# Peeking Behind Objects: Layered Depth Prediction from a Single Image

Helisa Dhamo[1], Keisuke Tateno[1,2], Iro Laina[1],
Nassir Navab[1], Federico Tombari[1]

[1] Technical University of Munich, Germany
{dhamo, tateno, laina, navab, tombari}@in.tum.de
[2] Canon Inc., Tokyo, Japan

**Abstract.** While conventional depth estimation can infer the geometry of a scene from a single RGB image, it fails to estimate scene regions that are occluded by foreground objects. This limits the use of depth prediction in augmented and virtual reality applications, that aim at scene exploration by synthesizing the scene from a different vantage point, or at diminished reality. To address this issue, we shift the focus from conventional depth map prediction to the regression of a specific data representation called Layered Depth Image (LDI), which contains information about the occluded regions in the reference frame and can fill in occlusion gaps in case of small view changes. We propose a novel approach based on Convolutional Neural Networks (CNNs) to jointly predict depth maps and foreground separation masks used to condition Generative Adversarial Networks (GANs) for hallucinating plausible color and depths in the initially occluded areas. We demonstrate the effectiveness of our approach for novel scene view synthesis from a single image.

## 1 Introduction

With the increasing availability of new hardware in the market of virtual reality and user-machine interaction, such as Head-Mounted Displays (HMD), various solutions have been explored to make traditional image-based content available for such devices. The advantage would be 3D immersive visualization of a huge amount of visual data which is either already available on web databases (Google, Flickr, Facebook) or commonly acquired by consumer cameras present on smartphones and tablets. One relevant and, recently, actively investigated direction is monocular depth prediction [1,2,3,4,5,6,7,8], that aims to regress the scene geometry from a single RGB frame, to provide 3D content that can be viewed via stereo HMDs.

Nevertheless, the scene geometry obtained via depth prediction is relative to the specific viewpoint associated to the RGB frame. When viewed through an HMD, viewpoint changes induced by the user's head motion would reveal holes in the predicted geometry due to the occlusion caused by foreground objects (Fig. 1, *right*). An alternative solution would be to predict a full 3D reconstruction of the scene from a single viewpoint. Nevertheless, this is a hard task to achieve
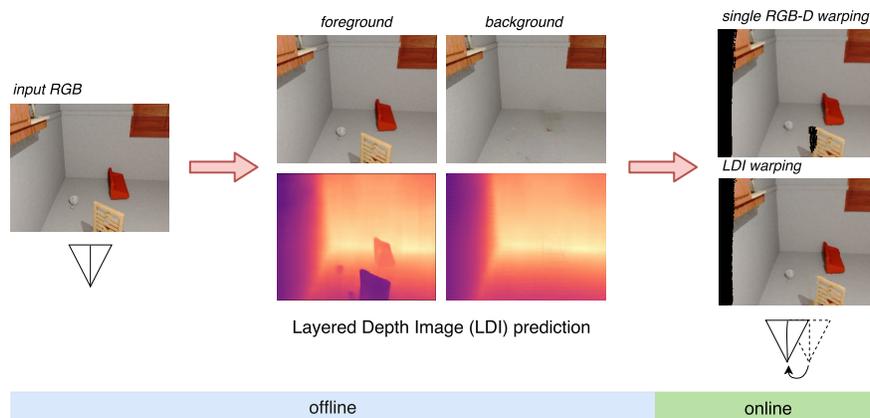
**Fig. 1.** (*Left*) Single RGB input in the original view. (*Center*) Layered Depth Image prediction using the proposed approach. The *foreground* layer consists of the original RGB image and monocular depth prediction; the *background* RGB-D layer describes the same scene after foreground object diminishing. (*Right*) One application of our method: viewpoint perturbation of the original image comparing two cases, simple warping of a single RGB-D and rendering from an LDI representation. The LDI generation is done offline, prior to the online simulation of the viewpoint changes.

due to the inherent ambiguity. Most related work is limited to reconstruction of objects only [9,10,11], or requires depth as that acquired by an active 3D sensor [12,13]. Also, the 3D volume of the scene is computation- and memory-intensive to process, which seems excessive for our goal, since for this type of applications we aim to solve the missing geometry only for small viewpoint changes.

An efficient compromise between 2.5D and 3D, would be an RGB-D structure that contains two layers, namely the first visible color and depth, and the next values along the same ray that are not part of the same object. This data representation was originally introduced by Shade et al. [14] as Layered Depth Image (LDI). While rendering from a target view, the RGB and depths from the second LDI layer are utilized to fill in the holes that become present after perturbing the viewpoint. Recently, Hedman et al. [15] improve upon the original LDI structure to construct panoramic 3D photography from multiple views. However, in practical situations the only input source is a color image, for which neither tracking nor depth data is available. Hence, we propose a novel approach to generate an LDI for a similar application, but restricting the input to a single RGB, as illustrated in Fig. 1. To deal with this added ambiguity, we make a *simple scene assumption*, meaning that there is no more than one level of occlusion in every pixel of any view. In particular, an advantage of our approach for applications such as HMD-based virtual reality is *efficiency*, as the LDI can be computed once (offline) and online warping is accomplished by directly processing the LDI data with low computational effort on the fly. This differs from the

field of view synthesis, where regions occluded from the available vantage point require inpainting online for every new viewpoint.

To the best of our knowledge, we are the first to infer an LDI representation from a single RGB image. As a first step of our pipeline (Fig. 2) we learn a standard depth map and a foreground–background mask using a fully–convolutional network. Second, using the latter prediction, we remove the foreground in the RGB and depth images. This results in incomplete RGB-D background images, thus we use a GAN-based approach to inpaint the missing regions. Our contribution is not defined by specific deep learning tools within our framework. In fact, our method is agnostic to individual inpainting and depth prediction architectures. However, we introduce a *pair discriminator* to encourage inter-domain consistency (here RGB-D). Besides view synthesis, the generated RGB-D with removed foreground, offers the potential for a range of AR applications, such as rendering new content in a diminished scene.

## 2   Related Work

To our knowledge, there is no research work on LDI prediction from a single RGB image. However, the work of this paper is built upon certain methods and data representations, which we want to visit in detail.

**Layered Depth Images (LDI)** encompass a scene representation that contains multiple depth values along each ray line in a single camera view. It was first introduced by Shade et al. [14], as an efficient image-based rendering method, emerged in times of limited computational power. When rendering in a perturbed view, some scene structures that are occluded in the original camera frame become visible – known as disocclusion – and therefore populate the target depth map with additional information. Applications include video view interpolation [16] and 3D photography [15]. Recently, Hedman et al. [15] present a 3D photo capturing from a set of hand–held camera images, that incorporates a cost term into a plane–sweep MVS, to penalize close depths. The different views are stitched together in a panorama LDI of two layers, which can be subject of geometry aware effects. Liu et al. [17] decompose a scene in depth layers (given an RGB-D input), with no color information being inferred in occluded regions.

In contrast to all previous approaches, our application assumes a single input image, which makes the problem of obtaining an LDI more challenging.

**Single-view depth estimation** One of the very first approaches for monocular depth estimation involved hand-engineered features and inference via Markov Random Fields [18,19]. Data-driven methods were later proposed that involved transferring and warping similar candidates from a database [20]. Recently, the application of CNNs for depth map prediction (2.5D) from a single RGB image is widespread. In the pioneering work of Eigen et al. [1], a multi-scale CNN is proposed for coarse and refined maps, further extended to normals and semantics [2]. Methods based on deeper networks further boost performance. Laina et al. [3] introduce a ResNet–based [21] fully convolutional network and a reverse Huber loss. Kendall and Yal [22] propose a Bayesian fully convolutional DenseNet

[23] for capturing uncertainty. In contrast, regression forests with shallow CNNs as tree nodes are proposed in [24]. Another family of methods [5,6,8] use Conditional Random Field regularizers to encourage geometrical consistency.

Besides 2.5D, CNNs have been also applied to 3D model prediction. Choy et al. [9] map an RGB image to a 3D volume, using a 2D encoder and 3D decoder architecture, as well as a recurrent update component in case of multiple inputs. Fan et al. [10] instead, predict a 3D point cloud from an RGB, by generating multiple 3D shapes. Wu et al. [11] (3D–VAE–GAN), builds a 3D model from the latent vector of an image. All these approaches output an object 3D model only, as opposed to our goal of representing whole scenes. Recently, Tulsiani et al. [25] introduce learning of a 3D scene representation that consists of a room layout and a set of object shapes. This method does not provide color information for the occluded background, as required in our proposed VR application.

Our LDI prediction lies on 2.75D, between these two major branches of CNN-based 3D perception.

**View synthesis** The task of predicting the depth or 3D structure of a scene is further related to novel view synthesis and we later demonstrate our LDI prediction within this challenging application. View synthesis has been extensively addressed in prior work. A recent line of work on synthesizing new views directly minimizes a reconstruction loss between the reference and the target image in an end-to-end manner [26,27,28]. Besides that, depth prediction often arises as an implicit, intermediate representation in such frameworks when using view pairs as input/supervision [28,29,30,31].

**CNN-based image inpainting** The goal of image inpainting (completion) is to fill in missing parts of an image. In the deep learning era, inpainting has been typically done with CNNs [32,33]. More recently, Generative Adversarial Networks (GANs) [34] have also been used for a variety of image prediction tasks, among which inpainting [35,36,37], editing [38] and 3D shape filling [39], due to their success in producing realistic looking samples. In essence, a GAN consists of a generator $G$ and a discriminator $D$, trained with conflicting objectives. G aims to generate realistic data, whereas D distinguishes samples from the real data distribution. Conditional GANs [40,41] use additional priors, such as a lower resolution image [42], or an incomplete image [35] to control the generation. Isola et al. [36] propose a general–purpose GAN for a variety of image translation tasks, including inpainting. In this work, we extend the image completion task to RGB-D data, which to our knowledge has not been tackled before.

## 3   Method

Our goal is to learn a mapping from a single RGB image to an LDI, i.e. a two–layer RGB-D representation, as depicted in Fig. 1. In Section 3.1 we describe the acquisition of the data which will be used for the purpose of this work, as well as the accompanying assumptions. The proposed pipeline is presented in Fig. 2. To break down the ambiguity of the prediction task, we suggest splitting the LDI prediction task into two subsequent steps. First, as explained in Section 3.2, we
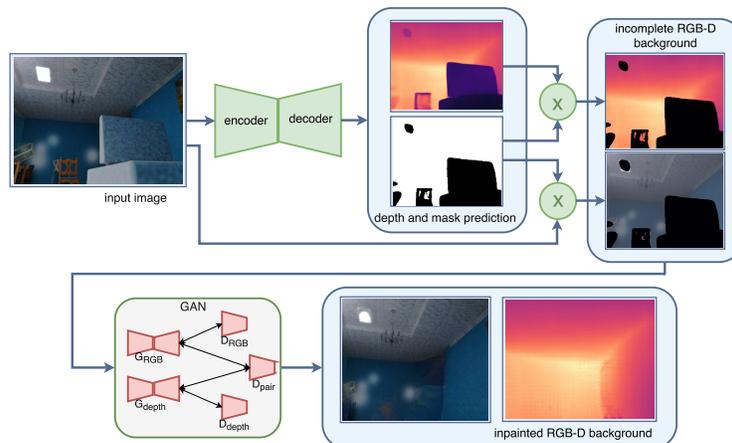
**Fig. 2.** General pipeline of our proposed method. (*Top*) A single layer depth map and segmentation mask is predicted, using a fully–convolutional CNN. The segmentation mask is applied to the RGB input as well as the predicted depth, so that the respective foreground pixels are discarded. (*Bottom*) The incomplete RGB-D information is inpainted using a GAN approach

train a network to jointly predict a conventional depth map and a foreground–background segmentation mask. Then, the mask is applied on both the RGB image and the predicted depth, so that only the visible background regions are fed to the next stage, while foreground regions are discarded. Finally, we perform GAN-based background completion conditioned on both the incomplete RGB as well as depth samples, as detailed in Section 3.3.

### 3.1   Dataset generation

Since learning to predict layered depth maps from a single image is a novel task, publicly available datasets containing LDI representations are quite limited for deep learning purposes. For instance, the authors of [15], currently provide a dataset of only 20 panoramic 3D photos [3], which have an LDI structure. Hence, we generate our own data based on existing datasets that provide RGB-D sequences with associated camera poses.

We construct the LDI samples by projecting multiple views of the same scene onto a reference frame. This populates the color and depth representations of the reference frame with new depth values, which correspond to background regions that were not visible in the original depth map. To obtain a robust warping from a source frame to a target frame, the camera pose at each frame has to be known in advance. In our experiments, we used around 30000 RGB-D frames with the associated camera poses. In addition, we build our dataset under the assumption of *no self–occlusion*, i.e. an object is not occluding parts of itself, to

---

[3] http://visual.cs.ucl.ac.uk/pubs/casual3d/datasets.html

avoid additional ambiguity in the prediction task. For this reason we also require object instance annotations, to make sure that at a specific coordinate, RGB-D values that correspond to the same object are not stored in both foreground and background layer. Given all these requirements, we generate our dataset based on SceneNet [43] data, that consists of synthesized photorealistic environments.

The LDI extraction procedure is as follows. For each subset of $N = 20$ consecutive RGB-D frames extracted from the same scene, we define the middle frame as reference, $f_{ref}$. Then, we warp every pixel $\boldsymbol{u}$ of the supportive frames $f_i$ into the reference view

$$\mathbf{u}_{\text{warped}} = \pi \left( \mathbf{T}_i^{\text{ref}} \pi^{-1} \left( \mathbf{u}_i \right) \right) \ , \tag{1}$$

where $\mathbf{T}_i^{\text{ref}}$ is the transformation from $f_i$ to $f_{ref}$, and $\pi()$ denotes perspective projection. We use $\mathbf{u}_{\text{warped}}$ as the reference frame coordinates of the re-projected colors, instance annotations and depths. The instance annotations are particularly useful during warping to keep track of the scene entities responsible for occlusions.

Along with pixel coordinates, the perspective projection gives the new depth values, corresponding to metric distance from the reference camera plane. Note that SceneNet [43] provide ray lengths instead of depths, therefore the following conversion is necessary

$$d(\mathbf{u}) = \frac{r(\mathbf{u})}{||\mathbf{K}^{-1}\dot{\mathbf{u}}||_2} \tag{2}$$

between the ray length $r$ and depth $d$, where $\mathbf{K}$ is the camera intrinsic matrix and $\dot{\mathbf{u}}$ denotes a pixel in homogeneous coordinates.

After computing the depth values for the occluded background regions of the reference view, we want to populate the background layer. The following validity conditions apply:

1. The candidate depth value is larger compared to the respective foreground pixel at that coordinate.

$$d_{\text{warped}}(\mathbf{u}_{\text{warped}}) > d_{\text{ref}}^{\text{FG}}(\mathbf{u}_{\text{warped}}) \tag{3}$$

2. The candidate background pixel and the respective original foreground pixel do not share the same object *id*. Here, we rely on the assumption of *no self–occlusion*.

$$id_{\text{warped}}(\mathbf{u}_{\text{warped}}) \neq id_{\text{ref}}^{\text{FG}}(\mathbf{u}_{\text{warped}}) \tag{4}$$

3. The candidate pixel is a potential background layer pixel, only if the object instance that contains it does not occlude other objects at any point. In such a way, we make sure that the *simple scene assumption* is fulfilled, meaning that an object instance can not be part of both layers.

Finally, we store the warped RGB-D frame with the smallest depth among all valid candidates in the background layer. Correspondingly, we extract a foreground–background segmentation mask, utilizing the accumulated information on object instances. Namely, a pixel in the image is considered as background
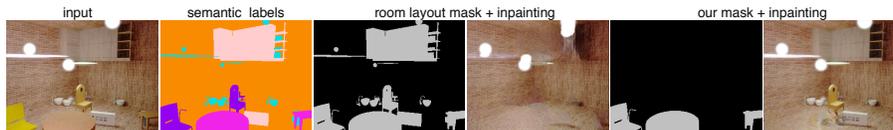
input | semantic labels | room layout mask + inpainting | our mask + inpainting

**Fig. 3.** Illustration of semantic room layout separation vs. our foreground segmentation method.

if its reference instance label is not found in the list of occluding instances, and foreground otherwise. As a result, as described in Section 3.3, the background inpainting will not be prone to undesired foreground context. We use the original *train/valid* split of SceneNet, to generate our training and test data. It is worth noting that the result of our method differs from a simple room layout separation, which would instead give a 3D "box" representation of the scene. In particular some of the object instances should still be considered as part of the background context, given an application that involves *small* viewpoint perturbations. For instance, in Fig. 3 we expect to see more parts of the occluded chair instead of plain floor, while exploring around the table in front of it.

Importantly, we observed that this LDI reconstruction approach does not work on real datasets like ScanNet [44], in which instance segmentation is not perfectly accurate. This leads to almost every pixel in the image being classified as foreground, as long as there is a single background pixel in the instance annotation maps, that is wrongly classified as foreground. Therefore, while for the current ground truth generation we rely on synthetic data, our method does also generalize to real-world data as shown later through experimental evaluation.

### 3.2 Joint depth map and segmentation mask prediction

The first stage of the proposed pipeline consists in standard depth map estimation as well as foreground–background segmentation from a single RGB image.

For this task, a wealth of CNN architectures for depth prediction and/or semantic segmentation exists in literature and could be employed [1,2,3,4,5,8]. In our work, we employ the fully convolutional ResNet-50 architecture proposed by Laina et al. [3] given its competitive performance in depth estimation and the publicly available implementation[4]. We modify the original network architecture by adding one more up–projection layer, so that the output depth preserves the input resolution. By converting the ray length images originally provided by SceneNet to metric distances, we get depth maps corresponding to the *visible* structures in the LDI representation of every frame, that is the *first layer* composing the LDI. We additionally aim to learn a segmentation map that would later allow to mask out the foreground entities in the current reference view, thus creating holes for the upcoming background completion. We recognize that this is not a deterministic task, especially in complex scenes, since there might

---

[4] https://github.com/iro-cp/FCRN-DepthPrediction

be multiple valid separation options. However, we train the same model for segmentation, considering the binary masks we generated in Section 3.1 as ground truth. This can be interpreted as *implicit* learning of an adaptive threshold for the separation, dependent on the distance from the camera, occlusion, as well as the structure continuity. For example, in Fig. 2, wall regions that are closer compared to the brown chair, are still predicted as background, since they are smoothly connected with the rest of the wall and do not cause occlusion.

We train our fully convolutional network on the depth prediction and foreground segmentation tasks, both separately and jointly. The latter shows superior performance for both tasks, which is intuitively also justified by the underlying relation between them. Our final model configuration allows the two tasks to share all network weights, except from the last layer. The loss function for the depth component is the reverse Huber $\mathcal{B}(d)$, following [3], whereas for the foreground segmentation $L_2$–norm $\mathcal{L}_2(s)$ outperformed the standard cross entropy loss. We train jointly by combining the depth loss and segmentation loss with equal weight:

$$\mathcal{L}_{\text{total}} = \mathcal{B}(d) + \mathcal{L}_2(s). \tag{5}$$

After obtaining the continuous segmentation predictions, we apply a threshold of 0.45 to distinguish between background and foreground. The threshold favors classification as foreground to prevent the subsequent background completion from getting confused by undesired foreground influence, while losing part of the background context does not have a negative effect on completion performance. Moreover, we observe that applying dilation on the resulting masks, further removes undesired foreground information, usually located around the object borders. Concretely, in our tests, we used cross–shaped dilation structure with a size of 5 pixels.

### 3.3  GAN-based RGB-D inpainting

Depth prediction *behind* occlusion brings additional ambiguity compared to conventional depth map prediction, since it implies inferring distance where the respective color context is not visible from a single view. Hence, learning depths for the unknown is rather a hallucinatory process, that involves creating plausible missing context conditioned on the visible regions as a prior.

In this work, we exploit the inpainting potential with a GAN-based approach. We start from a state-of-the-art model as a baseline and explore plausible ways of boosting the accuracy while extending to the RGB-D case. We use a similar architecture as in Isola et al. [36] for the generator. Our discriminators adopt the C64-C128-C256-C512 architecture as proposed in [36], where C denotes a Convolution–BatchNorm–ReLU block followed by the number of filters. The base loss in our inpainting GANs is

$$\mathcal{L}_{\text{inpaint}} = \mathbb{E}_{x,y}[logD(x,y)] + \mathbb{E}_x[log(1 - D(x,G(x)))] \tag{6}$$

where $x$ represents incomplete input and $y$ the corresponding full image. In this formulation, G optimizes the following objective

$$\hat{G} = \min_G \max_D \mathcal{L}_{\text{inpaint}} + \lambda_{\mathcal{L}_1} \mathcal{L}_1(y - G(x)). \qquad (7)$$

The $\mathcal{L}_1$ loss accounts for similarity between generated and ground truth samples.

Intuitively, one would argue that a joint learning between color and depth could potentially enhance the consistency between them. However, the semantic diversity between the two, encourages separate learning. We explore different architectures, to investigate the adverse responses of these two motivations.

First, we train a GAN for *combined RGB-D completion*. In this approach, the color and depth counterparts share all the generator and discriminator weights, leading to the lowest number of parameters. The objective function $\mathcal{L}_{\text{RGB}-\text{D}}$ has the form of Eq. 6, with $x = x_{RGB-D}$ and $y = y_{RGB-D}$. In addition, we explore *separate RGB and depth completion*, without any shared parameters between the respective RGB and depth inpainting, i.e. $\mathcal{L}_{\text{RGB}}$ and $\mathcal{L}_{\text{depth}}$ are updated independently from each other. Finally, we introduce an additional GAN model, built upon the latter and further enhanced with a multi-modal *pair discriminator* $D_{\text{pair}}$, whose role is to encourage inter-domain consistency between RGB and depth. We refer to it as *separate RGB and depth completion with pairing*. $D_{\text{pair}}$ takes an RGB-D input, either from the ground truth or the generator and distinguishes real RGB and depth *correspondences*. The respective generators $G_{\text{depth}}$ and $G_c$, receive feedback from their individual discriminators as well as from the combined domain, thus optimizing an additional term

$$\mathcal{L}_{\text{pair}} = \mathbb{E}_{x,y}[log D_{\text{pair}}(y_c, y_d)] + \mathbb{E}_x[log(1 - D_{\text{pair}}(G_c(x_c), G_d(x_d)))], \qquad (8)$$

where $x_c$ and $y_c$ refer to the color image, and $x_d$ and $y_d$ to depth. The final $\hat{G}_c$ and $\hat{G}_d$ objectives then become

$$\hat{G}_c = \min_{G_c} \max_{D_{RGB}} +\lambda_{\text{pair}} \min_{G_c} \max_{D_{pair}} \mathcal{L}_{\text{pair}} + \lambda_{\mathcal{L}_1} \mathcal{L}_1(y_c, G_c(x_c)) \qquad (9)$$

$$\hat{G}_d = \min_{G_d} \max_{D_d} +\lambda_{\text{pair}} \min_{G_d} \max_{D_{pair}} \mathcal{L}_{\text{pair}} + \lambda_{\mathcal{L}_1} \mathcal{L}_1(y_d, G_d(x_d)) \qquad (10)$$

with $\lambda_{\text{pair}} = 0.5$ and $\lambda_{L_1} = 100$. In all our GAN models at hand, after normalizing the input color and depth images separately to $[-1, 1]$, we set to $-2$ the values of inpainting interest. The GAN implicitly learns to identify and invalidate image regions assigned to those values.

## 4    Experiments

In this section, we present the experiments we conducted to evaluate the performance of our proposed method on SceneNet [43] (synthetic) and NYU depth v2 [45] (real) datasets. For SceneNet, we create our ground truth background data, as described in Section 3.1, therefore we perform both qualitative and quantitative measurements. Since this is not the case for NYU, we present qualitative results only, on the view synthesis application (Section 4.3), so that the reader can perceive the effect of the perturbed views in a real world context.
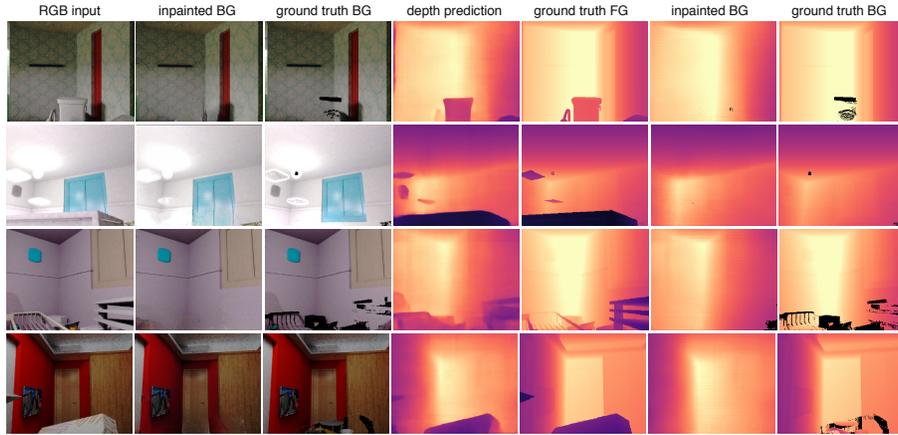
| RGB input | inpainted BG | ground truth BG | depth prediction | ground truth FG | inpainted BG | ground truth BG |
|---|---|---|---|---|---|---|



**Fig. 4.** Examples of foreground removal and background inpainting, for RGB and depth, on the SceneNet [43] dataset, with accompanying ground truth. Black indicates invalid pixels on the RGB and depth maps

### 4.1 Conventional depth and background mask

We train the joint depth map and foreground segmentation prediction task using a subset of around 30,000 RGB-D samples from the *train* partition of the SceneNet dataset. Further, we test on around 500 images, from the *valid* partition. For the depth prediction task we obtain a relative error of 0.184. The segmentation result is evaluated using the intersection over union metric (IoU), which is 0.71 for foreground and 0.93 for background pixels.

### 4.2 Ablation study on foreground diminishing

Table 1 shows the quantitative results on the predicted background RGB and depth maps of our proposed pipeline. In absence of previous similar work, we formulate this evaluation process as an ablation study, to consider our main investigated inpainting approaches. We compute the relative error (rel) and root mean square error (rms) for depth, and structural similarity index (SSIM) for RGB which gives an accuracy measure alongside the rms error. We measure the errors in two different cases, one is for the entire image while the other considers only the background information that was predicted by the inpainting stage. Please note that here we compare the quality of the depth completion against ground truth data, despite it being conditioned on the depth prediction results from the CNN. Thus, all results are products of our pipeline that operates on just a single image. In addition, Table 2 illustrates the same evaluations when directly inpainting incomplete ground truth depth maps, where ground truth segmentation masks are also used to separate the foreground.

The ablation study shows that the *combined RGB-D inpainting* approach performs the poorest in terms of depth prediction. In contrast, the accuracy of

| Method | Image area | RGB metrics | | Depth error | |
|---|---|---|---|---|---|
| | | ssim (higher) | rms (lower) | rel (lower) | rms (lower) |
| Combined RGB-D | *whole* | 0.935 | **17.59** | 0.156 | 0.574 |
| | *inpainted* | – | **54.46** | 0.197 | 0.704 |
| Separate RGB and depth | *whole* | 0.940 | 18.61 | 0.151 | 0.555 |
| | *inpainted* | – | 58.51 | 0.175 | 0.633 |
| Separate RGB and depth, paired | *whole* | **0.942** | 18.04 | **0.149** | **0.549** |
| | *inpainted* | – | 56.73 | **0.172** | **0.622** |

**Table 1.** Analysis of our GAN versions, on the SceneNet dataset [43]. Inpainting applied on *predicted FG-BG masks and depths*, from the CNN model from [3]

| Method | Image area | RGB metrics | | Depth error | |
|---|---|---|---|---|---|
| | | ssim (higher) | rms (lower) | rel (lower) | rms (lower) |
| Combined RGB-D | *whole* | 0.891 | **19.45** | 0.044 | 0.189 |
| | *inpainted* | – | **51.32** | 0.090 | 0.349 |
| Separate RGB and depth | *whole* | 0.900 | 20.09 | 0.018 | 0.100 |
| | *inpainted* | – | 53.98 | **0.041** | **0.193** |
| Separate RGB and depth, paired | *whole* | **0.903** | 19.76 | **0.017** | **0.095** |
| | *inpainted* | – | 52.70 | **0.041** | 0.197 |

**Table 2.** Analysis of our GAN versions, on the SceneNet dataset [43]. Inpainting module applied on *ground truth FG-BG masks and depths*

the RGB inpainting has a slight advantage compared to the other methods. One could argue, that this is due to the dominance of the RGB counterpart as textures and features from the color domain could have a negative impact on depth. Moreover, we believe that along with the incomplete RGB channels, the learning benefits from the incomplete depth information, as a better orientation for object separation. On the other end, the *separate RGB and depth* approaches, with and without pairing have a clear advantage in terms of depth inpainting, particularly on the newly added background region. Additionally, our pairing loss term brings an improvement in the depth inpainting task, particularly when learning from predicted depths. Since the *pair discriminator* encourages consistency between the complete RGB and depth maps, the depth map generator gets optimized to drive its output towards a plausible ground truth depth.

Comparing the results of Table 1 and 2, one can observe that the depth inpainting accuracy goes hand in hand with the accuracy of the predicted depth images. Notably, the relative error between the ground truth background and the completed map (Table 2) is in negligible range, whereas the relative errors in Table 1 are not far from the inherent depth prediction error, reported in Section 4.1. Conventional depth maps, usually have a good absolute scale, but bring high error values around the object borders, due to the lack of sharpness in edges.

This can explain why the relative errors from inpainting results, with removed foreground, are even lower than those of original depth map prediction. Interestingly, this means that LDI prediction from an RGB image, is an additional asset that can be easily incorporated into a regular depth prediction, without affecting the accuracy. In particular, recent advances in CNN depth prediction lead to more accurate LDIs, with no major change in the proposed method.

Fig. 4 illustrates background inpainting examples, using the CNN predicted depths and masks. As can be seen, the background ground truth is not always available, since during the warping procedure of Section 3.1, some regions have been not visible from any of the available viewpoints. Nevertheless, GAN inpainting covers the whole image. From the second example, we observe that although the inpainted image has hallucinated a door as opposed to the ground truth window, the result is equally plausible, considering the uncertainty.

### 4.3   View synthesis evaluation

To demonstrate the visual effect of a support background layer, we simulate a view perturbation scenario. In this experiment, the reference view corresponds to the original input RGB image. For the sake of simplicity, we define the reference camera pose to be in the origin of the world coordinates, meaning identity matrix rotation $\mathbf{R}_{ref}$ and zero translation $\mathbf{t}_{ref}$. We are interested in seeing the same content, from a slightly perturbed view angle. Concretely, we modify $t_x$ or $t_y$ in the reference translation vector $\mathbf{t}_{ref}$ to obtain a target pose with horizontal or vertical shift respectively. Next, we utilize Eq. 1 to warp the color and depth images of both LDI layers, from the source view to the target views. The foreground layer is warped first, followed by a dilation and erosion to fill small holes. Afterwards, the corresponding background layer is warped, substituting the pixel values that are still void after the morphological operations. Note that the respective background layer in this experiment is the inpainted output of the learned depths, which exposes the rendering task to additional inaccuracy.

For every test frame, we simulate different perturbation levels, in all four directions (top, bottom, left and right). The rendering results on the SceneNet dataset [43] are shown in Fig. 5, whereas the equivalent results on NYU [45] are to be seen in Fig. 6. One can observe the added value in completing the target view with additional pixels, that are occluded in the reference view. The completion task adapts to the background context, even when edges or relatively complex structures and patterns are present in the occluded background.

In almost every example, one can observe missing pixels around the image border, on the side towards which the viewpoint change was made. This is due to the fact that the current method does not perform inpainting outside the image borders. In some cases, background color is available outside the view, and it appears in the rendered image even though it is out of place. This is to be witnessed for instance in the last example of Fig. 6. Notably, dealing with inpainting outside the borders is not an easy task, since one has to hallucinate foreground information alongside with the background counterpart. Please refer to the paper supplement for more qualitative evaluations (images, video).
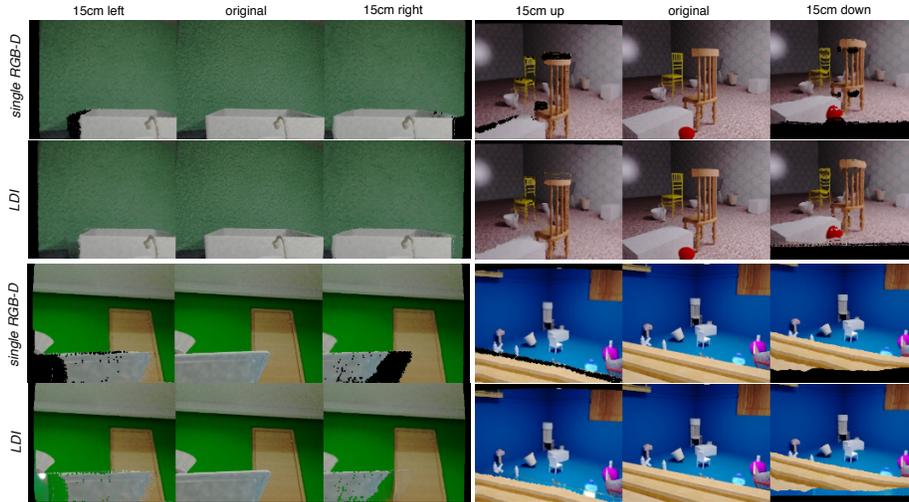
**Fig. 5.** Results on SceneNet [43], LDIs rendered on perturbed views. (*Upper row*) for each scene presents the warping of a simple RGB-D layer. (*Lower row*) shows the warping of the two layer RGB-D, obtained with our proposed pipeline.
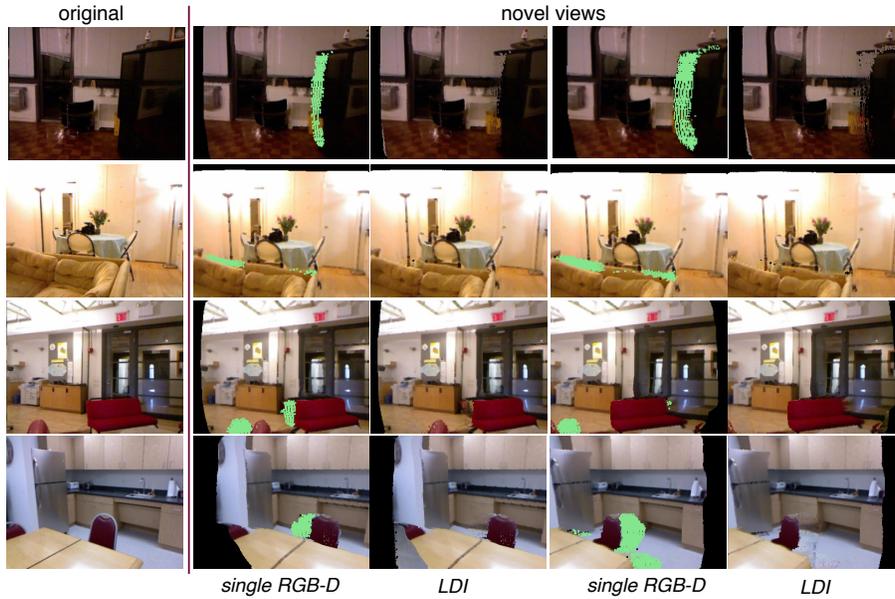


**Fig. 6.** Results on NYU dataset [45], LDIs rendered on perturbed views. (*Left*) Input RGB image. (*Right*) Two novel views obtained after rendering from a single RGB-D and our LDI representation.

We performed comparative evaluations to state-of-the-art view synthesis work. To the best of our knowledge, the most comparable method is Appearance Flow (AF) [27] as it considers a single input image and an arbitrary viewpoint transformation. Apart from that, we note that the majority of the view synthesis literature explores multi-view scenarios [26,28,29,30,31] or exploits scene geometry [29,30,31] (e.g. emerging depth via stereo) and is thus limited by the learned geometry (e.g. stereo baseline) for view generation. We trained their original model on SceneNet using the same partitions as for our own method. For fairness, we also trained AF based on a fully convolutional model (same as in our method [3])[5]. We report the comparisons in Table 3 and show that the application of LDI to view synthesis outperforms AF [27] in two different metrics, Mean Absolute Error (MAE) and SSIM. Moreover, from Fig. 7 one can see that our method preserves object shapes and aligns the image better with the target view.



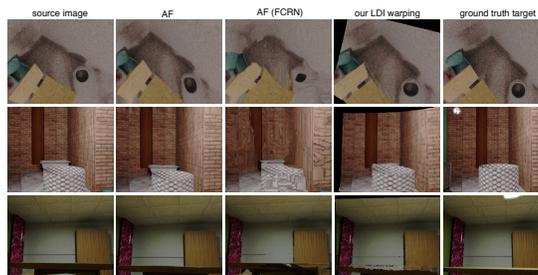| Method | MAE | SSIM |
|---|---|---|
| AF [27] | 0.200 | 0.537 |
| AF [27] (FCRN) | 0.185 | 0.534 |
| Ours | **0.147** | **0.617** |

**Table 3.** Quantitative comparison with view synthesis methods.

**Fig. 7.** Qualitative comparison of novel view synthesis.

## 5   Conclusion

We proposed a method to regress a layered depth map from a single RGB image. We illustrated how such maps can be regressed via a pipeline built over a CNN and a GAN. In addition we also demonstrate how the additional information included in a layered depth map can be useful for an enhanced user experience of the 3D content of a scene with respect to depth prediction, e.g. by improving view synthesis under occlusion. Importantly, the quality of the LDI prediction goes along with the accuracy of the individual components of our pipeline, such as CNN-based depth prediction and GAN-based inpainting. Future work aims at learning to regress depth representations that support more than one level of occlusion, thus capable to overcome the *simple scene assumption*.

---

[5] In their public repository, Zhou et al. [27] report better performance when switching to fully convolutional networks.

# References

1. Eigen, D., Puhrsch, C., Fergus, R.: Depth map prediction from a single image using a multi-scale deep network. In: NIPS. (2014)
2. Eigen, D., Fergus, R.: Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In: ICCV. (2015)
3. Laina, I., Rupprecht, C., Belagiannis, V., Tombari, F., Navab, N.: Deeper depth prediction with fully convolutional residual networks. In: 3DV. (2016)
4. Liu, F., Shen, C., Lin, G.: Deep convolutional neural fields for depth estimation from a single image. (2015)
5. Li, B., Shen, C., Dai, Y., van den Hengel, A., He, M.: Depth and surface normal estimation from monocular images using regression on deep features and hierarchical crfs. In: CVPR. (2015)
6. Liu, M., Salzmann, M., He, X.: Discrete-continuous depth estimation from a single image. In: CVPR. (2014) 716–723
7. Wang, P., Shen, X., Lin, Z., Cohen, S., Price, B., Yuille, A.L.: Towards unified depth and semantic prediction from a single image. In: CVPR. (2015)
8. Xu, D., Ricci, E., Ouyang, W., Wang, X., Sebe, N.: Multi-scale continuous crfs as sequential deep networks for monocular depth estimation. In: CVPR. (2017)
9. 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)
10. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: CVPR. (2017)
11. Wu, J., Zhang, C., Xue, T., Freeman, W.T., Tenenbaum, J.B.: Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling. In: NIPS. (2016) 82–90
12. Firman, M., Mac Aodha, O., Julier, S., Brostow, G.J.: Structured Prediction of Unobserved Voxels From a Single Depth Image. In: CVPR. (2016)
13. Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., Funkhouser, T.: Semantic scene completion from a single depth image. CVPR (2017)
14. Shade, J., Gortler, S., He, L.w., Szeliski, R.: Layered depth images. SIGGRAPH '98 (1998) 231–242
15. Hedman, P., Alsisan, S., Szeliski, R., Kopf, J.: Casual 3D Photography. (2017)
16. Zitnick, C.L., Kang, S.B., Uyttendaele, M., Winder, S., Szeliski, R.: High-quality video view interpolation using a layered representation. In: ACM SIGGRAPH. (2004)
17. Liu, C., Kohli, P., Fukurawa, Y.: Layered scene decomposition via the occlusion-crf. In: CVPR. (2016)
18. Saxena, A., Chung, S.H., Ng, A.Y.: Learning depth from single monocular images. In: Advances in neural information processing systems. (2006) 1161–1168
19. Saxena, A., Sun, M., Ng, A.Y.: Make3d: Learning 3d scene structure from a single still image. PAMI (2009)
20. Karsch, K., Liu, C., Kang, S.B.: Depth transfer: Depth extraction from videos using nonparametric sampling. In: Dense Image Correspondences for Computer Vision. Springer (2016) 173–205
21. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. CVPR (2016)
22. Kendall, A., Gal, Y.: What uncertainties do we need in bayesian deep learning for computer vision? In: NIPS. (2017) 5580–5590

23. Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: CVPR. (2017)
24. Roy, A., Todorovic, S.: (Monocular depth estimation using neural regression forest)
25. Tulsiani, S., Gupta, S., Fouhey, D., Efros, A.A., Malik, J.: Factoring shape, pose, and layout from the 2d image of a 3d scene. In: CVPR. (2018)
26. Flynn, J., Neulander, I., Philbin, J., Snavely, N.: Deepstereo: Learning to predict new views from the world's imagery. In: CVPR. (2016)
27. Zhou, T., Tulsiani, S., Sun, W., Malik, J., Efros, A.A.: View synthesis by appearance flow. In: ECCV, Springer (2016) 286–301
28. Xie, J., Girshick, R., Farhadi, A.: Deep3d: Fully automatic 2d-to-3d video conversion with deep convolutional neural networks. In: ECCV. (2016)
29. Garg, R., BG, V.K., Carneiro, G., Reid, I.: Unsupervised cnn for single view depth estimation: Geometry to the rescue. In: ECCV. (2016)
30. Godard, C., Mac Aodha, O., Brostow, G.J.: Unsupervised monocular depth estimation with left-right consistency. In: CVPR. Volume 2. (2017)  7
31. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: CVPR. (2017)
32. Ren, J.S., Xu, L., Yan, Q., Sun, W.: Shepard convolutional neural networks. In Cortes, C., Lawrence, N.D., Lee, D.D., Sugiyama, M., Garnett, R., eds.: NIPS. Curran Associates, Inc. (2015) 901–909
33. Yang, C., Lu, X., Lin, Z., Shechtman, E., Wang, O., Li, H.: High-resolution image inpainting using multi-scale neural patch synthesis. In: CVPR. (2017)
34. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: NIPS. (2014)
35. Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A.: Context encoders: Feature learning by inpainting. In: CVPR. (2016)
36. Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In CVPR (2017)
37. Iizuka, S., Simo-Serra, E., Ishikawa, H.: Globally and locally consistent image completion. ACM Transactions on Graphics (TOG) (2017)
38. Zhu, J.Y., Krähenbühl, P., Shechtman, E., Efros, A.A.: Generative visual manipulation on the natural image manifold. In: ECCV. (2016)
39. Wang, W., Huang, Q., You, S., Yang, C., Neumann, U.: Shape inpainting using 3d generative adversarial network and recurrent convolutional networks. arXiv preprint arXiv:1711.06375 (2017)
40. Mirza, M., Osindero, S.: Conditional generative adversarial nets. arXiv preprint arXiv:1411.1784 (2014)
41. Radford, A., Metz, L., Chintala, S.: Unsupervised representation learning with deep convolutional generative adversarial networks. arXiv preprint arXiv:1511.06434 (2015)
42. Ledig, C., Theis, L., Huszar, F., Caballero, J., Aitken, A.P., Tejani, A., Totz, J., Wang, Z., Shi, W.: Photo-realistic single image super-resolution using a generative adversarial network. CoRR **abs/1609.04802** (2016)
43. McCormac, J., Handa, A., Leutenegger, S., J.Davison, A.: Scenenet rgb-d: 5m photorealistic images of synthetic indoor trajectories with ground truth. ICCV (2017)
44. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: CVPR. (2017)
45. Nathan Silberman, Derek Hoiem, P.K., Fergus, R.: Indoor segmentation and support inference from rgbd images. In: ECCV. (2012)