

Points2Pix: 3D Point-Cloud to Image Translation using conditional GANs

Stefan Milz^{1 *}, Martin Simon^{1 *}, Kai Fischer^{1 *}, Maximillian Pöpperl¹, and
Horst-Michael Gross²

¹ Valeo Schalter und Sensoren GmbH, Germany

² Ilmenau University of Technology, Germany

Abstract. We present the first approach for 3D point-cloud to image translation based on conditional Generative Adversarial Networks (cGAN). The model handles multi-modal information sources from different domains, i.e. raw point-sets and images. The generator is capable of processing three conditions, whereas the point-cloud is encoded as raw point-set and camera projection. An image background patch is used as constraint to bias environmental texturing. A global approximation function within the generator is directly applied on the point-cloud (PointNet). Hence, the representative learning model incorporates global 3D characteristics directly at the latent feature space. Conditions are used to bias the background and the viewpoint of the generated image. This opens up new ways in augmenting or texturing 3D data to aim the generation of fully individual images. We successfully evaluated our method on the KITTI and SunRGBD dataset with an outstanding object detection inception score.

1 Introduction

Domain translation is a well known and widely applied problem. It is typically treated in computer graphics or computer vision. Most research focuses on image-to-image translation [7,34,27]. Examples are Semantic-Labels to Image (e.g. Labels to Street-Scene, Labels to Facades) or Image conversions (e.g. Day to Night, Black-and-White to Color). Those techniques deal with real domain translation problems, since they convert semantic sensor-independent context into realistic RGB image data or vice versa. However, domain translation is performed on top of images. Both domains encode the information as RGB values in pictures with a spatial dependency. We call that single mode domain translation:

$$\begin{aligned} G_{x \rightarrow y} : x &\rightarrow y & x, y &\in \mathbb{R}^{w \times h \times 3} \\ G_{y \rightarrow x} : y &\rightarrow x \end{aligned} \tag{1}$$

Whereas, $G_{x \rightarrow y}, G_{y \rightarrow x}$ describe the translation functions between both image domains x, y with fixed image sizes: h (height), w (width).

* Equal contribution

We propose a novel multi-modal domain translation model using the example of 3D point-cloud to image translation. The treated problem is formally known as:

$$G_{p \rightarrow y} : p \rightarrow y \quad p \in \mathbb{R}^{n \times 3} \quad y \in \mathbb{R}^{w \times h \times 3} \quad (2)$$

Here, n describes the number of points within the point-set. Our work is limited to $G_{p \rightarrow y}$ (not $G_{y \rightarrow p}$). Therefore an extensive new architecture is presented as combination of a typical encoder-decoder for image segmentation (UNet: [21]) as proposed by [7]. More important is the model's second input, where the architecture incorporates the real point-set to add 3D characteristics into the global feature space for constraint based individual image generation. We put conditions as viewpoint dependent projection and background image patches for fully individual image generation in compliance with 3D specifications (conditions: background, shape, distance, viewpoint).

2 Related Work

2.1 Image generation

Handcrafted Losses As image generation could be reduced to per-pixel classification/regression with a wide application area it turns out to have a long tradition [23,30,31,6]. Those applications suppose a conditionally unstructured loss applied on the output space, i.e. a pixel independence in terms of semantic relationship is supposed. The performance of those approaches strongly depends on the loss design, e.g. semantic segmentation [15].

Conditional GANs Conditional GANs (cGAN) instead learn structured losses that affect the overall output in form of a joint improvement [7]. In common, the cGAN is applied in a conditional setting. For image generation researchers were setting variable conditions: e.g. discrete labels [14,4], text [20] and images [7,34,27].

In general the cGAN performs a mapping function G , called generator, based on a condition c and a random noise vector z to generate an image y :

$$G : \{c, z\} \rightarrow y \quad y \in \mathbb{R}^{w \times h \times 3} \quad c = \begin{cases} \in \mathbb{R} & \rightarrow \text{label-to-image} \\ \in \mathbb{R}^t & \rightarrow \text{text-to-image} \\ \in \mathbb{R}^{w \times h \times 3} & \rightarrow \text{image-to-image} \end{cases} \quad (3)$$

For image-to-image translation [7] proposes a *U-Net* like structure for G . To create realistic images at higher resolutions (e.g. 1024 x 2048) [27] recommend a pyramidal approach for G similar to a *PSPNet* [32].

In general, the cGAN is composed by G and a competing discriminator D , which distinguishes between real images and created fake ones. A well-established discriminator network is the *Patchgan* [10] proposed by [7]. Derived from that the competing objective of the cGAN could be described by its loss L_{cGAN} :

$$L_{\text{cGAN}}(G, D) = \mathbb{E}_{c,y} \{\log(D(c, y))\} + \mathbb{E}_{c,z} \{\log(1 - D(c, G(c, z)))\} \quad (4)$$

2.2 Point-cloud processing

High requirements for perception tasks of robotic applications enforced the usage of 3D sensors, e.g. RGBD-cameras [26], Lidar (Valeo SCALA). Research progress in the field of 3D point-cloud processing received a boost in the recent years. In principle, point-clouds have specific properties that clearly distinguish them from images. Hence, specific processing models are needed. Points usually are not ordered, there is no grid that encodes the 3D position as an image does. The overall category of a point-set is influenced by the interaction of points among others. Only the global sum of the points forms a shape with a meaning. Last, point-sets are invariant to basic transformations like translation or rotation. Therefore the combination of 3D points clouds and machine learning is indispensable. The processing type could be categorized into the following three classes.

Real 3D Point-cloud processing [16] proposed the first neural network architecture *Point-Net* that handles natural points sets for classification and segmentation tasks with outperforming segmentation results on ShapeNet [2]: **mIoU 83.7**. The model does not use convolutional layers, but fully connected ones and directly processes the coordinates of the point-set (size n): $p = x_1, \dots, x_n$ with $p \in \mathbb{R}^{n \times 3}$. A chain of local transformations h on the point-set followed by a global max-pooling layer is used to create an overall feature space, i.e. a global approximation function:

$$\begin{aligned} f(x_1, \dots, x_n) &\approx g(h(x_1), \dots, h(x_n)) & f : \mathbb{R}^{n \times 3} &\rightarrow \mathbb{R} \\ & & h : \mathbb{R}^{n \times 3} &\rightarrow \mathbb{R}^K \\ & & g : \{\mathbb{R}_1^K \times \dots \times \mathbb{R}_n^K\} &\rightarrow \mathbb{R} \end{aligned} \quad (5)$$

I.e. The overall meaning f (e.g. object class) of a point set p is approximated by g . The advantage of the architecture is that it is robust against unordered point-clouds and transformations. The independence from the viewpoint variance helps to train with less training samples. Due to disadvantage for learning global features of large point sets the authors developed a second version *Point-Net++* [18].

Voxelization Voxelization approaches make use of the findings performing CNNs on images. Therefore, 3D data is converted to voxels or grid cells. After pre-processing standard machine learning architectures are applied. Unordered point-sets are avoided. Famous applications are 3D object detectors like [24,3,13,9]

Combined models Combined models have often shown most robust results (e.g. 3D object detection) and mostly make use of different sensor types. [33] investigated a method based on many local *Point-Nets* followed by a global 3D

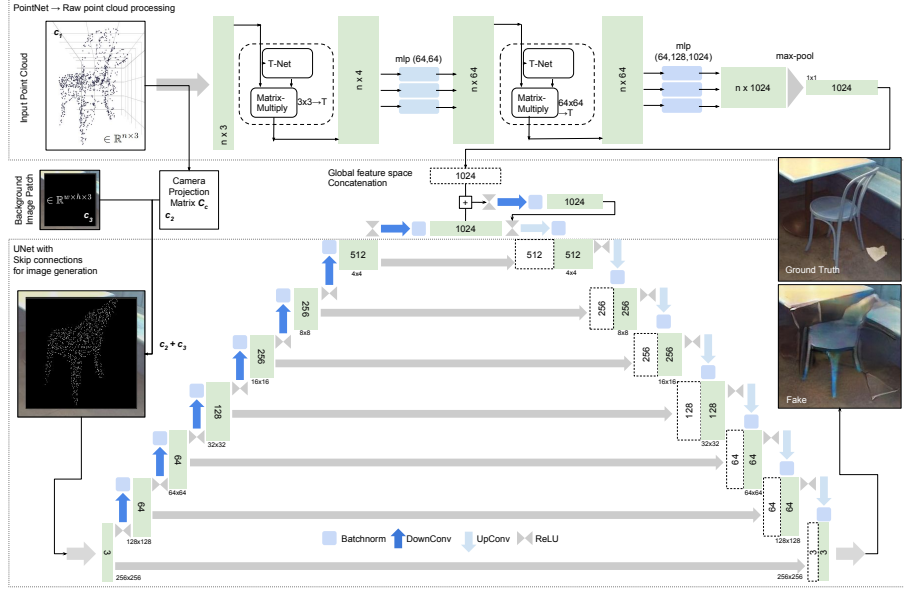


Fig. 1. Points2Pix Generator Architecture. The figure outlines the overall pipeline of the Points2Pix generator. In general we split the design into three areas: Top: the PointNet for a raw point-cloud processing; Bottom: Unet with skip connections for image generation; Middle: Global feature space concatenation from point-set (top) and image processing (bottom) pipeline. The model needs only a raw point-set as input, which acts as condition c_1 . The point-sets coordinates will be directly processed by *PointNet*. A projection into the image plane is used for UNet as input, whereas as the camera projection matrix C_c works as condition c_2 . Additionally an arbitrary background patch c_3 is used for background generation.

CNN. [17] architectures works the other way around. With the aid of a camera frustum points are filtered using a camera object detector. The filtered points are processed for 3D object detection with only one *Point-Net* [16] up to the last global max pooling layer ending in a 1×1024 general feature space. A 8 bit depth projection using the given camera projection matrix (c_2)

Generative models Point Cloud GAN by [11] is a famous approach for point-cloud generation. They do not perform any translation task, but they show that the common discriminator is not suitable for raw point-clouds. [1] performs label to point-cloud translation by using representative learning and introduce several 3D GAN derivatives. A similar study is published by [29] with focus on latent space analysis. However, learning 3D representations to generate viewpoint based images is missing within the research community. Therefore, we propose our novel technique Points2Pix.

3 Points2Pix

We propose a novel cGAN architecture for generating photo-realistic images from pure point-clouds. Additionally, we describe conditions to bias the view-point, distance, shape and the background within the latent space. Therefore, we introduce the network architecture consisting of a generator G (converting points to images), a discriminator D and the specific loss.

3.1 Generator

The objective of our generator $G(c_1, c_2, c_3)$ is to translate point-clouds into realistic-looking images, while using three conditions c_1, c_2, c_3 . The whole architecture is shown in Figure 1. The design is inspired by [7], which serves as base.

Condition one First, c_1 as raw point-cloud $c_1 = \{x_1, \dots, x_n\}$ is processed by *PointNet* [16]. The model samples $n = 1024$ points as input, applying an input transformation and aggregates global point features using fully connected layers and a generic max pooling (see equation 5):

$$f(c_1) \approx g(h(x_1), \dots, h(x_n)) \quad n = \{1 \dots 1024\} \quad (6)$$

However, in contrast to the basic *PointNet* pipeline, the proposed model incorporates the the global 3D feature space ($n \times 1024$) using concatenation at the innermost part within the Image encoder-decoder (*UNet*). Hence, $h(x_1), \dots, h(x_n)$ are applied by the *PointNet* part, g is implicitly performed with the aid of *UNets* decoder (see Fig. 1).

Condition two The second condition denoted as $c_2 = \mathbf{C}_c$ is an image projection of the point-cloud using a perspective projection matrix \mathbf{P}

$$\mathbf{P} = \begin{bmatrix} s & 0 & 0 & 0 \\ 0 & s & 0 & 0 \\ 0 & 0 & -\frac{f_c}{f_c - n_c} & -1 \\ 0 & 0 & -\frac{f_c \cdot n_c}{f_c - n_c} & 0 \end{bmatrix} \quad s = \frac{1}{\tan(\frac{fov}{2} \cdot \frac{\pi}{180})} \quad \mathbf{x}_{\text{pixel}} = \underbrace{\mathbf{P} \cdot \mathbf{T}_{\text{pc}}^{\text{cam}}}_{\mathbf{C}_c} x_i \quad (7)$$

$$\mathbf{y}_g(\mathbf{x}_{\text{pixel}}) = \frac{|\mathbf{x}_{\text{cam}} - x_i|}{d_{\text{max}}} \quad \mathbf{y}_b(\mathbf{x}_{\text{pixel}}) = I(x_i) \quad (8)$$

with a scaling s according to the horizontal field of view denoted as fov in degrees, near clipping plane denoted as n_c and far clipping plane f_c . We encode radial depth ($\mathbf{y}_g(\mathbf{x}_{\text{pixel}}) \rightarrow$ green channel) with a normalized depth d_{max} and intensities $I(x_i)$ of the measured reflectance for each point falling into the projection image ($\mathbf{y}_b(\mathbf{x}_{\text{pixel}}) \rightarrow$ blue channel). Before applying P , all points are transformed into the camera coordinate system using the extrinsic calibration $\mathbf{T}_{\text{pc}}^{\text{cam}}$. In this way we ensure the consistent viewpoint during training compared to the raw ground truth rgb image.

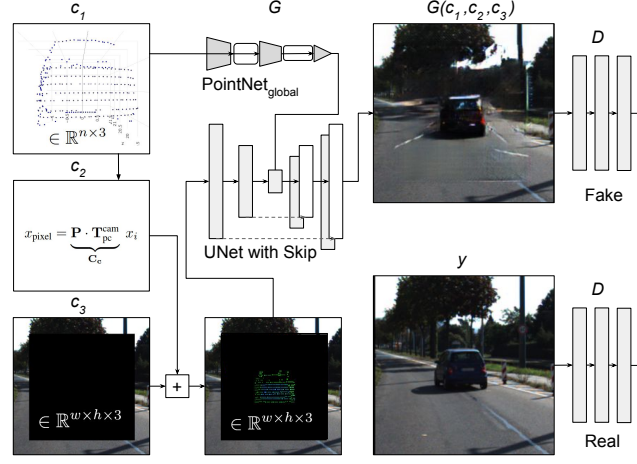


Fig. 2. Training Points2Pix: The figure outlines the competing training structure. The generators (G) output is a fake image based on its three conditions c_1, c_2, c_3 . The discriminator D has to distinguish between fake $D[G(c_1, c_2, c_3)]$ and real images $D[y]$.

Condition three Finally, the third condition c_3 is an arbitrary image background patch constraining environmental texturing. A surrounding image patch of the object cropped from the data set centered at the object origin up to a size of 256×256 is extracted. During training the image background patch is compliant to the ground truth. In test-mode, background patches can be randomly mixed with point-clouds.

Both, c_2 and c_3 are combined to an 256×256 input image, which is fed into a UNet with skip connections. At the innermost part, down sampled input features are concatenated with the global 3D feature space from c_1 . After up-sampling the output is a generated image with 256×256 pixels. Since, we use a cGAN for training, there is no need for an unstructured Loss. The assessment of the output is performed by the discriminator. As a note, we do not use a random noise vector z (3). Noise is only incorporated as dropout similar to [7].

3.2 Discriminator

We use the Markovian discriminator *PatchGAN* [7] that tries to distinguish between fake $D[G(c_1, c_2, c_3)]$ and real images $D[y]$ at the scale of $N \times N$ patches as well as possible. In contrast to [7] we do not take the condition c_1, c_2, c_3 into account. The output depends only on the generated image. Therefore, it consists of an $L1$ term to force low-frequency correctness [34] and is applied convolutionally across the image, averaging all responses. We only use 5 convolutional layers with batch and instance normalization. In this way, it effectively solves the problem to be able to model high- and low frequency structures at once.

3.3 Loss

The objective of a basic GAN can be explained by an additive combination of the generative network L_G loss and the discriminative network L_D loss. In order to iteratively improve results during training L_G should be reduced while L_D grows ideally. Consequently the basic cGAN loss can be described as follows assuming the three input conditions (c_1, c_2, c_3) :

$$L_{cGAN}(G, D) = \mathbb{E}_y[\log(D(c_2, c_3, y))] + \mathbb{E}_{c_1, c_2, c_3}[\log(1 - D(c_2, c_3, G(c_1, c_2, c_3)))] \quad (9)$$

Random noise z (3) is only realized using dropout. Compared to the typical cGAN loss L_{cGAN} (4), the model does not involve all conditions into the discriminator. However, we implicitly force conditions to be compliant within the output by using a weighted L_1 term [27] in the overall loss. This part describes the L_1 difference between the output and the ground truth. The final loss can be written as:

$$L_{\text{Points2Pix}} = L_{cGAN}(G, D) + \lambda_{L1} \cdot \mathbb{E}_{c_1, c_2, c_3, y}[\|y - G(c_1, c_2, c_3)\|_1] \quad (10)$$

4 Experiments

We conduct experiments on KITTI [5] for outdoor and SunRGBD [25] for indoor scenarios to explore the general validity of the method. Additionally, we show that the approach works for both, Lidar generated point-clouds and point-clouds coming from by RGB-D sensors. Following the recommendations of [7], the quality of the synthesized images is evaluated using an object based inception score. Furthermore, classification and diversity scores are added as additional assessment. Finally, we present some insights into our architecture decisions with additional ablation experiments.

4.1 Metrics

To assess the realism of the produced images, YOLOv3 [19] is used for validation. It is an off-the-shelf state of the art 2D object detector pre-trained on ImageNet and fine-tuned on the MS-Coco [12] data-set. This model includes overlapping classes in comparison to our experiments, e.g. cars (for KITTI) and chair (for SunRGBD). For the quantitative metrics we follow the instructions recommended by [28].

Classification Score With the aid of YOLOv3 the number of correct detected classes is measured. This could be achieved due to object centered image patches in our experiments. The classification score S_c ratio is then given by the detection ratio of fake images and ground truth ($TP \rightarrow$ true positives). The score could be directly affected by adjusting the confidence rate of the 2D object detector: $S_c = TP_{\text{fake}}/TP_{\text{real}}$.

Object based Inception Score¹ For positive results in terms of classification we measure the intersection over union $IoU_{\text{Points2Pix}}$ of the predicted bounding box BB coming out of YOLOv3 for the ground truth and the accompanied fake image.

$$IoU_{\text{Points2Pix}} = \frac{BB_{\text{fake}} \cap BB_{\text{real}}}{BB_{\text{fake}} \cup BB_{\text{real}}} \quad | \quad S_c = 1 \quad (11)$$

Diversity Score We measure the diversity ability of our cGAN to produce a wide spread of different output features using a diversity score. Our objective is to bias the shape, distance and 3D characteristics of the object. We collect randomly ten different background image patches, while keeping the point-cloud constant ($c_1 = \text{const}$ and $c_2 = \text{const}$, $c_3 \rightarrow \{1 \dots 10\}$). This leads to different output images that all should have the same 3D object inside. Therefore we compare the ground truth YOLOv3 results and all the fake images with the aid of calculating the mean S_c and the mean IoU .

4.2 Training Details

We train the network on both data sets separately for 100 epochs from scratch each, using the ADAM optimizer [8], with a learning rate of 0.0002 and momentum parameters $\beta_1 = 0.5$, $\beta_2 = 0.999$ such as $\lambda_{L1} = 100$. For our background condition c_3 , we use image patches with a border width of 15 pixels. We found using objects containing at least 700 points in their point-cloud as a good trade-off for minimum point density as well as object size.

Kitti: In a pre-processing step, we split the 7481 training examples of the 3D object detection benchmark and use 3712 samples for training and 3769 for evaluation. Therefore, we generate more than 20k training images for the class *car* only using $d_{\text{max}} = 60m$. Thus, each camera image is cropped centered at one labeled object with 256×256 pixels. At the same time, strongly occluded or truncated objects are skipped.

SunRGBD: We extract 3267 images from the SunRGBD data-set containing the following classes: *chair*, *table*, *desk*, *pillow*, *sofa* and *garbage bin*. The split for training and validation is a 90/10 ratio. Image patches are extracted at the object center from the cameras point of view with a size of 256×256 pixels with $d_{\text{max}} = 4m$. The depth information comes from either MS kinect v1 or v2 and the Intel real-sense. Since, those sensors do not measure a reflectance, we only encode the radial depth inside the projection of c_2 . Hence, the projection image contains one channel only.

4.3 Results

In Fig. 3 we show qualitative results for both data-sets and four different classes. Widely distributed output images are produced by alternating the background

¹ We call it inception score, because its similar to the proposal of [22]. We do not use an inception model.

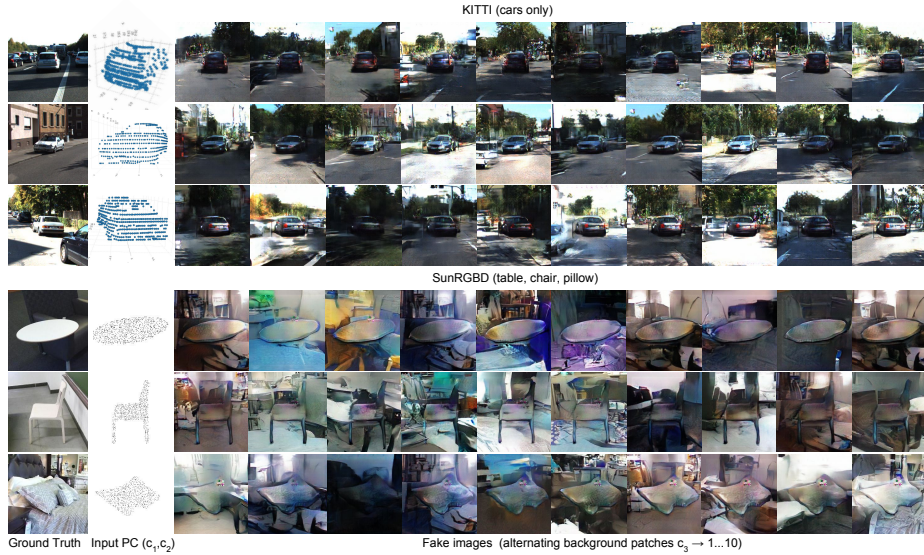


Fig. 3. Qualitative results of Points2Pix: The figure shows four different classes (cars (top) \rightarrow 3 samples of KITTI; table, chair, pillow (bottom) \rightarrow in each case one sample of SunRGBD). The results are taken from the test-set and never seen during training. The left column shows the ground truth image and the corresponding point-cloud in the second column. Fake images are generated based on a constant point-cloud $c_1, c_2 = const$ and ten alternating background patches c_3 (column 3-10). The model retains 3D characteristics of the objects.

while keeping the point-cloud constant. An interesting point is, that our model learns 3D characteristics. This could be proven with different outputs (backgrounds) where the objects geometry stays constant. Note, even the objects color stays the same apart from slight differences in reflections and illuminations. This means, the model associates a color with a specific 3D shape represented within the 3D latent feature space. Hence, alternating backgrounds do not affect the objects representation (geometry, color).

Tables 4.3 and 4.3, as well as Fig. 4 show quantitative results based on our metrics described in 4.1. We achieve extreme positive results for KITTI (S_c , IoU) and sufficient values for SunRGBD. SunRGBD includes a higher number of occlusions which drastically affects the scores. Additionally, there are far less samples on each class compared to *cars* in KITTI. Qualitative results of the inception score are shown in Fig. 5.

Ablation study

Architectural Review For completeness, we test two derivative architectures of our full pipeline (Fig. 6). In this way, we successfully show a point-cloud to image translation only based on the point cloud itself (*PointNet only*). Doing this, the

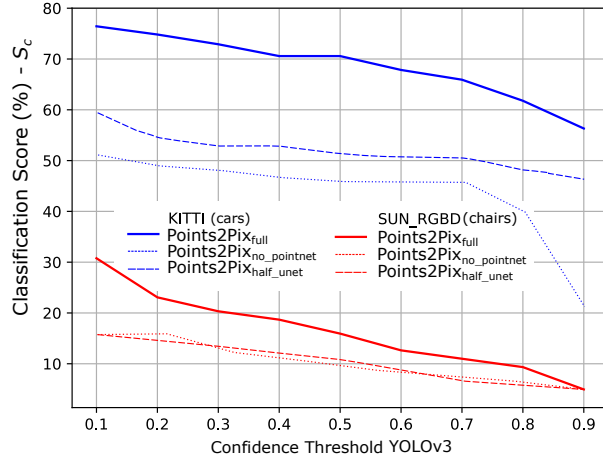


Fig. 4. Points2Pix classification score. The plot shows classification scores S_c for KITTI and SUN-RGBD of our full Points2Pix architecture as well as two derivative architectures (see Fig. 6) over confidence thresholds used for object detection with YOLOv3. The full architecture outperforms for KITTI as well as for SunRGBD.

dataset	class	$IoU_{\text{Points2Pix}}$		
		0.3	0.5	0.7
KITTI	car	0.76	0.77	0.77
		$IoU_{\text{Points2Pix}}$		
		0.1	0.2	0.3
SunRGBD	sofa	0.52	0.77	0.77
	table	0.70	-	-
	chair	0.60	0.58	0.58

Table 1. The object based inception score $IoU_{\text{Points2Pix}}$ is calculated on the test set for both data-sets. We show results for varying confidence thresholds, i.e. 0.3, 0.5, 0.7 for KITTI and 0.1, 0.2, 0.3 for SunRGBD.

dataset	class	$IoU_{\text{Points2Pix}}$		
		0.3	0.5	0.7
KITTI	car	0.71	0.70	0.68
		$IoU_{\text{Points2Pix}}$		
		0.1	0.2	0.3
SunRGBD	sofa	0.16	-	-
	table	0.24	0.22	-
	chair	0.45	0.37	0.33

Table 2. Diversity score is calculated on the test-set for both data-sets. Each sample is recomputed ten times with a random image background patches. A minus indicates no detections for the associated class.

whole training procedure runs much faster due to far less parameters to optimize. Nevertheless, sometimes a repeating noise with a high contrast similar to Moire effects appears, which indicates instabilities and uncertainties. Generated objects are in compliance with their 3D specifications, but in order to enlarge variance of the outputs and to control background conditions c_2 and c_3 are required. We found that the first part of the *UNet* and the view-point dependent projection c_2 especially help to reduce the mentioned noise effects. They provide additional information in 2D space and stabilize the network. As a fallback we addition-



Fig. 5. Qualitative object based inception results. The figure shows several generated cars and chairs (left) together with their accompanied real images (right). Green bounding boxes indicate detections on the real rgb image patches and red boxes visualize the corresponding ones on the fake images. The blue value obtains the IoU of both.

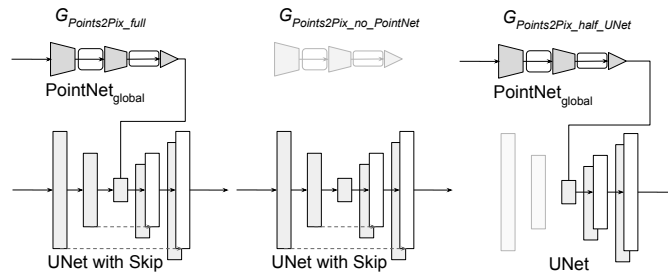


Fig. 6. Architectural review: Two derivatives from the basic Points2Pix $G(c_1, c_2, c_3)$ (left) generator are tested regarding their classification score (see Fig.4). One on hand a *Unet only* version $G(c_2, c_3)$ (middle), on the other hand a *PointNet only* $G(c_1)$ (right) version is tested. The full model outperforms the others.



Fig. 7. Learning 3D representations: The full Points2Pix architecture learns 3D representations. The model offers a high flexibility in generation of different view points by adjusting condition c_2 . The left part shows two examples of KITTI when rotating the point-cloud slightly by 20 degrees around the y-axis. The right (SunRGBD) shows the results when flipping the projection by 180 degrees around the x-axis.

ally test a *Unet only* version (Fig. 6). However, our full pipeline significantly outperforms the derivative architectures in terms of classification (Fig. 4).

Rotations To further emphasize the influence of c_2 and to show our models ability to constrain object view-points, we rotate all input points x_i for c_2 . We test that for KITTI with a rotation of 20 degrees around the y-axis and for SunRGBD with a rotation of 180 degrees around the x-axis (see Fig. 7). Note, that our point-cloud condition c_1 stays unmodified, because *PointNet* approximates a symmetric function to be invariant of rotations. The test shows that rotations can be implicitly learned. This offers many opportunities in generating 3D data.

5 Conclusion

In this work, we propose a novel approach for 3D point-cloud to image translation based on conditional GANs. Our network handles multi-modal sources from different domains and is capable of the translating unordered point-clouds to regular image grids. We use three conditions to generate a high diversity, while being flexible and keeping 3D characteristics. We prove that the model learns 3D characteristics, what even makes it possible to sample images from different viewpoints. Those networks are applicable in a wide variety of applications, especially 3D texturing.

References

1. Achlioptas, P., Diamanti, O., Mitliagkas, I., Guibas, L.J.: Representation learning and adversarial generation of 3d point clouds. CoRR **abs/1707.02392** (2017), <http://arxiv.org/abs/1707.02392>
2. Chang, A.X., Funkhouser, T.A., Guibas, L.J., 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. CoRR **abs/1512.03012** (2015), <http://arxiv.org/abs/1512.03012>
3. Chen, X., Ma, H., Wan, J., Li, B., Xia, T.: Multi-view 3d object detection network for autonomous driving. CoRR **abs/1611.07759** (2016), <http://arxiv.org/abs/1611.07759>
4. Denton, E.L., Chintala, S., Szlam, A., Fergus, R.: Deep generative image models using a laplacian pyramid of adversarial networks. CoRR **abs/1506.05751** (2015), <http://arxiv.org/abs/1506.05751>
5. Geiger, A.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: Proceedings of the 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 3354–3361. CVPR '12, IEEE Computer Society, Washington, DC, USA (2012), <http://dl.acm.org/citation.cfm?id=2354409.2354978>
6. Iizuka, S., Simo-Serra, E., Ishikawa, H.: Let there be Color!: Joint End-to-end Learning of Global and Local Image Priors for Automatic Image Colorization with Simultaneous Classification. ACM Transactions on Graphics (Proc. of SIGGRAPH 2016) **35**(4), 110:1–110:11 (2016)
7. Isola, P., Zhu, J., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. CoRR **abs/1611.07004** (2016), <http://arxiv.org/abs/1611.07004>
8. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. CoRR **abs/1412.6980** (2014), <http://arxiv.org/abs/1412.6980>
9. Ku, J., Mozifian, M., Lee, J., Harakeh, A., Waslander, S.: Joint 3d proposal generation and object detection from view aggregation. IROS (2018)
10. Li, B.: 3d fully convolutional network for vehicle detection in point cloud. CoRR **abs/1611.08069** (2016), <http://arxiv.org/abs/1611.08069>
11. Li, C., Zaheer, M., Zhang, Y., Póczos, B., Salakhutdinov, R.: Point cloud GAN. CoRR **abs/1810.05795** (2018), <http://arxiv.org/abs/1810.05795>
12. Lin, T., Maire, M., Belongie, S.J., Bourdev, L.D., Girshick, R.B., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: common objects in context. CoRR **abs/1405.0312** (2014), <http://arxiv.org/abs/1405.0312>
13. Luo, W., Yang, B., Urtasun, R.: Fast and furious: Real time end-to-end 3d detection, tracking and motion forecasting with a single convolutional net. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
14. Mirza, M., Osindero, S.: Conditional generative adversarial nets. CoRR **abs/1411.1784** (2014), <http://arxiv.org/abs/1411.1784>
15. Paszke, A., Chaurasia, A., Kim, S., Culurciello, E.: Enet: A deep neural network architecture for real-time semantic segmentation. CoRR **abs/1606.02147** (2016), <http://arxiv.org/abs/1606.02147>
16. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. arXiv preprint arXiv:1612.00593 (2016)
17. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from RGB-D data. CoRR **abs/1711.08488** (2017), <http://arxiv.org/abs/1711.08488>

18. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. CoRR **abs/1706.02413** (2017), <http://arxiv.org/abs/1706.02413>
19. Redmon, J., Farhadi, A.: Yolo3: An incremental improvement. CoRR **abs/1804.02767** (2018), <http://arxiv.org/abs/1804.02767>
20. Reed, S.E., Akata, Z., Yan, X., Logeswaran, L., Schiele, B., Lee, H.: Generative adversarial text to image synthesis. CoRR **abs/1605.05396** (2016), <http://arxiv.org/abs/1605.05396>
21. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. CoRR **abs/1505.04597** (2015), <http://arxiv.org/abs/1505.04597>
22. Salimans, T., Goodfellow, I.J., Zaremba, W., Cheung, V., Radford, A., Chen, X.: Improved techniques for training gans. CoRR **abs/1606.03498** (2016), <http://arxiv.org/abs/1606.03498>
23. Shelhamer, E., Long, J., Darrell, T.: Fully convolutional networks for semantic segmentation. IEEE Trans. Pattern Anal. Mach. Intell. **39**(4), 640–651 (2017). <https://doi.org/10.1109/TPAMI.2016.2572683>, <https://doi.org/10.1109/TPAMI.2016.2572683>
24. Simon, M., Milz, S., Amende, K., Gross, H.: Complex-yolo: Real-time 3d object detection on point clouds. CoRR **abs/1803.06199** (2018), <http://arxiv.org/abs/1803.06199>
25. Song, S., Lichtenberg, S.P., Xiao, J.: SUN RGB-D: A RGB-D scene understanding benchmark suite. Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. **07-12-June-2015**, 567–576 (2015). <https://doi.org/10.1109/CVPR.2015.7298655>
26. Song, S., Xiao, J.: Deep sliding shapes for amodal 3d object detection in RGB-D images. CoRR **abs/1511.02300** (2015), <http://arxiv.org/abs/1511.02300>
27. Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J., Catanzaro, B.: High-resolution image synthesis and semantic manipulation with conditional gans. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (2018)
28. Wang, X., Gupta, A.: Generative image modeling using style and structure adversarial networks. CoRR **abs/1603.05631** (2016), <http://arxiv.org/abs/1603.05631>
29. Wu, J., Wang, Y., Xue, T., Sun, X., Freeman, W.T., Tenenbaum, J.B.: MarrNet: 3D Shape Reconstruction via 2.5D Sketches. In: Advances In Neural Information Processing Systems (2017)
30. Xie, S., Tu, Z.: Holistically-nested edge detection. CoRR **abs/1504.06375** (2015), <http://arxiv.org/abs/1504.06375>
31. Yoo, D., Kim, N., Park, S., Paek, A.S., Kweon, I.: Pixel-level domain transfer. CoRR **abs/1603.07442** (2016), <http://arxiv.org/abs/1603.07442>
32. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. CoRR **abs/1612.01105** (2016), <http://arxiv.org/abs/1612.01105>
33. Zhou, Y., Tuzel, O.: Voxelnet: End-to-end learning for point cloud based 3d object detection. CoRR **abs/1711.06396** (2017), <http://arxiv.org/abs/1711.06396>
34. Zhu, J., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. CoRR **abs/1703.10593** (2017), <http://arxiv.org/abs/1703.10593>