, Plane	3D, T
SDF-SRN [33], TARS [8]	CK , S	ShapeNet, Bird, Car, Plane	3D, T
TexturedMeshGen [17]	CK , A (†)	ShapeNet, Bird, Car	3D, T
UCMR [12], IMR [50]	S , A (◇, †)	Animal, Car, Moto	3D, T , P
UMR [32]	S , A (↔, †)	Animal, Car, Moto	3D, T , P
RADAR [55]	S , A (‡)	Vase	3D, T , P
SMR [21]	S , A (†)	ShapeNet, Animal, Moto	3D, T , P
Unsup3D [56]	A (⊠, <, †)	Face	P , T , P
Henderson & Ferrari [16]	A (⊠, ∅)	ShapeNet	3D, P
Ours	None	ShapeNet, Animal, Car, Moto	3D, T , P

Table 1: **Comparison with selected works.** For each method, we outline the supervision and priors used (**3D**, **Multi-V**iews, **C**amera or **K**eypoints, **S**ilhouettes, **A**ssumptions like ◇ template shape, † symmetry, ‡ solid of revolution, ↔ semantic consistency, ⊠ no/limited background, < frontal view, ∅ no texture), which data it has been applied to and the model output (3D, **T**exture, **P**ose, **P**depth).

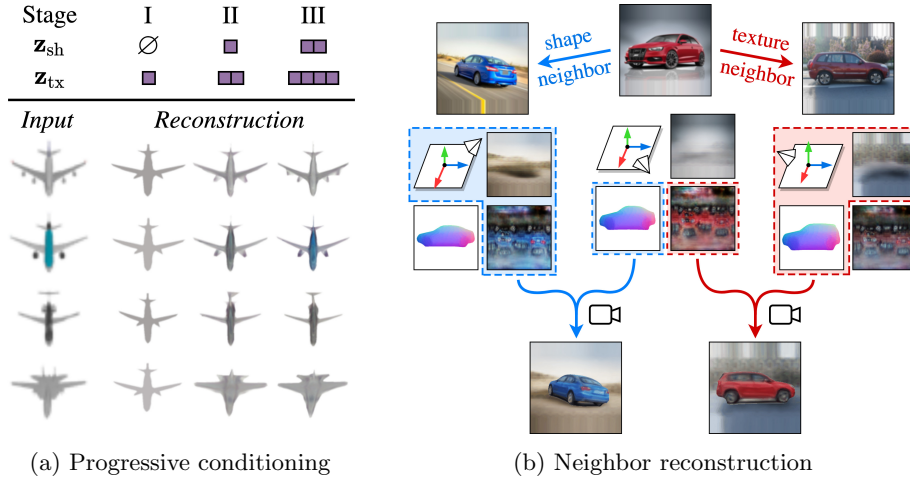


Fig. 2: **Leveraging cross-instance consistency.** (a) Progressive conditioning amounts to gradually increasing, in a multi-stage fashion, the size of the conditioning latent spaces, here associated to shape \mathbf{z}_{sh} and texture \mathbf{z}_{tx} . (b) We explicitly share the shape and texture models across neighboring instances by swapping their characteristics and applying a loss to associated neighbor reconstructions.

texture and background at the beginning of training and progressively allowing for more diversity (Fig. 2a). Second, we introduce a *neighbor reconstruction* loss, which explicitly enforces neighboring instances from different viewpoints to share the same shape or texture model (Fig. 2b). Note that these simple yet effective techniques are data-driven and not specific to any dataset. Our only remaining assumption is the knowledge of the semantic class of the depicted object.

We also provide two technical insights that we found critical to learn our model without viewpoint and silhouette annotations: (i) a differentiable rendering formulation inspired by layered image models [24, 38] which we found to perform better than the classical SoftRasterizer [34], and (ii) a new optimization strategy which alternates between learning a set of pose candidates with associated probabilities and learning all other components using the most likely candidate.

We validate our approach on the standard ShapeNet [4] benchmark, real image SVR benchmarks (Pascal3D+ Car [57], CUB [54]) as well as more complex real-world datasets (CompCars [60], LSUN Motorbike and Horse [63]). In all scenarios, we demonstrate high-quality textured 3D reconstructions.

Summary. We present UNICORN, a framework leveraging **UN**supervised cross-**I**nstance **C**onsistency for 3D **R**econstruction. Our main contributions are: 1) the first fully unsupervised SVR system, demonstrating state-of-the-art textured 3D reconstructions for both generic shapes and real images, and not requiring supervision or restrictive assumptions beyond a categorical image collection; 2) two data-driven techniques to enforce cross-instance consistency, namely *progressive conditioning* and *neighbor reconstruction*. Code and video results are available at imagine.enpc.fr/~monniet/UNICORN.

2 Related work

We first review deep SVR methods and mesh-based differential renderers we build upon. We then discuss works exploring cross-instance consistency and curriculum learning techniques, to which our progressive conditioning is related.

Deep SVR. There is a clear trend to remove supervision from deep SVR pipelines to directly learn 3D from raw 2D images, which we summarize in Table 1.

A first group of methods uses strong supervision, either paired 3D and images or multiple views of the same object. Direct 3D supervision is successfully used to learn voxels [6], meshes [53], parametrized surfaces [14] and implicit functions [36, 58]. The first methods using silhouettes and multiple views initially require camera poses and are also developed for diverse 3D shape representations: [51, 59] opt for voxels, [5, 30, 34] introduce mesh renderers, and [43] adapts implicit representations. Works like [23, 49] then introduce techniques to remove the assumption of known poses. Except for [66] which leverages GAN-generated images [13, 27], these works are typically limited to synthetic datasets.

A second group of methods aims at removing the need for 3D and multi-view supervision. This is very challenging and they hence typically focus on learning 3D from images of a single category. Early works [26, 51, 52] estimate camera poses with keypoints and minimize the silhouette reprojection error. The ability to predict textures is first incorporated by CMR [25] which, in addition to keypoints and silhouettes, uses symmetry priors. Recent works [8, 33] replace the mesh representation of CMR with implicit functions that do not require symmetry priors, yet the predicted texture quality is strongly deteriorated. [17] improves upon CMR and develops a framework for images with camera annotations that does not rely on silhouettes. Two works managed to further avoid the need for camera estimates but at the cost of additional hypothesis: [16] shows results with textureless synthetic objects, [56] models 2.5D objects like faces with limited background and viewpoint variation. Finally, recent works only require object silhouettes but they also make additional assumptions: [12, 50] use known template shapes, [32] assumes access to an off-the-shelf system predicting part semantics, and [55] targets solids of revolution. Other related works [11, 18, 21, 29, 44, 62] leverage in addition generative adversarial techniques to improve the learning.

In this work, we *do not* use camera estimates, keypoints, silhouettes, nor strong dataset-specific assumptions, and demonstrate results for both diverse shapes and real images. To the best of our knowledge, we present the first generic SVR system learned from raw image collections.

Mesh-based differentiable rendering. We represent 3D models as meshes with parametrized surfaces, as introduced in AtlasNet [14] and advocated by [50]. We optimize the mesh geometry, texture and camera parameters associated to an image using differentiable rendering. Loper and Black [35] introduce the first generic differentiable renderer by approximating derivatives with local filters, and [30] proposes an alternative approximation more suitable to learning neural networks. Another set of methods instead approximates the rendering function to allow differentiability, including SoftRasterizer [34, 45] and DIB-R [5]. We refer

the reader to [28] for a comprehensive study. We build upon SoftRasterizer [34], but modify the rendering function to learn without silhouette information.

Cross-instance consistency. Although all methods learned on categorical image collections implicitly leverage the consistency across instances, few recent works explicitly explore such a signal. Inspired by [31], the SVR system of [21] is learned by enforcing consistency between the interpolated 3D attributes of two instances and attributes predicted for the associated reconstruction. [61] discovers 3D parts using the inconsistency of parts across instances. Closer to our approach, [39] introduces a loss enforcing cross-silhouette consistency. Yet it differs from our work in two ways: (i) the loss operates on silhouettes, whereas our loss is adapted to image reconstruction by modeling background and separating two terms related to shape and texture, and (ii) the loss is used as a refinement on top of two cycle consistency losses for poses and 3D reconstructions, whereas we demonstrate results without additional self-supervised losses.

Curriculum learning. The idea of learning networks by “starting small” dates back to Elman [9] where two curriculum learning schemes are studied: (i) increasing the difficulty of samples, and (ii) increasing the model complexity. We respectively coin them *curriculum sampling* and *curriculum modeling* for differentiation. Known to drastically improve the convergence speed [2], curriculum sampling is widely adopted across various applications [1, 22, 47]. On the contrary, curriculum modeling is typically less studied although crucial to various methods. For example, [53] performs SVR in a coarse-to-fine manner by increasing the number of mesh vertices, and [37] clusters images by aligning them with transformations that increase in complexity. We propose a new form of curriculum modeling dubbed *progressive conditioning* which enables us to avoid bad minima.

3 Approach

Our goal is to learn a neural network that reconstructs a textured 3D object from a single input image. We assume we have access to a raw collection of images depicting objects from the same category, without any further annotation. We propose to learn in an analysis-by-synthesis fashion by autoencoding images in a structured way (Fig. 3). We first introduce our structured autoencoder (Sec. 3.1). We then present how we learn models consistent across instances (Sec. 3.2). Finally, we discuss one more technical contribution necessary to our system: an alternate optimization strategy for joint 3D and pose estimation (Sec. 3.3).

Notations. We use bold lowercase for vectors (e.g., \mathbf{a}), bold uppercase for images (e.g., \mathbf{A}), double-struck uppercase for meshes (e.g., \mathbb{A}), calligraphic uppercase for the main modules of our system (e.g., \mathcal{A}), lowercase indexed with generic parameters θ for networks (e.g., a_θ), and write $a_{1:N}$ the ordered set $\{a_1, \dots, a_n\}$.

3.1 Structured autoencoding

Overview. Our approach can be seen as a structured autoencoder: it takes an image as input, computes parameters with an encoder, and decodes them into

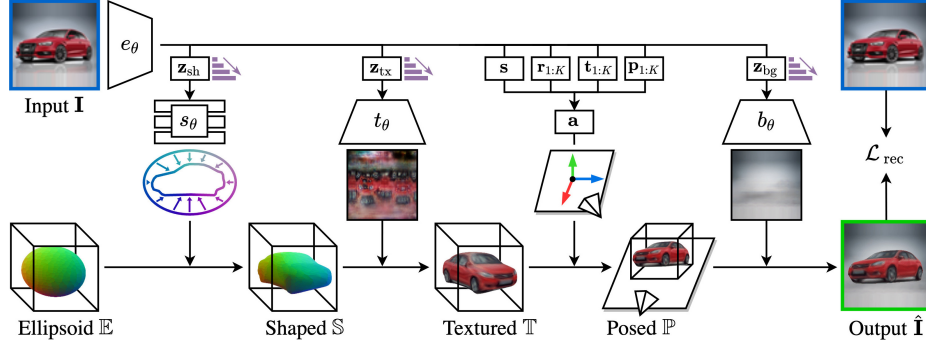



Fig. 3: **Structured autoencoding.** Given an **input**, we predict parameters that are decoded into 4 factors (shape, texture, pose, background) and composed to generate the **output**. Progressive conditioning is represented with .

explicit and interpretable factors that are composed to generate an image. We model images as the rendering of textured meshes on top of background images. For a given image I , our model thus predicts a shape, a texture, a pose and a background which are composed to get the reconstruction \hat{I} , as shown in Figure 3. More specifically, the image I is fed to convolutional encoder networks e_θ which output parameters $e_\theta(I) = \{z_{sh}, z_{tx}, a, z_{bg}\}$ used for the decoding part. a is a 9D vector including the object pose, while the dimension of the latent codes z_{sh} , z_{tx} and z_{bg} will vary during training (see Sec. 3.2). In the following, we describe the decoding modules using these parameters to build the final image by generating a shape, adding texture, positioning it and rendering it over a background.

Shape deformation. We follow [50] and use the parametrization of AtlasNet [14] where different shapes are represented as deformation fields applied to the unit sphere. We apply the deformation to an icosphere slightly stretched into an ellipsoid mesh E using a fixed anisotropic scaling. More specifically, given a 3D vertex x of the ellipsoid, our shape deformation module $\mathcal{S}_{z_{sh}}$ is defined by $\mathcal{S}_{z_{sh}}(x) = x + s_\theta(x, z_{sh})$, where s_θ is a Multi-Layer Perceptron taking as input the concatenation of a 3D point x and a shape code z_{sh} . Applying this displacement to all the ellipsoid vertices enables us to generate a shaped mesh $S = \mathcal{S}_{z_{sh}}(E)$. We found that using an ellipsoid instead of a raw icosphere was very effective in encouraging the learning of objects aligned w.r.t. the canonical axes. Learning surface deformations is often preferred to vertex-wise displacements as it enables mapping surfaces, and thus meshes, at any resolution. For us, the mesh resolution is kept fixed and such a representation is a way to regularize the deformations.

Texturing. Following the idea of CMR [25], we model textures as an image UV-mapped onto the mesh through the reference ellipsoid. More specifically, given a texture code z_{tx} , a convolutional network t_θ is used to produce an image $t_\theta(z_{tx})$, which is UV-mapped onto the sphere using spherical coordinates to associate a 2D point to every vertex of the ellipsoid, and thus to each vertex of the shaped mesh. We write $\mathcal{T}_{z_{tx}}$ this module generating a textured mesh $T = \mathcal{T}_{z_{tx}}(S)$.

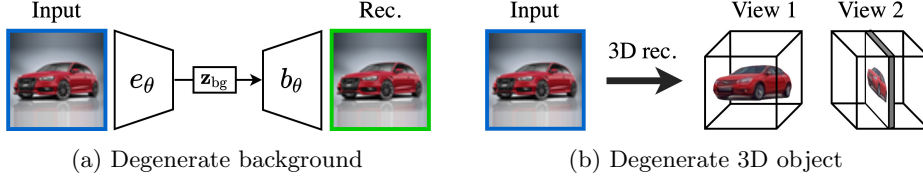


Fig. 4: **Degenerate solutions.** An SVR system learned by raw image autoencoding is prone to degenerate solutions through (a) the background or (b) the 3D object model. We alleviate the issue with cross-instance consistency.

Affine transformation. To render the textured mesh \mathbb{T} , we define its position w.r.t. the camera. In addition, we found it beneficial to explicitly model an anisotropic scaling of the objects. Because predicting poses is difficult, we predict K poses candidates, defined by rotations $\mathbf{r}_{1:K}$ and translations $\mathbf{t}_{1:K}$, and associated probabilities $\mathbf{p}_{1:K}$. This involves learning challenges we tackle with a specific optimization procedure described in Sec. 3.3. At inference, we select the pose with highest probability. We combine the scaling and the most likely 6D pose in a single affine transformation module \mathcal{A}_a . More precisely, \mathcal{A}_a is parametrized by $\mathbf{a} = \{\mathbf{s}, \mathbf{r}, \mathbf{t}\}$, where $\mathbf{s}, \mathbf{r}, \mathbf{t} \in \mathbb{R}^3$ respectively correspond to the three scales of an anisotropic scaling, the three Euler angles of a rotation and the three coordinates of a translation. A 3D point \mathbf{x} on the mesh is then transformed by $\mathcal{A}_a(\mathbf{x}) = \text{rot}(\mathbf{r})\text{diag}(\mathbf{s})\mathbf{x} + \mathbf{t}$ where $\text{rot}(\mathbf{r})$ is the rotation matrix associated to \mathbf{r} and $\text{diag}(\mathbf{s})$ is the diagonal matrix associated to \mathbf{s} . Our module is applied to all points of the textured mesh \mathbb{T} resulting in a posed mesh $\mathbb{P} = \mathcal{A}_a(\mathbb{T})$.

Rendering with background. The final step of our process is to render the mesh over a background image. The background image is generated from a background code \mathbf{z}_{bg} by a convolutional network b_θ . A differentiable module $\mathcal{B}_{\mathbf{z}_{bg}}$ renders the posed mesh \mathbb{P} over this background image $b_\theta(\mathbf{z}_{bg})$ resulting in a reconstructed image $\hat{\mathbf{I}} = \mathcal{B}_{\mathbf{z}_{bg}}(\mathbb{P})$. We perform rendering through soft rasterization of the mesh. Because we observed divergence results when learning geometry from raw photometric comparison with the standard SoftRasterizer [34, 45], we propose two key changes: a layered aggregation of the projected faces and an alternative occupancy function. We provide details in our supplementary material.

3.2 Unsupervised learning with cross-instance consistency

We propose to learn our structured autoencoder without any supervision, by synthesizing 2D images and minimizing a reconstruction loss. Due to the unconstrained nature of the problem, such an approach typically yields degenerate solutions (Fig. 4a and Fig. 4b). While previous works leverage silhouettes and dataset-specific priors to mitigate this issue, we instead propose two unsupervised data-driven techniques, namely *progressive conditioning* (a training strategy) and *neighbor reconstruction* (a training loss). We thus optimize the shape, texture and background by minimizing for each image \mathbf{I} reconstructed as $\hat{\mathbf{I}}$:

$$\mathcal{L}_{3D} = \mathcal{L}_{\text{rec}}(\mathbf{I}, \hat{\mathbf{I}}) + \lambda_{\text{nbr}}\mathcal{L}_{\text{nbr}} + \lambda_{\text{reg}}\mathcal{L}_{\text{reg}}, \quad (1)$$

where λ_{nbr} and λ_{reg} are scalar hyperparameters, and \mathcal{L}_{rec} , \mathcal{L}_{nbr} and \mathcal{L}_{reg} are respectively the reconstruction, neighbor reconstruction, and regularization losses, described below. In all experiments, we use $\lambda_{\text{nbr}} = 1$ and $\lambda_{\text{reg}} = 0.01$. Note that we optimize pose prediction using a slightly different loss in an alternate optimization scheme described in Section 3.3.

Reconstruction and regularization losses. Our reconstruction loss has two terms, a pixel-wise squared L_2 loss \mathcal{L}_{pix} and a perceptual loss [65] $\mathcal{L}_{\text{perc}}$ defined as an L_2 loss on the `relu3_3` layer of a pre-trained VGG16 [48], similar to [56]. While pixel-wise losses are common for autoencoders, we found it crucial to add a perceptual loss to learn textures that are discriminative for the pose estimation. Our full reconstruction loss can be written $\mathcal{L}_{\text{rec}}(\mathbf{I}, \hat{\mathbf{I}}) = \mathcal{L}_{\text{pix}}(\mathbf{I}, \hat{\mathbf{I}}) + \lambda_{\text{perc}} \mathcal{L}_{\text{perc}}(\mathbf{I}, \hat{\mathbf{I}})$ and we use $\lambda_{\text{perc}} = 10$ in all experiments. While our deformation-based surface parametrization naturally regularizes the shape, we sometimes observe bad minima where the surface has folds. Following prior works [5, 12, 34, 64], we thus add a small regularization term $\mathcal{L}_{\text{reg}} = \mathcal{L}_{\text{norm}} + \mathcal{L}_{\text{lap}}$ consisting of a normal consistency loss [7] $\mathcal{L}_{\text{norm}}$ and a Laplacian smoothing loss [40] \mathcal{L}_{lap} .

Progressive conditioning. The goal of *progressive conditioning* is to encourage the model to share elements (*e.g.*, shape, texture, background) across instances to prevent degenerate solutions. Inspired by the curriculum learning philosophy [9, 37, 53], we propose to do so by gradually increasing the latent space representing the shape, texture and background. Intuitively, restricting the latent space implicitly encourages maximizing the information shared across instances. For example, a latent space of dimension 0 (*i.e.*, no conditioning) amounts to learning a global representation that is the same for all instances, while a latent space of dimension 1 restricts all the generated shapes, textures or backgrounds to lie on a 1-dimensional manifold. Progressively increasing the size of the latent code during training can be interpreted as gradually specializing from category-level to instance-level knowledge. Figure 2a illustrates the procedure with example results where we can observe the progressive specialization to particular instances: reactors gradually appear/disappear, textures get more accurate. Because common neural network implementations have fixed-size inputs, we implement progressive conditioning by masking, stage-by-stage, a decreasing number of values of the latent code. All our experiments share the same 4-stage training strategy where the latent code dimension is increased at the beginning of each stage and the network is then trained until convergence. We use dimensions 0/2/8/64 for the shape code, 2/8/64/512 for the texture code and 4/8/64/256 for the background code. We provide real-image results for each stage in our supplementary material.

Neighbor reconstruction. The idea behind *neighbor reconstruction* is to explicitly enforce consistency between different instances. Our key assumption is that neighboring instances with similar shape or texture exist in the dataset. If such neighbors are correctly identified, switching their shape or texture in our generation model should give similar reconstruction results. For a given input image, we hence propose to use its shape or texture attribute in the image formation process of neighboring instances and apply our reconstruction loss

on associated renderings. Intuitively, this process can be seen as mimicking a multi-view supervision without actually having access to multi-view images by finding neighboring instances in well-designed latent spaces. Figure 2b illustrates the procedure with an example.

More specifically, let $\{\mathbf{z}_{\text{sh}}, \mathbf{z}_{\text{tx}}, \mathbf{a}, \mathbf{z}_{\text{bg}}\}$ be the parameters predicted by our encoder for a given input training image \mathbf{I} , let Ω be a memory bank storing the images and parameters of the last M instances processed by the network. We write $\Omega^{(m)} = \{\mathbf{I}^{(m)}, \mathbf{z}_{\text{sh}}^{(m)}, \mathbf{z}_{\text{tx}}^{(m)}, \mathbf{a}^{(m)}, \mathbf{z}_{\text{bg}}^{(m)}\}$ each of these M instances and associated parameters. We first select the closest instance from the memory bank Ω in the texture (respectively shape) code space using the L_2 distance, $m_t = \operatorname{argmin}_m \|\mathbf{z}_{\text{tx}} - \mathbf{z}_{\text{tx}}^{(m)}\|_2$ (respectively $m_s = \operatorname{argmin}_m \|\mathbf{z}_{\text{sh}} - \mathbf{z}_{\text{sh}}^{(m)}\|_2$). We then swap the codes and generate the reconstruction $\hat{\mathbf{I}}_{\text{tx}}^{(m_t)}$ (respectively $\hat{\mathbf{I}}_{\text{sh}}^{(m_s)}$) using the parameters $\{\mathbf{z}_{\text{sh}}^{(m_t)}, \mathbf{z}_{\text{tx}}, \mathbf{a}^{(m_t)}, \mathbf{z}_{\text{bg}}^{(m_t)}\}$ (respectively $\{\mathbf{z}_{\text{sh}}, \mathbf{z}_{\text{tx}}^{(m_s)}, \mathbf{a}^{(m_s)}, \mathbf{z}_{\text{bg}}^{(m_s)}\}$). Finally, we compute the reconstruction loss between the images $\mathbf{I}^{(m_t)}$ and $\hat{\mathbf{I}}_{\text{tx}}^{(m_t)}$ (respectively $\mathbf{I}^{(m_s)}$ and $\hat{\mathbf{I}}_{\text{sh}}^{(m_s)}$). Our full loss can thus be written:

$$\mathcal{L}_{\text{nbr}} = \mathcal{L}_{\text{rec}}(\mathbf{I}^{(m_t)}, \hat{\mathbf{I}}_{\text{tx}}^{(m_t)}) + \mathcal{L}_{\text{rec}}(\mathbf{I}^{(m_s)}, \hat{\mathbf{I}}_{\text{sh}}^{(m_s)}). \quad (2)$$

Note that we recompute the parameters of the selected instances with the current network state, to avoid uncontrolled effects of changes in the network state. Also note that, for computational reasons, we do not use this loss in the first stage where codes are almost the same for all instances.

To prevent latent codes from specializing by viewpoint, we split the viewpoints into V bins w.r.t. the rotation angle, uniformly sample a target bin for each input and look for the nearest instances only in the subset of instances within the target viewpoint range (see the supplementary for details). In all experiments, we use $V = 5$ and a memory bank of size $M = 1024$.

3.3 Alternate 3D and pose learning

Because predicting 6D poses is hard due to self-occlusions and local minima, we follow prior works [12, 16, 23, 50] and predict multiple pose candidates and their likelihood. However, we identify failure modes in their optimization framework (detailed in our supplementary material) and instead propose a new optimization that alternates between 3D and pose learning. More specifically, given an input image \mathbf{I} , we predict K pose candidates $\{(\mathbf{r}_1, \mathbf{t}_1), \dots, (\mathbf{r}_K, \mathbf{t}_K)\}$, and their associated probabilities $\mathbf{p}_{1:K}$. We render the model from the different poses, yielding K reconstructions $\hat{\mathbf{I}}_{1:K}$. We then alternate the learning between two steps: (i) the *3D-step* where shape, texture and background branches of the network are updated by minimizing $\mathcal{L}_{3\text{D}}$ using the pose associated to the highest probability, and (ii) the *P-step* where the branches of the network predicting candidate poses and their associated probabilities are updated by minimizing:

$$\mathcal{L}_{\text{P}} = \sum_k \mathbf{p}_k \mathcal{L}_{\text{rec}}(\mathbf{I}, \hat{\mathbf{I}}_k) + \lambda_{\text{uni}} \mathcal{L}_{\text{uni}}, \quad (3)$$

where \mathcal{L}_{rec} is the reconstruction loss described in Sec. 3.2, \mathcal{L}_{uni} is a regularization loss on the predicted poses and λ_{uni} is a scalar hyperparameter. More precisely,

we use $\mathcal{L}_{\text{uni}} = \sum_k |\bar{\mathbf{p}}_k - 1/K|$ where $\bar{\mathbf{p}}_k$ is the averaged probabilities for candidate k in a particular training batch. Similar to [16], we found it crucial to introduce this regularization term to encourage the use of all pose candidates. In particular, this prevents a collapse mode where only few pose candidates are used. Note that we do not use the neighbor reconstruction loss nor the mesh regularization loss which are not relevant for viewpoints. In all experiments, we use $\lambda_{\text{uni}} = 0.02$.

Inspired by the camera multiplex of [12], we parametrize rotations with the classical Euler angles (azimuth, elevation and roll) and rotation candidates correspond to offset angles w.r.t. reference viewpoints. Since in practice elevation has limited variations, our reference viewpoints are uniformly sampled along the azimuth dimension. Note that compared to [12], we do not directly optimize a set of pose candidates per training image, but instead learn a set of K predictors for the entire dataset. We use $K = 6$ in all experiments.

4 Experiments

We validate our approach in two standard setups. It is first quantitatively evaluated on ShapeNet where state-of-the-art methods use multiple views as supervision. Then, we compare it on standard real-image benchmarks and demonstrate its applicability to more complex datasets. Finally, we present an ablation study.

4.1 Evaluation on the ShapeNet benchmark

We compare our approach to state-of-the-art methods using multi-views, viewpoints and silhouettes as supervision. Our method is instead learned without supervision. For all compared methods, one model is trained per class. We adhere to community standards [30, 34, 43] and use the renderings and splits from [30] of the ShapeNet dataset [4]. It corresponds to a subset of 13 classes of 3D objects, each object being rendered into a 64×64 image from 24 viewpoints uniformly spaced along the azimuth dimension. We evaluate all methods using the standard Chamfer- L_1 distance [36, 43], where predicted shapes are pre-aligned using our gradient-based version of the Iterative Closest Point (ICP) [3] with anisotropic scaling. Indeed, compared to competing methods having access to the ground-truth viewpoint during training, we need to predict it for each input image in addition to the 3D shape. This yields to both shape/pose ambiguities (*e.g.*, a small nearby object or a bigger one far from the camera) and small misalignment errors that dramatically degrade the performances. We provide evaluation details as well as results without ICP in our supplementary.

We report quantitative results and compare to the state of the art in Table 2, where methods using multi-views are visually separated. We evaluate the pre-trained weights for SDF-SRN [33] and train the models from scratch using the official implementation for DVR [43]. We tried evaluating SMR [21] but could not reproduce the results. We do not compare to TARS [8] which is based on SDF-SRN and share the same performances. Our approach achieves results that are on average better than the state-of-the-art methods supervised with silhouette

Method	Ours	SDF-SRN [33]	DVR [43]	DVR [43]
MV				✓
CK		✓	✓	✓
S		✓	✓	✓
airplane	0.110	0.128	0.114	0.111
bench	0.159	-	0.255	0.176
cabinet	0.137	-	0.254	0.158
car	0.168	0.150	0.203	0.153
chair	0.253	0.262	0.371	0.205
display	0.220	-	0.257	0.163
lamp	0.523	-	0.363	0.281
phone	0.127	-	0.191	0.076
rifle	0.097	-	0.130	0.083
sofa	0.192	-	0.321	0.160
speaker	0.224	-	0.312	0.215
table	0.243	-	0.303	0.230
vessel	0.155	-	0.180	0.151
mean	0.201	-	0.250	0.166

Table 2: **ShapeNet comparison.** We report Chamfer- L_1 ↓ obtained after ICP, **best** results are highlighted. Supervisions are: **M**ulti-**V**iews, **C**amera or **K**eypoints, **S**ilhouettes.

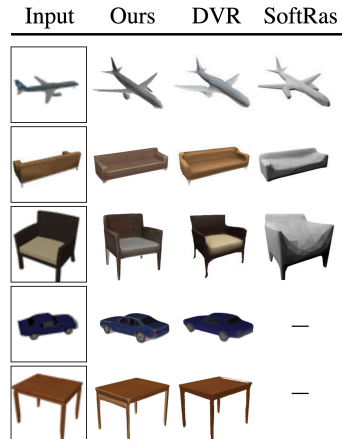


Fig. 5: **Visual comparisons.** We compare to DVR [43] and SoftRas [34] learned with full supervision (**MV**, **CK**, **S**).

and viewpoint annotations. This is a strong result: while silhouettes are trivial in this benchmark, learning without viewpoint annotations is extremely challenging as it involves solving the pose estimation and shape reconstruction problems simultaneously. For some categories, our performances are even better than DVR supervised with multiple views. This shows that our system learned on raw images generates 3D reconstructions comparable to the ones obtained from methods using geometry cues like silhouettes and multiple views. Note that for the lamp category, our method predicts degenerate 3D shapes; we hypothesize this is due to their rotation invariance which makes the viewpoint estimation ambiguous.

We visualize and compare the quality of our 3D reconstructions in Figure 5. The first three examples correspond to examples advertised in DVR [43], the last two corresponds to examples we selected. Our method generates textured meshes of high-quality across all these categories. The geometry obtained is sharp and accurate, and the predicted texture mostly corresponds to the input.

4.2 Results on real images

Pascal3D+ Car and CUB benchmarks. We compare our approach to state-of-the-art SVR methods on real images, where multiple views are not available. All competing methods use silhouette supervision and output meshes that are symmetric. CMR [25] additionally use keypoints, UCMR [12] and IMR [50] starts learning from a given template shape; we *do not* use any of these and directly learn from raw images. We strictly follow the community standards [12, 25, 50] and use the train/test splits of Pascal3D+ Car [57] (5000/220 images) and CUB-200-2011 [54] (5964/2874 images). Images are square-cropped around the object using bounding box annotations and resized to 64×64 .

Method	Supervision			Pascal3D+ Car			CUB-200-2011	
	CK	S	A	3D IoU \uparrow	Ch- L_1 \downarrow	Mask IoU \uparrow	PCK@0.1 \uparrow	Mask IoU \uparrow
CMR [25]	✓	✓	✓	64	-	-	48.3	70.6
IMR [50]		✓	✓	-	-	-	53.5	-
UMR [32]		✓	✓	62	-	-	58.2	73.4
UCMR [12]		✓	✓	67.3	0.172	73.7	-	63.7
SMR [21]		✓	✓	-	-	-	62.2	80.6
Ours				65.9	0.163	83.9	49.0	71.4

Table 3: **Real-image quantitative comparisons.** Supervision corresponds to **C**amera or **K**eypoints, **S**ilhouettes, **A**ssumptions (see Tab. 1 for details).

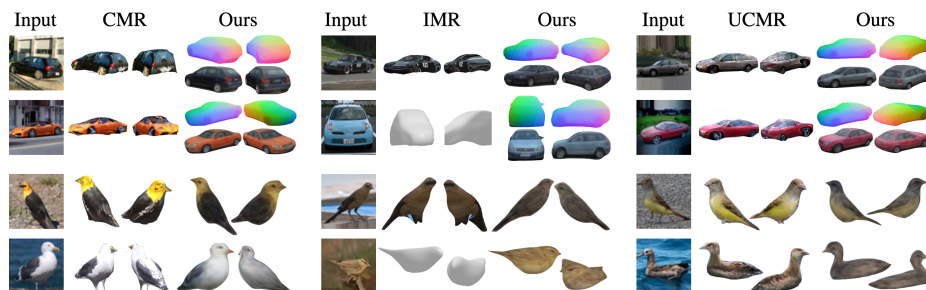


Fig. 6: **Real-image comparisons.** We show reconstructions on Pascal3D+ Cars (top) and CUB (bottom) and compare to CMR [25], IMR [50], UCMR [12].

A quantitative comparison is summarized in Table 3, where we report 3D IoU, Chamfer- L_1 (with ICP alignment), Mask IoU for Pascal3D+ Car, and Percentage of Correct Keypoints thresholded at $\alpha = 0.1$ (PCK@0.1) [31], Mask IoU for CUB. Our approach yields competitive results across all metrics although it does not rely on any supervision used by other works. On Pascal3D+ Car, we achieve significantly better results than UCMR for Chamfer- L_1 and Mask IoU, which we argue are less biased metrics than the standard 3D IoU [25, 51] computed on unaligned shapes (see supplementary). On CUB, our approach achieves reasonable results that are however slightly worse than the state of the art. We hypothesize this is linked to our pose regularization term encouraging the use of all viewpoints whereas these bird images clearly lack back views.

We qualitatively compare our approach to the state of the art in Figure 6. For each input, we show the mesh rendered from two viewpoints. For our car results, we additionally show meshes with synthetic textures to emphasize correspondences. Qualitatively, our approach yields results on par with prior works both in terms of geometric accuracy and overall realism. Although the textures obtained in [50] look more accurate, they are modeled as pixel flows, which has a clear limitation when synthesizing unseen texture parts. Note that we do not recover details like the bird legs which are missed by prior works due to coarse silhouette annotations. We hypothesize we also miss them because they are hardly consistent across instances, *e.g.*, legs can be bent in multiple ways.

Real-world datasets. Motivated by 3D-aware image synthesis methods learned in-the-wild [41, 42], we investigate whether our approach can be applied to real-world datasets where silhouettes are not available and images are not methodically

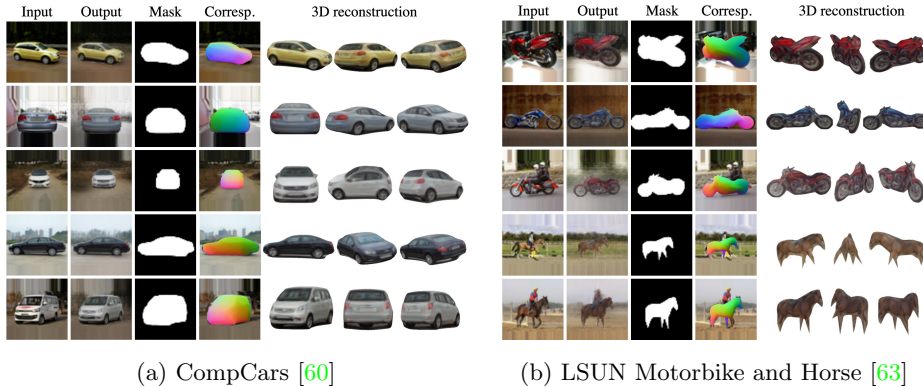


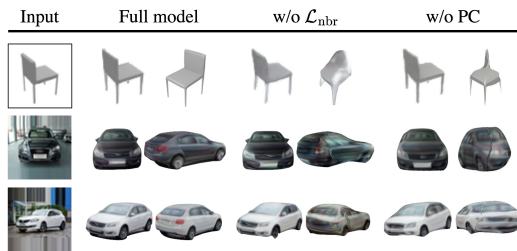
Fig. 7: **Real-world dataset results.** From left to right, we show: input and output images, the predicted mask, a correspondence map and the mesh rendered from 3 viewpoints. Note that for LSUN Horse, the geometry quality is low and outlines our approach limitations (see text). Best viewed digitally.

cropped around the object. We adhere to standards from the 3D-aware image synthesis community [41, 42] and apply our approach to 64×64 images of CompCars [60]. In addition, we provide results for the more difficult scenario of LSUN images [63] for motorbikes and horses. Because many LSUN images are noise, we filter the datasets as follows: we manually select 16 reference images with different poses, find the nearest neighbors from the first 200k images in a pre-trained ResNet-18 [15] feature space, and keep the top 2k for each reference image. We repeat the procedure with flipped reference images yielding 25k images.

Our results are shown in Figure 7. For each input image, we show from left to right: the output image, the predicted mask, a correspondence map, and the 3D reconstruction rendered from the predicted viewpoint and two other viewpoints. Although our approach is trained to synthesize images, these are all natural by-products. While the quality of our 3D car reconstructions is high, the reconstructions obtained for LSUN images lack some realism and accuracy (especially for horses), thus outlining limitations of our approach. However, our segmentation and correspondence maps emphasize our system ability to accurately localize the object and find correspondences, even when the geometry is coarse.

Limitations. Even if our approach is a strong step towards generic unsupervised SVR, we can outline three main limitations. First, the lack of different views in the data harms the results (*e.g.*, most CUB birds have concave backs); this can be linked to our uniform pose regularization term which is not adequate in these cases. Second, complex textures are not predicted correctly (*e.g.*, motorbikes in LSUN). Although we argue it could be improved by more advanced autoencoders, the neighbor reconstruction term may prevent unique textures to be generated. Finally, despite its applicability to multiple object categories and diverse datasets, our multi-stage progressive training is cumbersome and an automatic way of progressively specializing to instances is much more desirable.

Model	Full	w/o \mathcal{L}_{nbr}	w/o PC
airplane	0.110	0.124	0.107
bench	0.159	0.188	0.206
car	0.168	0.179	0.173
chair	0.253	0.319	0.527
table	0.243	0.246	0.598
mean	0.187	0.211	0.322

Table 4: **Ablation results on ShapeNet [4].**Fig. 8: **Ablation visual results.** For each input, we show the mesh rendered from two viewpoints.

4.3 Ablation study

We analyze the influence of our neighbor reconstruction loss \mathcal{L}_{nbr} and progressive conditioning (PC) by running experiments without each component.

First, we provide quantitative results on ShapeNet in Table 4. When \mathcal{L}_{nbr} is removed, the results are worse for almost all categories, outlining that it is important to the predicted geometry accuracy. When PC is removed, results are comparable to the full model for airplane and car but much worse for chair and table. Indeed, they involve more complex shapes and our system falls into a bad minimum with degenerate solutions, a scenario that is avoided with PC.

Second, we perform a visual comparison on ShapeNet and CompCars examples in Figure 8. For each input, we show the mesh rendered from the predicted viewpoint and a different viewpoint. When \mathcal{L}_{nbr} is removed, we observe that the reconstruction seen from the predicted viewpoint is correct but it is either wrong for chairs and degraded for cars when seen from the other viewpoint. Indeed, the neighbor reconstruction explicitly enforces the unseen reconstructed parts to be consistent with other instances. When PC is removed, we observe degenerate reconstructions where the object seen from a different viewpoint is not realistic.

5 Conclusion

We presented UNICORN, an unsupervised SVR method which, in contrast to all prior works, learns from raw images only. We demonstrated it yields high-quality results for diverse shapes as well as challenging real-world image collections. This was enabled by two key contributions aiming at leveraging consistency across different instances: our *progressive conditioning* training strategy and *neighbor reconstruction* loss. We believe our work includes both an important step forward for unsupervised SVR and the introduction of a valuable conceptual insight.

Acknowledgements. We thank François Darmon for inspiring discussions; Robin Champenois, Romain Loiseau, Elliot Vincent for feedback on the manuscript; and Michael Niemeyer, Shubham Goel for details on the evaluation. This work was supported in part by ANR project EnHerit ANR-17-CE23-0008, project Rapid Tabasco, gifts from Adobe and HPC resources from GENCI-IDRIS (2021-AD011011697R1, 2022-AD011013538).

References

1. Bengio, S., Vinyals, O., Jaitly, N., Shazeer, N.: Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks. In: NIPS (2015) 5
2. Bengio, Y., Louradour, J., Collobert, R., Weston, J.: Curriculum learning. In: ICML (2009) 5
3. Besl, P., McKay, N.D.: A method for registration of 3-D shapes. TPAMI 14(2) (1992) 10
4. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., Xiao, J., Yi, L., Yu, F.: ShapeNet: An Information-Rich 3D Model Repository. arXiv:1512.03012 [cs] (2015) 3, 10, 14
5. Chen, W., Gao, J., Ling, H., Smith, E.J., Lehtinen, J., Jacobson, A., Fidler, S.: Learning to Predict 3D Objects with an Interpolation-based Differentiable Renderer. In: NeurIPS (2019) 2, 4, 8
6. Choy, C.B., Xu, D., Gwak, J., Chen, K., Savarese, S.: 3D-R2N2: A Unified Approach for Single and Multi-view 3D Object Reconstruction. In: ECCV (2016) 2, 4
7. Desbrun, M., Meyer, M., Schröder, P., Barr, A.H.: Implicit fairing of irregular meshes using diffusion and curvature flow. In: SIGGRAPH (1999) 8
8. Duggal, S., Pathak, D.: Topologically-Aware Deformation Fields for Single-View 3D Reconstruction. In: CVPR (2022) 2, 4, 10
9. Elman, J.L.: Learning and development in neural networks: The importance of starting small. Cognition (1993) 5, 8
10. Finger, S.: Origins of neuroscience: a history of explorations into brain function. Oxford University Press (1994) 2
11. Gadelha, M., Maji, S., Wang, R.: 3D Shape Induction from 2D Views of Multiple Objects. In: 3DV (2017) 4
12. Goel, S., Kanazawa, A., Malik, J.: Shape and Viewpoint without Keypoints. In: ECCV (2020) 2, 4, 8, 9, 10, 11, 12
13. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative Adversarial Nets. In: NIPS (2014) 4
14. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: AtlasNet: A Papier-Mâché Approach to Learning 3D Surface Generation. In: CVPR (2018) 2, 4, 6
15. He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. In: CVPR (2016) 13
16. Henderson, P., Ferrari, V.: Learning single-image 3D reconstruction by generative modelling of shape, pose and shading. IJCV (2019) 2, 4, 9, 10
17. Henderson, P., Tsiminaki, V., Lampert, C.H.: Leveraging 2D Data to Learn Textured 3D Mesh Generation. In: CVPR (2020) 2, 4
18. Henzler, P., Mitra, N., Ritschel, T.: Escaping Plato’s Cave: 3D Shape From Adversarial Rendering. In: ICCV (2019) 4
19. Hoiem, D., Efros, A.A., Hebert, M.: Geometric Context from a Single Image. In: ICCV (2005) 2
20. Hoiem, D., Efros, A.A., Hebert, M.: Putting Objects in Perspective. IJCV (2008) 2
21. Hu, T., Wang, L., Xu, X., Liu, S., Jia, J.: Self-Supervised 3D Mesh Reconstruction From Single Images. In: CVPR (2021) 2, 4, 5, 10, 12
22. Ilg, E., Mayer, N., Saikia, T., Keuper, M., Dosovitskiy, A., Brox, T.: FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks. In: CVPR (2017) 5
23. Insafutdinov, E., Dosovitskiy, A.: Unsupervised Learning of Shape and Pose with Differentiable Point Clouds. In: NIPS (2018) 2, 4, 9

24. Jojic, N., Frey, B.J.: Learning Flexible Sprites in Video Layers. In: CVPR (2001) [3](#)
25. Kanazawa, A., Tulsiani, S., Efros, A.A., Malik, J.: Learning Category-Specific Mesh Reconstruction from Image Collections. In: ECCV (2018) [2](#), [4](#), [6](#), [11](#), [12](#)
26. Kar, A., Tulsiani, S., Carreira, J., Malik, J.: Category-Specific Object Reconstruction from a Single Image. In: CVPR (2015) [2](#), [4](#)
27. Karras, T., Laine, S., Aila, T.: A Style-Based Generator Architecture for Generative Adversarial Networks. In: CVPR (2019) [4](#)
28. Kato, H., Beker, D., Morariu, M., Ando, T., Matsuoka, T., Kehl, W., Gaidon, A.: Differentiable Rendering: A Survey. arXiv:2006.12057 [cs] (2020) [5](#)
29. Kato, H., Harada, T.: Learning View Priors for Single-view 3D Reconstruction. In: CVPR (2019) [4](#)
30. Kato, H., Ushiku, Y., Harada, T.: Neural 3D Mesh Renderer. In: CVPR (2018) [2](#), [4](#), [10](#)
31. Kulkarni, N., Gupta, A., Tulsiani, S.: Canonical Surface Mapping via Geometric Cycle Consistency. In: ICCV (2019) [5](#), [12](#)
32. Li, X., Liu, S., Kim, K., De Mello, S., Jampani, V., Yang, M.H., Kautz, J.: Self-supervised Single-view 3D Reconstruction via Semantic Consistency. In: ECCV (2020) [2](#), [4](#), [12](#)
33. Lin, C.H., Wang, C., Lucey, S.: SDF-SRN: Learning Signed Distance 3D Object Reconstruction from Static Images. In: NeurIPS (2020) [2](#), [4](#), [10](#), [11](#)
34. Liu, S., Li, T., Chen, W., Li, H.: Soft Rasterizer: A Differentiable Renderer for Image-based 3D Reasoning. In: ICCV (2019) [2](#), [3](#), [4](#), [5](#), [7](#), [8](#), [10](#), [11](#)
35. Loper, M.M., Black, M.J.: OpenDR: An Approximate Differentiable Renderer. In: ECCV 2014, vol. 8695 (2014) [4](#)
36. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy Networks: Learning 3D Reconstruction in Function Space. In: CVPR (2019) [2](#), [4](#), [10](#)
37. Monnier, T., Groueix, T., Aubry, M.: Deep Transformation-Invariant Clustering. In: NeurIPS (2020) [5](#), [8](#)
38. Monnier, T., Vincent, E., Ponce, J., Aubry, M.: Unsupervised Layered Image Decomposition Into Object Prototypes. In: ICCV (2021) [3](#)
39. Navaneet, K.L., Mathew, A., Kashyap, S., Hung, W.C., Jampani, V., Babu, R.V.: From Image Collections to Point Clouds with Self-supervised Shape and Pose Networks. In: CVPR (2020) [5](#)
40. Nealen, A., Igarashi, T., Sorkine, O., Alexa, M.: Laplacian mesh optimization. In: GRAPHITE (2006) [8](#)
41. Nguyen-Phuoc, T., Li, C., Theis, L., Richardt, C., Yang, Y.L.: HoloGAN: Unsupervised learning of 3D representations from natural images. In: ICCV (2019) [12](#), [13](#)
42. Niemeyer, M., Geiger, A.: GIRAFFE: Representing Scenes as Compositional Generative Neural Feature Fields. In: CVPR (2021) [12](#), [13](#)
43. Niemeyer, M., Mescheder, L., Oechsle, M., Geiger, A.: Differentiable Volumetric Rendering: Learning Implicit 3D Representations without 3D Supervision. In: CVPR (2020) [2](#), [4](#), [10](#), [11](#)
44. Pavllo, D., Spinks, G., Hofmann, T., Moens, M.F., Lucchi, A.: Convolutional Generation of Textured 3D Meshes. In: NeurIPS (2020) [4](#)
45. Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3D Deep Learning with PyTorch3D. arXiv:2007.08501 [cs] (2020) [4](#), [7](#)
46. Saxena, A., Min Sun, Ng, A.: Make3D: Learning 3D Scene Structure from a Single Still Image. TPAMI (2009) [2](#)

47. Schroff, F., Kalenichenko, D., Philbin, J.: FaceNet: A unified embedding for face recognition and clustering. In: CVPR (2015) 5
48. Simonyan, K., Zisserman, A.: Very Deep Convolutional Networks for Large-Scale Image Recognition. In: ICLR (2015) 8
49. Tulsiani, S., Efros, A.A., Malik, J.: Multi-view Consistency as Supervisory Signal for Learning Shape and Pose Prediction. In: CVPR (2018) 2, 4
50. Tulsiani, S., Kulkarni, N., Gupta, A.: Implicit Mesh Reconstruction from Unannotated Image Collections. arXiv:2007.08504 [cs] (2020) 2, 4, 6, 9, 11, 12
51. Tulsiani, S., Zhou, T., Efros, A.A., Malik, J.: Multi-view Supervision for Single-view Reconstruction via Differentiable Ray Consistency. In: CVPR (2017) 2, 4, 12
52. Vicente, S., Carreira, J., Agapito, L., Batista, J.: Reconstructing PASCAL VOC. In: CVPR (2014) 2, 4
53. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2Mesh: Generating 3D Mesh Models from Single RGB Images. In: ECCV (2018) 2, 4, 5, 8
54. Welinder, P., Branson, S., Mita, T., Wah, C., Schroff, F., Belongie, S., Perona, P.: Caltech-UCSD Birds 200. Tech. rep., California Institute of Technology (2010) 3, 11
55. Wu, S., Makadia, A., Wu, J., Snavely, N., Tucker, R., Kanazawa, A.: De-rendering the World’s Revolutionary Artefacts. In: CVPR (2021) 2, 4
56. Wu, S., Rupprecht, C., Vedaldi, A.: Unsupervised Learning of Probably Symmetric Deformable 3D Objects from Images in the Wild. In: CVPR (2020) 2, 4, 8
57. Xiang, Y., Mottaghi, R., Savarese, S.: Beyond PASCAL: A benchmark for 3D object detection in the wild. In: WACV (2014) 3, 11
58. Xu, Q., Wang, W., Ceylan, D., Mech, R., Neumann, U.: DISN: Deep Implicit Surface Network for High-quality Single-view 3D Reconstruction. In: NeurIPS (2019) 4
59. Yan, X., Yang, J., Yumer, E., Guo, Y., Lee, H.: Perspective Transformer Nets: Learning Single-View 3D Object Reconstruction without 3D Supervision. In: NeurIPS (2016) 2, 4
60. Yang, L., Luo, P., Loy, C.C., Tang, X.: A large-scale car dataset for fine-grained categorization and verification. In: CVPR (2015) 3, 13
61. Yao, C.H., Hung, W.C., Jampani, V., Yang, M.H.: Discovering 3D Parts from Image Collections. In: ICCV (2021) 5
62. Ye, Y., Tulsiani, S., Gupta, A.: Shelf-Supervised Mesh Prediction in the Wild. arXiv:2102.06195 [cs] (2021) 4
63. Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., Xiao, J.: LSUN: Construction of a Large-scale Image Dataset using Deep Learning with Humans in the Loop. arXiv:1506.03365 [cs] (2016) 3, 13
64. Zhang, J.Y., Yang, G., Tulsiani, S., Ramanan, D.: NeRS: Neural Reflectance Surfaces for Sparse-view 3D Reconstruction in the Wild. In: NeurIPS (2021) 8
65. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In: CVPR (2018) 8
66. Zhang, Y., Chen, W., Ling, H., Gao, J., Zhang, Y., Torralba, A., Fidler, S.: Image GANs meet Differentiable Rendering for Inverse Graphics and Interpretable 3D Neural Rendering. In: ICLR (2021) 2, 4