# Novel View Synthesis from Single Images via Point Cloud Transformation

Hoang-An Le[1]
hoang-an.le@uva.nl

Thomas Mensink[2]
mensink@google.com

Partha Das[1]
p.das@uva.nl

Theo Gevers[1]
th.gevers@uva.nl

[1] Computer Vision Lab,
University of Amsterdam

[2] Google Research,
Amsterdam

arXiv:2009.08321v2 [cs.CV] 18 Sep 2020

## Abstract

In this paper the argument is made that for true novel view synthesis of objects, where the object can be synthesized from any viewpoint, an explicit 3D shape representation is desired. Our method estimates point clouds to capture the geometry of the object, which can be freely rotated into the desired view and then projected into a new image. This image, however, is sparse by nature and hence this coarse view is used as the input of an image completion network to obtain the dense target view. The point cloud is obtained using the predicted pixel-wise depth map, estimated from a single RGB input image, combined with the camera intrinsics. By using forward warping and backward warping between the input view and the target view, the network can be trained end-to-end without supervision on depth. The benefit of using point clouds as an explicit 3D shape for novel view synthesis is experimentally validated on the 3D ShapeNet benchmark. Source code and data are available at https://github.com/lhoangan/pc4novis

## 1 Introduction

Novel view synthesis aims to infer the appearance of an object from unobserved points of view. The synthesis of unseen views of objects could be important for image-based 3D object manipulation [18], robot traversability [12], or 3D object reconstruction [29]. Generating a coherent view of unseen parts of an object requires a non-trivial understanding of the object's inherent properties such as (3D) geometry, texture, shading, and illumination.

Different algorithms make use of provided source images in different ways. Model-based approaches use similar-look open stock 3D models [18], or through user interactive construction [2, 26, 31]. Image-based methods [23, 24, 28, 29, 32] assume an underlying parametric model of object appearances conditioned on viewpoints and try to learn it using statistical frameworks. Despite their differences, both approaches use 3D information in predicting object new views. The former imposes stronger assumptions on the full 3D structure and shifts the paradigm to obtain the full models, while the latter captures the 3D information in latent space to cope with (self) occlusion.
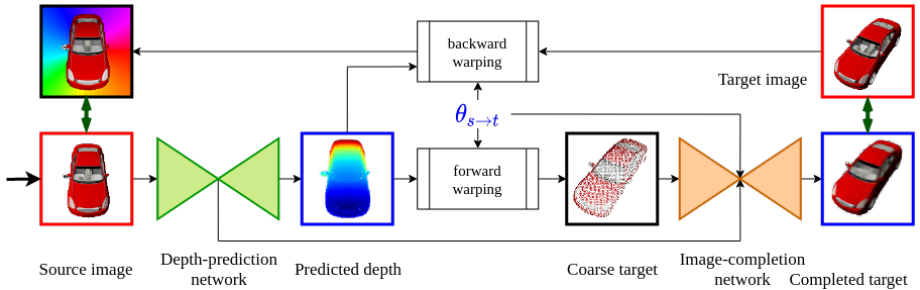
Figure 1: Overview of the proposed model for training and inference. From a single input image, the pixel-wise depth map is predicted. The depth map is subsequently used to compute a coarse novel view (forward warping), and trained by making use of backward warping (from the target view back to the source view). The model is trained end-to-end.

The principle is that the generation of a new view of an object is composed of (1) relocating pixels in source images that will be visible to the corresponding positions in the target view, (2) removing the pixels that will be occluded, and (3) adding disoccluded pixels that are not seen in the source and will be visible in the target view [24]. With the advance of convolution neural networks (CNNs) and generative adversarial networks (GANs), [24, 28, 32] show that (1) and (2) can be done by learning an appearance flow field that "flows" pixels from a source image to the corresponding positions in the target view, and (3) can be done by a completion network with an adversarial loss.

In this paper, we leverage the explicit use of geometry information in synthesizing novel views. We argue that (1) and (2) can be done in a straightforward manner by obtaining access to the geometry of the objects. The appearance flow [24, 28, 32] which associates pixels of the source view to their positions in the target view, is the projection of the 3D displacement of objects' points before and after transformation. Occluded object parts can be identified based on the orientation of the object surface normals and the view directions. The argument can also be extended for multiple input images. In this paper, we show that the geometry of an object provides an explicit and natural basis to the problem of novel view synthesis.

In contrast to geometry-based methods, the proposed approach does not require 3D supervision. The method predicts a depth map in a self-supervised manner by formulating the depth estimation problem in the context of novel view synthesis. The predicted depth is used to partly construct the target views and to assist the completion network.

The main contributions of this paper are: (1) a novel methodology for novel view synthesis using explicit transformations of estimated point clouds; (2) an integrated model combining self-supervised monocular depth estimation and novel view synthesis, which can be trained end-to-end; (3) natural extensions to multi-view inputs and full point cloud reconstruction from a single image; and (4) experimental benchmarking to validate the proposed method, which outperforms the current state-of-the art methods for novel view synthesis.

# 2    Related Work

## 2.1    Geometry-based view synthesis

**View synthesis via 3D models**     Full models (textured meshes or colored point clouds) of objects or scenes are constructed from multiple images taken from various viewpoints

[4, 20, 27] or are given and aligned interactively by users [18, 26]. The use of 3D models allows for extreme pose estimation, re-texturing and flexible (re-)lighting by applying rendering techniques [20, 22]. However, obtaining complete 3D models of objects or scenes is a challenging task in itself. Therefore, these approaches require additional user input to identify objects boundaries [2, 31], select and align 3D models with image views [18, 26], or use simple textured-mapped 3-planar billboard models [13]. In contrast, the proposed method makes use of objects partial point clouds constructed from a given source view and does not require a predefined (explicit) 3D model.

**View synthesis via depth** Methods using 3D models assume a coherent structure between the desired objects and the obtained 3D models [2, 18]. Synthesis using depth obtains an intermediate representation from depth information. The intermediate representation captures hidden surfaces from one or multiple viewpoints. [37] proposes to use layered depth images, [5] creates 3D plane sweep volumes by projecting images onto target viewpoints at different depths, [34] uses multi-plane images at fix-distances to the camera, and [3] estimates depth probability volumes to leverage depth uncertainty in occluded regions.

In contrast, the proposed method estimates depth directly from monocular views to partially construct the target views. Self-supervised depth estimation using deep neural networks using photometric re-projection consistency has been researched by several authors [7, 9, 10, 17, 33]. In this paper, we train a self-supervised depth prediction network with novel view synthesis in an end-to-end system.

## 2.2 Image-based view synthesis

Requiring explicit geometrical structures of objects or scenes as a precursor severely limits the applicability of a method. With the advance of neural networks (CNNs), generative adversarial networks [11] (GANs) achieve impressive results in image generation, allowing view synthesis without explicit geometrical structures of objects or scenes.

**View synthesis via embedded geometry** Zhou *et al.* [32] proposes learning a flow field that maps pixels in input images to their corresponding locations in target views to capture latent geometrical information. [23] learns a volumetric representation in a transformable bottleneck layer, which can generate corresponding views for arbitrary transformations. The former explicitly utilizes input (source) image pixels in constructing new views, either fully [32], or partly with the rest being filled by a completion network [24, 28]. The latter explicitly applies transformations on the volumetric representation in latent space and generates new views by means of pixel generation networks.

The proposed method takes the best of both worlds. By directly using object geometry the source pixels are mapped to their target positions based on given transformation parameters, hence making the best use of the given information synthesizing new views. Our approach is fundamentally different from [24]: we estimate the object point cloud using self-supervised depth predictions and obtain coarse target views from purely geometrical transformations, while [24] learns mappings from input images and ground truth occluded regions to generate coarse target views using one-hot encoded transformation vectors.

**View synthesis directly from image** Since the introduction of image-to-image translation [14], there is a paradigm shift towards pure image-based approaches [29]. [36] synthesizes bird view images from a single frontal view image, while [25] generates cross-views of aerial and street-view images. The networks can be trained to predict all the views in an orbit from a single-view object [17, 19], or generate a view in an iterative manner [6]. Additional features can be embedded such as view-independent intrinsic properties of objects [30]. In

this paper, we employ GANs to generate complete views, which is conditioned on the geometrical features and the relative poses between source and target views. Our approach can be interpreted as a reverse and end-to-end process of [17]: we estimate objects' arbitrary new views via point clouds constructed from self-supervised depth maps, while [17] predict objects' fixed orbit views for 3D reconstruction.

# 3    Method

## 3.1    Point-cloud based transformations

The core of the proposed novel view synthesis method is to use point clouds for geometrically aware transformations. Using the pinhole camera model and known intrinsics $\mathbf{K}$, the point cloud can be reconstructed when the pixel-wise depth map (D) is available. The camera intrinsics can be obtained by camera calibration, yet for the synthetic data used in our experiments, $\mathbf{K}$ is given. A pixel on the source image plane $p_s = [u, v, 1]$ (using homogeneous coordinates), corresponds to a point $P_s = [X, Y, Z]$ in the source camera space:

$$D_s \, p_s^\top = \mathbf{K} \, P_s^\top \qquad\qquad P_s^\top = \mathbf{K}^{-1} \, D_s \, p_s^\top \qquad (1)$$

Rigid transformations can be obtained by matrix multiplications. The relative transformation to the *target* viewpoint from the *source* camera, is given by

$$\theta_{s \to t} = \left[ \begin{array}{c|c} \mathbf{R} & \mathbf{t} \\ \hline 0 & 1 \end{array} \right] \qquad (2)$$

where $\mathbf{R}$ denotes the desired rotation matrix and $\mathbf{t}$ the translation vector. Points in the target camera view are given by $P_t = \theta_{s \to t} P_s$. This can also be regarded as an image-based flow field $\phi : \mathbb{R}^2 \to \mathbb{R}^2$ parameterized by $\theta_{s \to t}$ (c.f. [24, 28, 32]). The flow field $\phi(p_s; \theta_{s \to t})$ returns the homogeneous coordinates of pixels in the target image for each pixel in the source image:

$$\phi(p_s; \theta_{s \to t}) = \mathbf{K} \, \theta_{s \to t} \, \mathbf{K}^{-1} \, D_s p_s^\top \qquad (3)$$

By observing that $\phi(p_s; \theta_{s \to t}) = D_t p_t$, the Cartesian pixel coordinates in the target view can be extracted. The advantage of the flow field interpretation is that it provides a direct mapping between the image planes of the source view and the target view.

**Forward warping**    The flow field is used to generate the target view from the source:

$$\tilde{I}_t(\phi(p_s; \theta_{s \to t})) = I_s(p_s). \qquad (4)$$

The resulted image is sparse due to the discrete pixel coordinates and (dis)occluded regions, see Fig. 2 (*top-right*). It is used as input to the image completion network (Sec. 3.2).

**Backward warping**    The flow field is used to generate the source view from the target:

$$\tilde{I}_s(p_s) = I_t(\phi(p_s; \theta_{s \to t})). \qquad (5)$$

The process assigns a value to every pixel $(u, v)$ in $\tilde{I}_s$ resulting in a dense image, as illustrated in Fig. 2 (*bottom-right*). The generated source view may contain artifacts due to (dis)occlusion in the target view. To sample $\phi(p_s; \theta_{s \to t})$ from $I_t$, a differentiable bi-linear sampling layer [15] is used. The generated source view is used for self-supervised monocular depth prediction (Sec. 3.3).
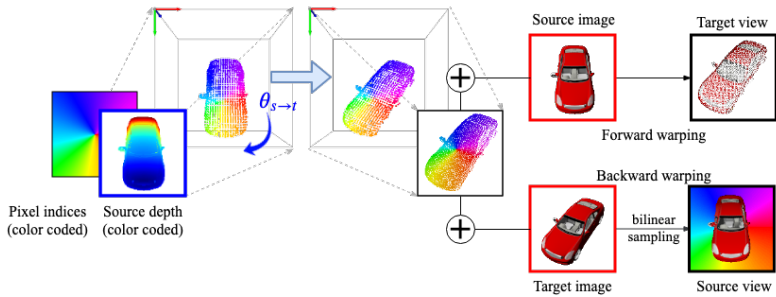
Figure 2: Illustration of the forward and backward warping operation of point clouds. The forward warping is used to generate a coarse target view, while the backward warping is used to reconstruct the source view from a target view for self-supervised depth estimation.

## 3.2 Novel view synthesis

The point-cloud-based forward warping relocates the visible pixels of the object in the source view to their corresponding positions in the target view. For novel view synthesis, however, two more steps are required: (1) obtaining the target coarse view by discarding occluded pixels, and (2) filling in the pixels that are not seen in the source view.

**Coarse view construction** The goal is to remove the pixels which are seen in the source view yet should not be visible in the target view, due to occlusion. To this end, pixels that have surface normals (after transformation) pointing away from the viewing direction are removed, similarly to [24]. Surface normals are obtained from normalized depth gradients.

An illustration of the coarse view construction is shown in Fig. 3 for different target views. The first row depicts the target views, the second row indicates the visible parts from the input image (third column). The third and fourth row show the coarse view with and without occlusion removal (or backface culling). Finally, the fifth row shows an enhanced version of the coarse view, where the object is assumed to be left-right symmetric [24]. The proposed method directly identifies and removes occlusion pixels from the input view using *estimated* depth, which contrasts to [24], where ground truth visibility mask are required for each target view to train a visibility prediction network.

**View completion** The obtained coarse view is already in the target viewpoint, but it remains sparse. To synthesize the final dense image, an image completion network is used.

The completion network uses the hour-glass architecture [21]. Following [24], we concatenate the depth bottleneck features and embedded transformation to the completion network bottleneck. By conditioning the completion network on the input features and the desired transformation $\theta_{s \to t}$, the network can fix artifacts and errors due to estimated depth and cope better with extreme pose transformations, *i.e.* when coarse view image is near empty (*e.g.* columns 9-11 in Fig. 3).

The image completion network is trained in a GAN-manner by using a generator $\mathcal{G}$, a discriminator $\mathcal{D}$, an input image $I_s$ and a target image $I_t$. The combination of losses that are used is given by

$$\mathcal{L}_{\mathcal{D}} = (\mathcal{D}(I_s) - 1)^2 + \mathcal{D}(\mathcal{G}(I_s))^2, \qquad \text{LS-GAN Discriminator loss} \qquad (6)$$

$$\mathcal{L}_{\mathcal{G}} = [1 - \text{SSIM}(I_t, \mathcal{G}(I_s))] + \|I_t - \mathcal{G}(I_s)\|_1, \qquad \text{Generator loss} \qquad (7)$$

$$\mathcal{L}_{Perc} = \left\|\mathcal{F}_{I_t}^{\mathcal{D}} - \mathcal{F}_{\mathcal{G}(I_s)}^{\mathcal{D}}\right\|_2 + \left\|\mathcal{F}_{I_t}^{\text{VGG}} - \mathcal{F}_{\mathcal{G}(I_s)}^{\text{VGG}}\right\|_2, \qquad \text{Perceptual loss} \qquad (8)$$
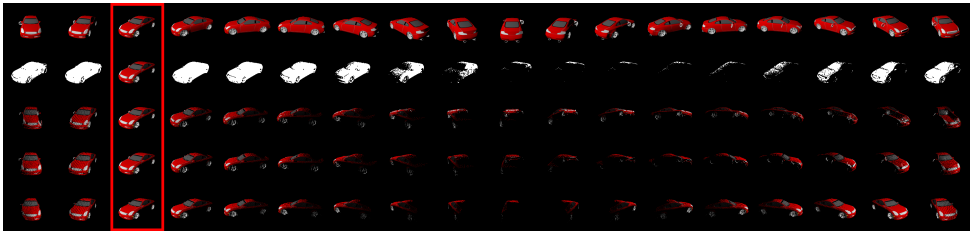
Figure 3: Image coarse views for different target viewpoints. The input image is depicted in the third column (red box). From top to bottom: (1) target views, (2) source region visible in each target viewpoint, coarse view (3) naive (4) with occlusion removal, and (5) with occlusion removal and symmetry.

where the perceptual loss uses $\mathcal{F}^{\mathcal{D}}$ ($\mathcal{F}^{\text{VGG}}$) to denote features extracted from image $I_t$ and $\mathcal{G}(I_s)$ from the discriminator network and pre-trained VGG network respectively, *c.f.* [16]. SSIM denotes the structural similarity index measure, see Sec. 4. The total loss is given by:

$$\mathcal{L}_c = w_1 \mathcal{L}_{\mathcal{D}} + w_2 \mathcal{L}_{\mathcal{G}} + w_3 \mathcal{L}_{Perc}, \tag{9}$$

where $w$ denotes the weighting of the losses ($w_1 = 1, w_2 = 100$, and $w_3 = 100$, *c.f.* [24]).

## 3.3 Self-supervised Monocular Depth estimation

The discussion so far has assumed that pixel-wise depth maps are available. In this section, the method used to estimate depth from a single *RGB* image is detailed. In order to make the minimum assumption about the training data, self-supervised methods are considered, which do not require ground-truth depth [7, 9, 10, 17, 33].

For the depth prediction an encoder-decoder network with bottleneck architecture is used, similar to [10]. The network is optimised using a set of (reconstruction) losses between the source image $I_s$ and its synthesized version $\tilde{I}_s$, using the backward warping, Eq. (5), from a second (target) image $I_t$ and the predicted depth map. The underlying rationale is that a more realistic depth map will have a lower reconstruction loss.

The losses are given by:

$$\mathcal{L}_p\left(I_s, \tilde{I}_s\right) = \frac{\alpha}{2}\left[1 - \text{SSIM}\left(I_s, \tilde{I}_s\right)\right] + (1-\alpha)\left\|I_s - \tilde{I}_s\right\|_1, \quad \text{Photometric loss} \tag{10}$$

$$\mathcal{L}_s(d) = \left|\partial_x d\right| e^{-\left|\partial_x I_s\right|} + \left|\partial_y d\right| e^{-\left|\partial_y I_s\right|}, \quad \text{Smoothness loss [9]} \tag{11}$$

$$\mathcal{L}_d\left(I_s, \tilde{I}_s, d\right) = \mu \mathcal{L}_p\left(I_s, \tilde{I}_s\right) + w_d \mathcal{L}_s, \quad \text{Total loss} \tag{12}$$

where $\alpha = 0.85$, and $d = \frac{\bar{D}}{D}$ is the mean-normalized inverse depth, $w_d = 10^{-3}$, and $\mu$ is an indicator function which equals 1 iff the photometric loss $\mathcal{L}_p(I_s, \tilde{I}_s) < \mathcal{L}_p(I_s, I_t)$, see [10] for more details. The smoothness loss encourages nearby pixels to have similar depths, while the artifacts due to (dis)occlusion are excluded by the per-pixel minimum-projection mechanism.

Table 1: Quality of coarse and completed view (*a*) and ablation study (*b*). The proposed transformation within estimated point clouds generates better coarse images over image-based predicted flow fields. Subsequently, the completed view quality is improved. The ablation study shows the best results are obtained by using LSGAN (LS), perceptual loss (PL), symmetry (Sym), bottleneck inter-connection (IC), and SSIM loss (SL).

(a) Coarse vs Completed

|  | Coarse view | | Completed view | |
|---|---|---|---|---|
|  | L1 ($\downarrow$) | SSIM ($\uparrow$) | L1 ($\downarrow$) | SSIM ($\uparrow$) |
| DOAFN [24] | .220 | .876 | .121 | .910 |
| M2NV [28] | .226 | .879 | .154 | .906 |
| Ours | **.203** | **.882** | **.118** | **.924** |

(b) Ablation study

| LS | PL | Sym | IC | SL | L1 ($\downarrow$) | SSIM ($\uparrow$) |
|---|---|---|---|---|---|---|
|  | ✓ | ✓ |  |  | .118 | .924 |
| ✓ | ✓ | ✓ |  |  | .101 | .939 |
| ✓ |  | ✓ |  |  | .103 | .939 |
| ✓ | ✓ |  |  |  | .107 | .933 |
| ✓ | ✓ | ✓ | ✓ |  | **.097** | .939 |
| ✓ | ✓ | ✓ | ✓ | ✓ | **.097** | .942 |

# 4 Experiments

In this section, the proposed method is analysed on the 3D ShapeNet benchmark including an ablation to study the effects of the different components and a state-of-the-art comparison.

**Dataset** We use the object-centered car and chair images rendered from the 3D ShapeNet models [1] using the same render engine[1] and set up as in [23, 24, 28, 32]. Specifically, there are 7497 car and 698 chair models with high-quality textures, split by 80%/20% for training and test. The images are rendered at 18 azimuth angles (in $[0, 340]$, $20°$-separation) and 3 elevation angles ($0°, 10°, 20°$). Input and output images are of size $256 \times 256$.

**Metrics** We evaluate the generated images using the standard $L_1$ pixel-wise error (normalized to the ranged $[0, 1]$, lower is better) and the structural similarity index measure (SSIM) [35] (value range of $[-1, 1]$, higher is better). $L_1$ indicates the proximity of pixel values between a completed image and the target, while SSIM measures the perceived quality and structural similarity between the images.

**Baseline** We compare the results of our method with the following state-of-the-art methods: AFN [32], TVSN [24], M2NV [28], and TBN [23].

## 4.1 Initial experiments

**Comparison to image-based completion** In this section, we compare the intermediate views generated by the forward warping using estimated point clouds and those by image-based flow field prediction by DOAFN [24] and M2NV [28]. For this experiment, the coarse view after occlusion removal and left-right symmetric enhancements are used. The image completion network is the basis variant, using DCGAN, without bottleneck inter-connections. The results are shown in Table 1(a). The transformation of estimated point clouds provides coarse views which are closer to the target view, and these help to obtain a higher quality of completed views.

**Ablation Study** We analyze the effects of the different component of the proposed pipeline. The results are shown in Table 1(b). The use of the LSGAN loss shows a relative large improvement over the traditional DCGAN. The drop of performance by removing symmetry

---

[1]The specific render engine and setup is to guarantee fair comparison with reported methods as none of the authors-provided weights perform at the similar level on images rendered with different rendering setups.

| Methods | cars | | chairs | |
|---|---|---|---|---|
| | L1 ($\downarrow$) | SSIM ($\uparrow$) | L1 ($\downarrow$) | SSIM ($\uparrow$) |
| *Same elevation* | | | | |
| AFN [32] | .148 | .877 | .229 | .871 |
| TVSN [24] | .119 | .913 | .202 | .889 |
| M2VN [28] | .098 | .923 | .181 | .895 |
| TBN [23] | .091 | .927 | .178 | .895 |
| ours | **.096** | **.945** | **.175** | **.914** |
| *Cross-elevation* | | | | |
| TBN | .199 | .910 | .215 | .902 |
| ours | **.122** | **.934** | **.207** | **.905** |

Table 2: Quantitative comparison with state-of-the-art methods on novel view synthesis: our method consistently performs (slightly) better than the other methods for both categories where target views have the same or different elevation angles with input views.

assumption shows the importance of prior knowledge on target objects, which is intuitive. The inter-connection from the depth network and the embedded transformation to the completion network allow the model to not rely solely on intermediate views. This is important for overcoming errors and artifacts which occur in the coarse images (due to inevitable uncertainties in depth prediction) and generate in general higher quality images. The SSIM loss, first employed by [23], shows improvement in SSIM metric, which is intuitive as training objectives are closer to evaluation metrics.

## 4.2    Comparison to State-of-the-Art

In this section, the proposed method is compared with state-of-the-art methods. The quantitative results are shown in Table 2. The proposed method performs consistently performs (slightly) better on both evaluation metrics for both types of objects. Quantitative results are shown in Fig. 4 where challenging cases are shown in the last 2 rows. Notice the better ability in retaining objects' textures (such as color patterns and texts on cars) of methods that explicitly use input pixel values in generating new views to that of TBN. The results of cars are constantly higher than that of chairs due to the intricate structures of chairs. However, by having access to object geometry, geometrical assumptions such as symmetry and occlusion can be applied directly to intermediate views (instead of having to learn from annotated data *c.f.* [24]), which creates better views for near-to symmetry targets. High-quality qualitative results and more analyses can be found in the supplementary materials.

Table 2 also shows the evaluation when target viewpoints are from different elevation angles. Methods such AFN, TVSN, and M2NV encode transformation as one-hot vectors and thus, are limited to operate within a pre-defined set of transformations (18 azimuth angles, same elevation). This is not the case for our method and TBN which apply direct transformation. We use the same azimuth angles as in the standard test set while randomly sample new elevation angles for input images in $(0°, 10°, 20°)$. The results are shown with networks trained with the regular fixed-elevation settings. The new transformations produces different statistics from what the networks have been trained, resulting in a performance drop for both methods. Nevertheless, the proposed method can still maintain high quality image synthesis.
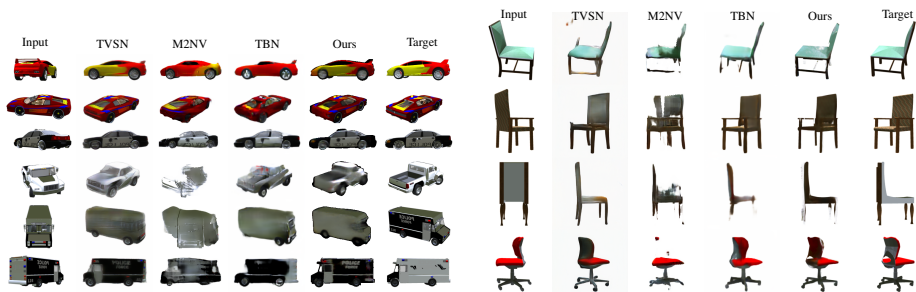
Figure 4: Qualitative comparisons of synthesized cars (*left*) and chairs (*right*), given a single input image (*first column*) and a given target view (*last column*). The last two rows show a more challenging examples. The proposed method captures better the geometry of the object and the fine (texture) details. More examples are provided in the supplementary materials.

## 4.3 Multi-View Synthesis and Point Cloud Reconstruction

**Multi-view inputs** The proposed method can be naturally extended to use multi-view inputs as follows: for each image depth is predicted independently and combined into a single point cloud. The resulting coarse target image will be denser when more images are used, and is passed through the image completion network.

In this experiment, the model trained for single-view prediction is used and evaluated using multiple (1 to 8) input images. The results in Table 3 show that the quality of the coarse view increases, as expected, when more input images are used and hence the point clouds are denser. Surprisingly, however, the image completion network only marginally improves, indicating that the coarse view contains enough information for the image completion network to synthesis a high quality target image.

**Point cloud reconstruction** In this final experiment, the aim is to reconstruct a full dense point cloud from a single image, using the models trained for novel view synthesis. In order to do so, 360°-views are generated from a single view of an object, see Fig. 5 (top). Each of these views are fed to the depth estimation network and the obtain estimated depth is used to generate a partial point cloud. These point clouds are stitched together, using corresponding transformations, resulting in a high quality dense point cloud, as shown in Fig. 5 (bottom).

| No. | Coarse | | Final | |
| views | L1 ($\downarrow$) | SSIM ($\uparrow$) | L1 ($\downarrow$) | SSIM ($\uparrow$) |
|---|---|---|---|---|
| 1 | .203 | .882 | .090 | .945 |
| 2 | .188 | .888 | .089 | .945 |
| 4 | .152 | .906 | .085 | .946 |
| 8 | **.111** | **.907** | **.084** | **.947** |

Table 3: Performance by extending single-view-trained networks for multi-view inputs.

## 4.4 Results on real-world imagery

We apply the trained car model to the car images of the real-imagery ALOI dataset [8], consisting of 100 objects, captured at 72 viewing angles. We use 4 cars for fine-tuning only the depth network, which requires no ground truths, while the image-completion network is left untouched. The quantivative resuls on the remaining 3 cars are shown in Fig. 6.
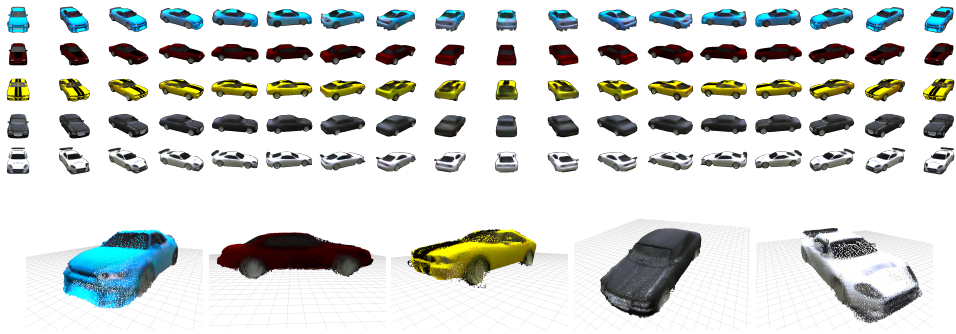
Figure 5: *Top* 360°-views generation from single input images (first column), *bottom* point cloud reconstruction using estimated depths of each generated views. Depth estimation trained on real images can perform well on synthesized ones.
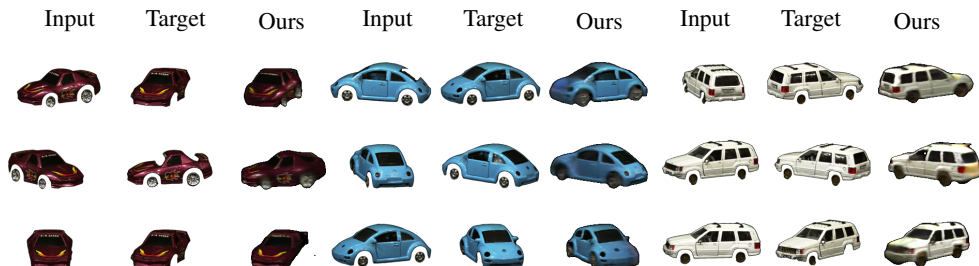


Figure 6: Quantitative results on real-imagery ALOI [8] dataset. The inputs and targets are shown with provided object masks, while the synthesized images are with predicted masks. The completion network does not need to be finetuned, yet provide competent results.

## 5    Conclusion

In this paper partial point clouds are estimated from a single image, by a self-supervised depth prediction network and used to obtain a coarse image in the target view. The final image is produced by an image completion network which uses the coarse image as input. Experimentally the proposed method outperforms any of the current SOTA methods on the ShapeNet Benchmark on novel view synthesis. Qualitative results show high quality and dense point clouds, obtained from a single image, by synthesizing and combining 360° views. Based on these results, we conclude that point clouds are a suitable, geometry aware representation for true novel view synthesis.

## References

[1] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical report, ArXiV:1512.03012, 2015. 7

[2] Tao Chen, Zhe Zhu, Ariel Shamir, Shi-Min Hu, and Daniel Cohen-Or. 3-Sweep: Extracting Editable Objects from a Single Photo. *ACM Trans. Graph.*, 32(6), 2013. 1, 3

[3] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme View Synthesis. In *ICCV*, 2019. 3

[4] Paul E Debevec, Camillo J Taylor, and Jitendra Malik. Modeling and Rendering Architecture from Photographs: A Hybrid Geometry- and Image-Based Approach. In *SIGGRAPH*, 1996. 3

[5] J Flynn, I Neulander, J Philbin, and N Snavely. Deep Stereo: Learning to Predict New Views from the World's Imagery. In *CVPR*, pages 5515–5524, 2016. 3

[6] Ysbrand Galama and Thomas Mensink. IterGANs: Iterative GANs to learn and control 3D object transformation. *Computer Vision and Image Understanding*, 2019. 3

[7] Ravi Garg, B G Vijay Kumar, Gustavo Carneiro, and Ian Reid. Unsupervised CNN for single view depth estimation: Geometry to the rescue. In *ECCV*, 2016. 3, 6

[8] Jan-Mark Geusebroek, Gertjan J. Burghouts, and Arnold W.M. Smeulders. The Amsterdam Library of Object Images. *International Journal of Computer Vision*, 61, 2005. 9, 10

[9] Clément Godard, Oisin Mac Aodha, and Gabriel J Brostow. Unsupervised Monocular Depth Estimation with Left-Right Consistency. In *CVPR*, 2017. 3, 6

[10] Clément Godard, Oisin Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into Self-Supervised Monocular Depth Prediction. In *ICCV*, 2019. 3, 6

[11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative Adversarial Nets. In *NeurIPS*. 2014. 3

[12] N Hirose, A Sadeghian, F Xia, R Martín-Martín, and S Savarese. VUNet: Dynamic Scene View Synthesis for Traversability Estimation Using an RGB Camera. *IEEE Robotics and Automation Letters*, 4(2):2062–2069, 2019. 1

[13] Derek Hoiem, Alexei A Efros, and Martial Hebert. Automatic Photo Pop-Up. In *SIGGRAPH*, 2005. 3

[14] Phillip Isola, Jun-Yan Yan Zhu, Tinghui Zhou, Alexei A Efros, and Berkeley Ai Research. Image-to-Image Translation with Conditional Adversarial Networks. In *CVPR*, volume 2017-Janua, pages 5967–5976, 2017. ISBN 9781538604571. 3

[15] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial Transformer Networks. In *NeurIPS*, 2015. 4

[16] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 6

[17] A Johnston and G Carneiro. Single View 3D Point Cloud Reconstruction using Novel View Synthesis and Self-Supervised Depth Estimation. In *2019 Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–8, 2019. 3, 4, 6

[18] Natasha Kholgade, Tomas Simon, Alexei Efros, and Yaser Sheikh. 3D Object Manipulation in a Single Photograph Using Stock 3D Models. *ACM Trans. Graph.*, 33(4): 127:1—-127:12, 2014. 1, 3

[19] B Kicanaoglu, R Tao, and A W M Smeulders. Estimating small differences in car-pose from orbits. In *British Machine Vision Conference*, 2018. 3

[20] M Meshry, D B Goldman, S Khamis, H Hoppe, R Pandey, N Snavely, and R Martin-Brualla. Neural Rerendering in the Wild. In *CVPR*, pages 6871–6880, 2019. 3

[21] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016. 5

[22] Thu Nguyen-Phuoc, Chuan Li, Stephen Balaban, and Yong-Liang Yang. RenderNet: A deep convolutional network for differentiable rendering from 3D shapes. In *NeurIPS*, 2018. 3

[23] Kyle Olszewski, Sergey Tulyakov, Oliver Woodford, Hao Li, and Linjie Luo. Transformable Bottleneck Networks. *ICCV*, 2019. 1, 3, 7, 8

[24] Eunbyung Park, Jimei Yang, Ersin Yumer, Duygu Ceylan, and Alexander C Berg. Transformation-grounded image generation network for novel 3D view synthesis. In *CVPR*, 2017. 1, 2, 3, 4, 5, 6, 7, 8

[25] Krishna Regmi and Ali Borji. Cross-View Image Synthesis Using Conditional GANs. In *CVPR*, 2018. 3

[26] K Rematas, C H Nguyen, T Ritschel, M Fritz, and T Tuytelaars. Novel Views of Objects from a Single Image. *IEEE Trans. on PAMI*, 39(8):1576–1590, 2017. 1, 3

[27] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. A Comparison and Evaluation of Multi-View Stereo Reconstruction Algorithms. In *CVPR*, 2006. 3

[28] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, Jan Kautz, and Jan Kautz Nvidia. PWC-Net: CNNs for Optical Flow Using Pyramid, Warping, and Cost Volume. In *CVPR*, 2018. 1, 2, 3, 4, 7, 8

[29] M Tatarchenko, A Dosovitskiy, and T Brox. Multi-view 3D Models from Single Images with a Convolutional Network. In *ECCV*, 2016. 1, 3

[30] Xiaogang Xu, Ying-Cong Chen, and Jiaya Jia. View Independent Generative Adversarial Network for Novel View Synthesis. In *ICCV*, 2019. 3

[31] Youyi Zheng, Xiang Chen, Ming-Ming Cheng, Kun Zhou, Shi-Min Hu, and Niloy J Mitra. Interactive Images: Cuboid Proxies for Smart Image Manipulation. *ACM Trans. Graph.*, 31(4), 2012. 1, 3

[32] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. View Synthesis by Appearance Flow. In *ECCV*, 2016. 1, 2, 3, 4, 7, 8

[33] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised Learning of Depth and Ego-Motion from Video. In *CVPR*, 2017. 3, 6

[34] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo Magnification: Learning View Synthesis using Multiplane Images. In *SIGGRAPH*, 2018. 3

[35] Zhou Wang, A C Bovik, H R Sheikh, and E P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. on Image Processing*, 13(4): 600–612, 2004. 7

[36] Xinge Zhu, Zhichao Yin, Jianping Shi, Hongsheng Li, and Dahua Lin. Generative Adversarial Frontal View to Bird View Synthesis. In *2018 International Conference on 3D Vision (3DV)*, pages 454–463. IEEE, 2018. 3

[37] C Lawrence Zitnick, Sing Bing Kang, Matthew Uyttendaele, Simon Winder, and Richard Szeliski. High-Quality Video View Interpolation Using a Layered Representation. In *SIGGRAPH*, 2004. 3

# 6 Supplementary materials

In the supplementary materials a more elaborate qualitative comparison is provided between the proposed method and other state-of-the-art methods. For this the synthesized target views for 36 car images are shown in Figure 7- 9 and the syntesized target views for 21 chair images are shown in Figure 10- 12.

From the results we observe that in line with the overall quality metrics (Table 2 of the main paper), our method synthesize in general higher quality (better geometrical shape and matching texture) compared to the other methods.

A few observations:

- Observe that TVSN, M2NV and the proposed method all have a similar inter-connection network architecture, in contrast to TBN. The inter-connection allow for explicit use of the input image pixels in constructing the generated views, and thus these models can retain the object textures in the generated views. Examples include row 4, 5 of Figure 8 and row 2, 5, 10 of Figure 9, where the specific color patterns or texts on the input views are retained in the generated views.

- Note also row 4 of Figure 8, which is indeed a failure, yet legitimate, case. The target is posed at an extreme angle with the input, but the unseen back of the truck has a different texture/color. The methods based on the assumption of object symmetry to make an "educated" guess of the unseen view, thus fail when the object texture does not follow such assumption.

- The main difference of our method to that of other state-of-the-arts is the explicit use of object geometry in reasoning of occlusion, symmetry and in generating new views. TVSN and M2NV use occlusion and symmetry in creating annotated data and in training the network to predict the coarse view, while the proposed method impose such assumption directly on the coarse view. This has the benefit when the target pose is close to the symmetric pose of the input. Examples can be seen in the generated chair images, specifically row 2-7 of Figure 10 and row 2-4 of Figure 11. Despite the intricate structure of chairs, these examples standout in quality compared to other methods.

Figure 7: Qualitative examples of cars generated by our point-cloud-based method and other related work. Our method retains better the objects' geometrical structure and detailed textures, *e.g.* row 3 and 5.

Figure 8: Qualitative examples of cars generated by our point-cloud-based method and other related work. Our method retains better the objects' geometrical structure and detailed textures, *e.g.* row 5 and 11.
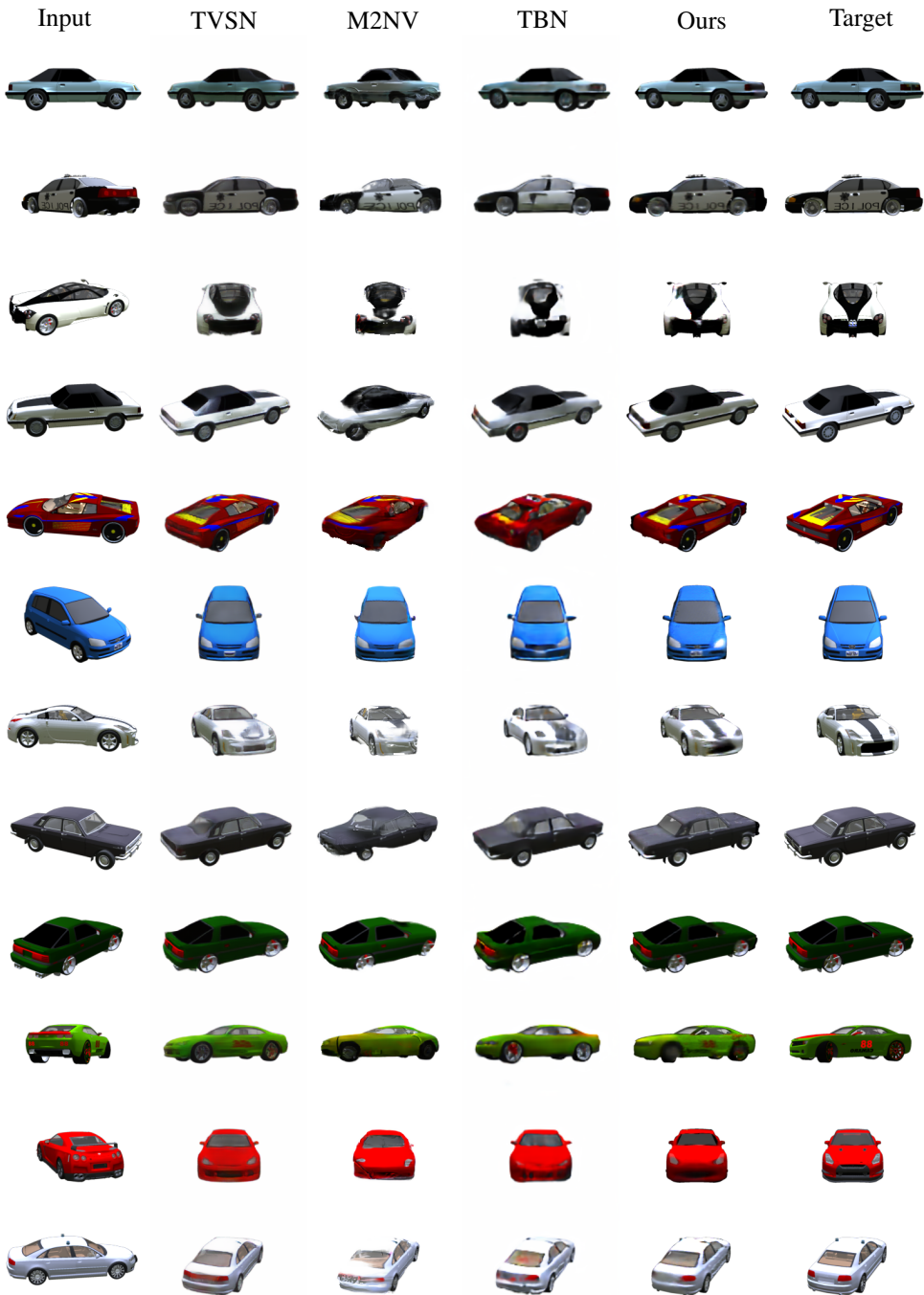
Figure 9: Qualitative examples of cars generated by our point-cloud-based method and other related work. Our method retains better the objects' geometrical structure and detailed textures, *e.g.* row 2 and 7.

Figure 10: Qualitative examples of chairs generated by our point-cloud-based method and other related work. Our method retains better the objects' geometrical structure and detailed textures, *e.g.* row 4 and 7.

Figure 11: Qualitative examples of chairs generated by our point-cloud-based method and other related work. Our method retains better the objects' geometrical structure and detailed textures, *e.g.* row 3 and 4.

Figure 12: Qualitative examples of chairs generated by our point-cloud-based method and other related work. Our method retain better the objects' geometrical structure and detailed textures, *e.g.* row 6 and 7.