# Deep Generative Modeling of LiDAR Data

Lucas Caccia[1,2], Herke van Hoof[1,4], Aaron Courville[2,3], Joelle Pineau[1,2,3]

*Abstract*— **Building models capable of generating structured output is a key challenge for AI and robotics. While generative models have been explored on many types of data, little work has been done on synthesizing lidar scans, which play a key role in robot mapping and localization. In this work, we show that one can adapt deep generative models for this task by unravelling lidar scans into a 2D point map. Our approach can generate high quality samples, while simultaneously learning a meaningful latent representation of the data. We demonstrate significant improvements against state-of-the-art point cloud generation methods. Furthermore, we propose a novel data representation that augments the 2D signal with absolute positional information. We show that this helps robustness to noisy and imputed input; the learned model can recover the underlying lidar scan from seemingly uninformative data.**

## I. INTRODUCTION

One of the main challenges in mobile robotics is the development of systems capable of fully understanding their environment. This non-trivial task becomes even more complex when sensor data is noisy or missing. An intelligent system that can replicate the data generation process is much better equipped to tackle inconsistency in its sensor data. There is significant potential gain in having autonomous robots equipped with data generation capabilities which can be leveraged for reconstruction, compression, or prediction of the data stream.

In autonomous driving, information from the environment is captured from sensors mounted on the vehicle, such as cameras, radars, and lidars. While a significant amount of research has been done on generating RGB images, relatively little work has focused on generating lidar data. These scans, represented as an array of three dimensional coordinates, give an explicit topography of the vehicle's surroundings, potentially leading to better obstacle avoidance, path planning, and inter-vehicle spatial awareness.

To this end, we leverage recent advances in deep generative modeling, namely variational autoencoders (VAE) [1] and generative adversarial networks (GAN) [2], to produce a generative model of lidar data. While the VAE and GAN approaches have different objectives, they can be used in conjunction with Convolutional Neural Networks (CNN) [3] to extract local information from nearby sensor points.

Unlike some approaches for lidar processing, we do not convert the data to voxel grids [4], [5]. Instead, we build off existing work [6] which projects the lidar scan into a 2D spherical point map. We show that this representation is fully

[1] MILA, McGill University
[2] MILA, Université de Montréal
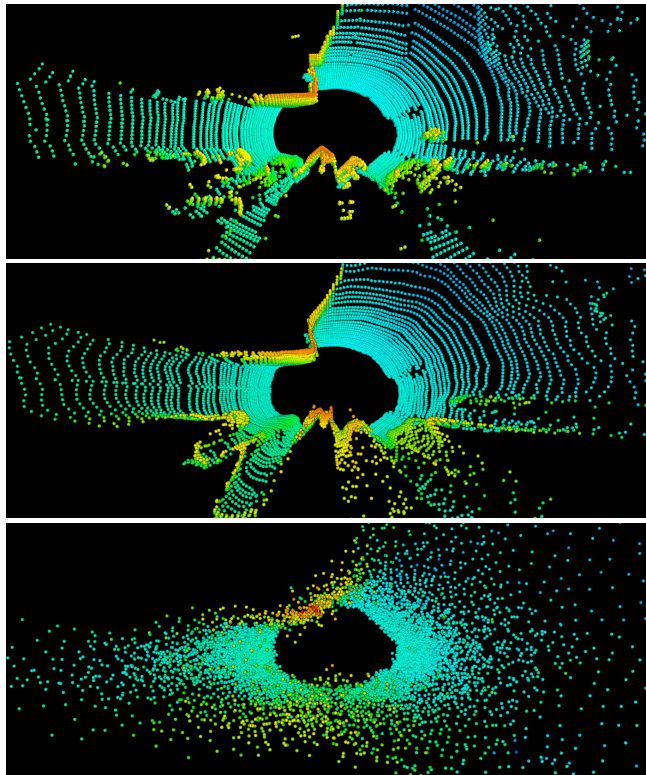[3] CIFAR Fellow
[4] University of Amsterdam

Fig. 1: Best viewed in color. Top: real LiDAR sample from the test set. Middle: reconstruction from our proposed model. Bottom: reconstruction from the baseline model.

compatible with deep architectures previously designed for image generation. Moreover, we investigate the robustness of this approach to missing or noisy data, a crucial property for real world applications. We propose a simple, yet effective way to improve the model's performance when the input is degraded. Our approach consists of augmenting the 2D map with absolute positional information, through extra $(x, y, z)$ coordinate channels. We validate these claims through a variety of experiments on the KITTI [7] dataset.

Our contributions are the following:

- We provide a fully unsupervised method for both *conditional* and *unconditional* lidar generation.
- We establish an evaluation framework for lidar reconstruction, allowing the comparison of methods over a spectrum of different corruption mechanisms.
- We propose a simple technique to help the model process noisy or missing data.

## II. RELATED WORK

### A. Lidar processing using Deep Learning

The majority of papers applying deep learning methods to lidar data present discriminative models to extract relevant information from the vehicle's environment. Dewan et al. [8] propose a CNN for pointwise semantic segmentation to distinguish between static and moving obstacles. Caltagirone et al. [9] use a similar approach to perform pixel-wise classification for road detection. To leverage the full 3D structure of the input, Bo Li [10] uses 3D convolutions on a voxel grid for vehicle detection. However processing voxels is computationally heavy, and does not leverage the sparsity of LiDAR scans. Engelcke et al. [5] propose an efficient 3D convolutional layer to mitigate these issues.

Another popular approach [6], [11]–[13] to avoid using voxels relies on the inherent two-dimensional nature of lidars. It consists of a bijective mapping from 3D point cloud to a 2D point map, where $(x, y, z)$ coordinates are encoded as azimuth and elevation angles measured from the origin. This can also be seen as projecting the point cloud onto a 2D spherical plane. Using such a bijection lies at the core of our proposed approach for generative modeling of lidar data.

### B. Grid-based lidar generation

An alternative approach for generative modeling of lidar data is from Ondruska et al [14]. They train a Recurrent Neural Network for semantic segmentation and convert their input to an occupancy grid. More relevant to our task, they train their network to also predict future occupancy grids, thereby creating a generative model for lidar data. Their approach differs from ours, as the occupancy grid used assigns a constant area (400 cm$^2$) to every slot, whereas we operate directly on projected coordinates. This not only reduces preprocessing time, but also allows us to efficiently represent data with non-uniform spatial density. We can therefore run our model at a much higher resolution, while remaining computationally efficient.

Concurrent with our work, Tomasello et al. [15] explore conditional lidar synthesis from RGB images. The authors use the same 2D spherical mapping proposed in [6]. Our approach differs on several points. First, we do not require any RGB input for generation, which may not always be available (e.g. in poorly lit environments). Second, we explore ways to augment the lidar representation to increase robustness to corrupted data. Finally, we look at generative modeling of lidar data (compared to a deterministic mapping in their case).

### C. Point Cloud Generation

A recent line of work [16]–[19] considers the problem of generating point clouds as *unordered sets* of $(x, y, z)$ coordinates. This approach does not define an ordering on the points, and must therefore be invariant to permutations. To achieve this, they use a variant of PointNet [20] to encode a variable-length point cloud into a fixed-length representation. This latent vector is then decoded back to a point cloud, and the whole network is trained using permutation invariant losses such as the *Earth-Mover's Distance* or the *Chamfer Distance* [19]. While these approaches work well for arbitrary point clouds, we show that they give suboptimal performance on lidar, as they do not leverage the known structure of the data.

### D. Improving representations through extra coordinate channels

In this work, we propose to augment the 2D spherical signal with Cartesian coordinates. This can be seen as a generalization of the CoodConv solution [21]. The authors propose to add two channels to the image input, corresponding to the $(i, j)$ location of every pixel. They show that this enables networks to learn either complete translation invariance or varying degrees of translation dependence, leading to better performance on a variety of downstream tasks.

## III. TECHNICAL BACKGROUND : GENERATIVE MODELING

The underlying task of generative models is density estimation. Formally, we are given a set of $d$-dimensional *i.i.d* samples $X = \{x_i \in \mathbb{R}^d\}_{i=1}^m$ from some unknown probability density function $p_{\text{real}}$. Our objective is to learn a density $p_\theta$ where $\theta \in \mathcal{F}$ represents the parameters of our estimator and $\mathcal{F}$ a parametric family of models. Training is done by minimizing some distance $\mathcal{D}$ between $p_{\text{real}}$ and $p_\theta$. The choice of both $\mathcal{D}$ and the training algorithm are the defining components of the density estimation procedure. Common choices for $\mathcal{D}$ are either $f$-divergences such as the Kullback-Liebler (KL) divergence, or Integral Probability Metrics (IPMs), such as the Wasserstein metric [22]. These similarity metrics between distributions often come with specific training algorithms, as we describe next.

### A. Maximum Likelihood Training

Maximum likelihood estimation (MLE) aims to find model parameters that maximize the likelihood of $X$. Since samples are *i.i.d*, the optimization criterion can be viewed as :

$$\max_{\theta \in \mathcal{F}} \mathbb{E}_{x \sim p_{\text{real}}} \log(p_\theta(x)). \tag{1}$$

It can be shown that training with the MLE criteria converges to a minimization of the KL-divergence as the sample size increases [23]. From Eqn (1) we see that any model admitting a differentiable density $p_\theta(x)$ can be trained via backpropagation. Powerful generative models trained via MLE include Variational Autoencoders [1] and autoregressive models [24]. In this work, we focus on the former, as the latter have slow sampling speed, limiting their potential use for real world applications.
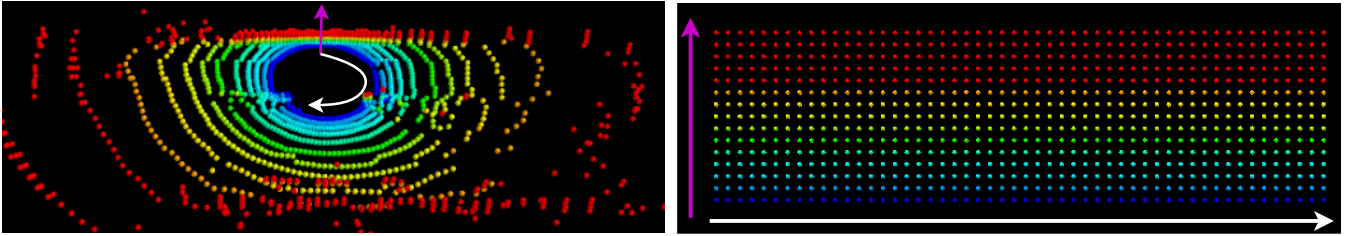
Fig. 2: Best viewed in color. Our proposed ordering of points from 3D space (left) into a 2D grid (right). Points sampled from the same elevation angle share the same color. The ordering of every row is obtained by *unrolling* points in increasing azimuth angle. The showed lidar was downsampled for visual purposes.

*1) Variational Autoencoders (VAE):* The VAE [1] is a regularized version of the traditional autoencoder (AE). It consists of two parts: an inference network $\phi_{enc} \equiv q(z|x)$ that maps an input x to a posterior distribution of latent codes $z$, and a generative network $\psi_{dec} \equiv p(x|z)$ that aims to reconstruct the original input conditioned on the latent encoding. By imposing a prior distribution $p(z)$ on latent codes, it enforces the distribution over $z$ to be smooth and well-behaved. This property enables proper sampling from the model via ancestral sampling from latent to input space.

The full objective of the VAE is then:

$$\mathcal{L}(\theta; x) = \mathbb{E}_{q_{(z|x)}} \log p(x|z) - \text{KL}(q(z|x)||p(z)) \leq \log p(x),$$
(2)

which is a valid lower bound on the true likelihood, thereby making Variational Autoencoders valid generative models. For a more in depth analysis of VAEs, see [25].

### B. Generative Adversarial Network (GAN)

The GAN [2] formulates the density estimation problem as a minimax game between two opposing networks. The *generator* $G(z)$ maps noise drawn from a prior distribution $p_{\text{noise}}$ to the input space, aiming to fool its adversary, the *discriminator* $D(x)$. The latter then tries to distinguish between real samples $x \sim p_{\text{real}}$ and fake samples $x' \sim G(z)$. In practice, both models are represented as neural networks. Formally, the objective is written as

$$\min_G \max_D \mathbb{E}_{x \sim p_{\text{real}}} \log(D(x)) + \mathbb{E}_{z \sim p_{\text{noise}}} \log(1 - D(G(z))). \quad (3)$$

GANs have shown the ability to produce more realistic samples [26] than their MLE counterparts. However, the optimization process is notoriously difficult; stabilizing GAN training is still an open problem. In practice, GANs can also suffer from *mode collapse* [27], which happens when the generator overlooks certain modes of the target distribution.

### IV. PROPOSED APPROACH FOR LIDAR GENERATION

We next describe the proposed deep learning framework used for generative modeling of lidar scans.

### A. Data Representation

Our approach relies heavily on 2D convolutions, therefore we start by converting a lidar scan containing $N$ $(x, y, z)$ coordinates into a 2D grid. We begin by clustering together points emitted from the same elevation angle into $H$ clusters.

Second, for every cluster, we sort the points in increasing order of azimuth angle. In order to have a proper grid with a fixed amount of points per row, we divide the $360°$ plane into $W$ bins. This yields a $H \times W$ grid, where for each cell we store the average $(x, y, z)$ coordinate, such that we can store all the information in a $H \times W \times 3$ tensor. We note that the default ordering in most lidar scanners is the same as the one obtained after applying this preprocessing. Therefore, sorting is not required in practice, and the whole procedure can be executed in $\mathcal{O}(N)$. Figure 2 provides a visual representation of this mapping. This procedure yields the same ordering of points as the projection discussed in II-A. The latter would then return a grid of $H \times W \times 2$, where the $(x, y)$ channels are compressed as $d = \sqrt{x^2 + y^2}$. We will refer to the two representations above as *Cartesian* and *Polar* respectively. While this small change in representation seems innocuous, we show that when the input is noisy or incomplete, this compression can lead to suboptimal performance.

### B. Training Phase

*1) VAEs:* In practice, both encoder $\phi$ and decoder $\psi$ are represented as neural networks with parameters $\theta_{enc}$ and $\theta_{dec}$ respectively.

Similar to a traditional AE, the training procedure first encodes the data $x$ into a latent representation $z = \phi(x; \theta_{enc})$. The variational aspect is introduced by interpreting $z$ not as a vector, but as parameters of a posterior distribution. In our work we choose a Gaussian prior and posterior, and therefore $z$ decomposes as $\mu_x, \sigma_x$.

We then sample from this distribution $\tilde{z} \sim \mathcal{N}(\mu_x, \sigma_x)$ and pass it through the decoder to obtain $\tilde{x} = \psi(\tilde{z}; \theta_{dec})$. Using the reparametrization trick [1], the network is fully deterministic and differentiable w.r.t its parameters $\theta_{enc}$ and $\theta_{dec}$, which are updated via stochastic gradient descent (SGD).

*2) GANs:* Training alternates between updates for the generator and discriminator, with parameters $\theta_{gen}$ and $\theta_{dis}$. Similarly to the VAE, samples are obtained by ancestral sampling from the prior through the generator. In the original GAN, the networks are updated according to Eqn. 3. In practice, we use the Relativistic Average GAN (RaGAN) objective [28], which is easier to optimize. Again, $\theta_{gen}$ and $\theta_{dis}$ are updated using SGD. For a complete hyperparameter list, we refer the reader to our publicly available source

code.[1]

### C. Model Architecture

Deep Convolutional GANs (DCGANs) [29] have shown great success in generating images. They use a symmetric architecture for the two networks: The generator consists of five transpose convolutions with stride two to upsample at each layer, and ReLU activations. The discriminator uses stride two convolutions to downsample the input, and Leaky ReLU activations. In both networks, Batch Normalization [30] is interleaved between convolution layers for easier optimization. We use this architecture for all our models: The VAE encoder setup is simply the first four layers of the discriminator, and the decoder's architecture replicates the DCGAN generator. We note that for both models, more sophisticated architectures [31], [32] are fully compatible with our framework. We leave this line of exploration as future work.

## V. Experiments

This section provides a thorough analysis of the performance of our framework fulfilling a variety of tasks related to generative modeling. First, we explore *conditional* generation, where the model must compress and reconstruct a (potentially corrupted) lidar scan. We then look at *unconditional* generation. In this setting, we are only interested in producing realistic samples, which are not explicitly tied to a real lidar cloud.

### A. Dataset

We consider the point clouds available in the KITTI dataset [7]. We use the train/validation/test set split proposed by [33], which yields 40 000, 80 and 700 samples for train, validation and test sets. We use the preprocessing described in section IV-A to get a $40 \times 256$ grid. For training we subsample from 10 Hz to 3 Hz since temporally adjacent frames are nearly identical.

### B. Baseline Models

Since, to the best of our knowledge, no work has attempted generative modeling of raw lidar clouds, we compare to our method models that operate on arbitrary point clouds. We first choose AtlasNet [17], which has shown strong modeling performance on the Shapenet [34] dataset. This network first encodes point clouds using a shared *MLP* network that operates on each point *individually*. A max-pooling operation is performed on the point axis to obtain a fixed-length global representation of the point cloud. In other words, the encoder treats each point independently of other points, without assuming an ordering on the set of coordinates. This makes the feature extraction process invariant to permutations of points. The decoder is given the encoder output along with $(x, y)$ coordinates of a 2D-grid, and attempts to *fold* this 2D-grid into a three-dimensional surface. The decoder also uses a *MLP* network shared across all points.

[1] https://www.github.com/pclucas14/lidar_generation

Similar to AtlasNet, we compare our model with the one from Achlioptas et al [16]. Only its decoder differs from AtlasNet; the model does not deform a 2D grid, but rather uses fully-connected layers to convert the latent vector into a point cloud, making it less parameter efficient.

Both networks are trained end-to-end using the *Chamfer Loss* [19], defined as

$$d_{\text{CH}} = \sum_{x \in S_1} \min_{y \in S_2} ||x - y||_2^2 + \sum_{y \in S_2} \min_{x \in S_1} ||x - y||_2^2, \quad (4)$$

where $S_1$ and $S_2$ are two sets of $(x, y, z)$ coordinates. We note again that this loss is invariant to the ordering of the output points. For both autoencoders, we regularize their latent space using a Gaussian prior to get a valid generative model.

### C. Conditional Generation

We proceed to test our approach in a conditional generation task. In this setting, we do not evaluate the GAN, as this family of model -in their original formulation- does not have an inference mechanism. In other words, we consider four models: our approach, using either the Cartesian or the Polar representation, and the two baselines above. Since we are not sampling, but rather reconstructing an input, we consider both VAE and AE variants of every model, and report the best performing one.

Formally, given a lidar cloud, we evaluate a model's ability to reconstruct it from a compressed encoding. More relevant to real word applications, we look at how robust the model's latent representation is to input perturbation. Specifically, we look at the two following corruption mechanisms:

- **Additive Noise** : we add Gaussian noise drawn from $\mathcal{N}(0, \sigma)$ to the $(x, y, z)$ coordinates of the lidar cloud. For this process, we normalize each of the three dimension independently prior to noise addition. We experiment with varying levels of $\sigma$.
- **Data Removal** : We remove random points from the input lidar scan. Specifically, the probability of removing a point is modeled as a Bernoulli distribution parametrized by $p$. We consider different values for $p$.

### D. Unconditional Generation

For this section, we consider the GAN model introduced in section IV-C. Our goal is to train a model that can produce realistic samples. Having access to such a generator can lead to better simulator development, which are heavily used to train self-driving agents [35]. In this use case, an agent operating in an environment that lacks *crispness* will likely result in poor skill transfer to real world navigation. Since the use of GANs has been shown to produce more realistic samples than MLE based models on images [36], we hope to see similar results with our model in the case of LiDAR data.

**Evaluation criteria**: Rigorous quantitative evaluation of samples produced by GANs and generative models is an open research question. GANs trained on images have been
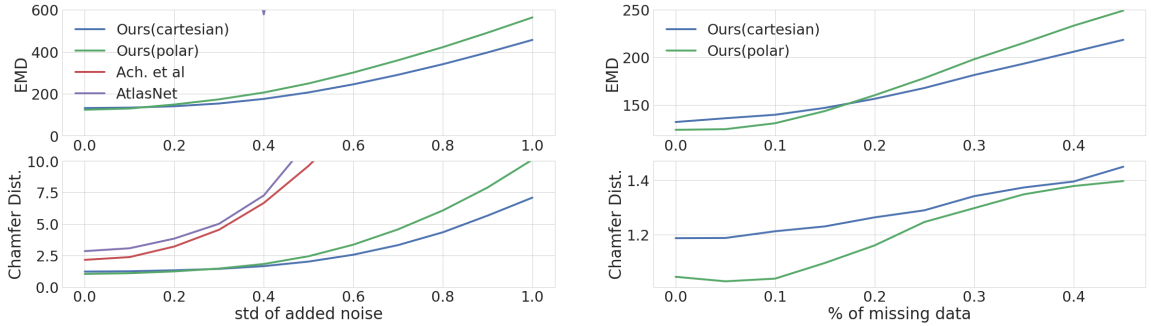
Fig. 3: EMD and Chamfer Distance under varying levels of added noise (left) and missing data (right). We remove models with poor performance for clarity. For both metrics lower is better.

evaluated by the Inception Score [27] and the Frechet Inception Distance (FID) [37]. Since there exists no standardized metric for unconditional generation of lidar clouds, we rely on visual inspection of samples for quality assessment.

*1) Evaluation criteria:* To measure how close the reconstructed output is to the original point cloud, we use the *Earth-Mover's Distance* [19]. It is defined as

$$d_{\text{EMD}}(S_1, S_2) = \min_{\gamma: S_1 \longrightarrow S_2} \sum_{x \in S_1} ||x - \gamma(x)||_2 \qquad (5)$$

where $\gamma$ is a bijection between the two sets.

The EMD gives the solution to the optimal transportation problem, which attempts to transform one point cloud into the other. Recent work [16] has shown that this metric correlates well with human evaluation, and does so better than the Chamfer Distance. Moreover, the Earth Mover's Distance is sensitive to both global and local structure, and does not require points to be ordered. Additionally, training and evaluating models on the same metric can result in models overfitting to this criterion, at the expense of sample quality [38]. Nevertheless, we also provide results measured by the Chamfer Distance for completeness.

*2) Training Protocol:* For every model considered, we perform the same hyperparameter search. We randomly select the learning rate, the latent dimension and the batch size from a predetermined set of values. This set of values is the same for all models to ensure fairness. This process is repeated for 10 different configurations, from which we choose the one obtaining the best performance on the validation set. We then proceed to evaluate this configuration on the test set according to the metrics described above. All models are trained end-to-end on the same dataset.

## VI. RESULTS

In this section, we will first discuss results for conditional generation and subsequently evaluate results for unconditional generation of lidar images.

### A. Conditional

In all conditional tasks, our proposed approach beats available baselines by a significant margin, both in terms of EMD, Chamfer Distance and visual inspection.
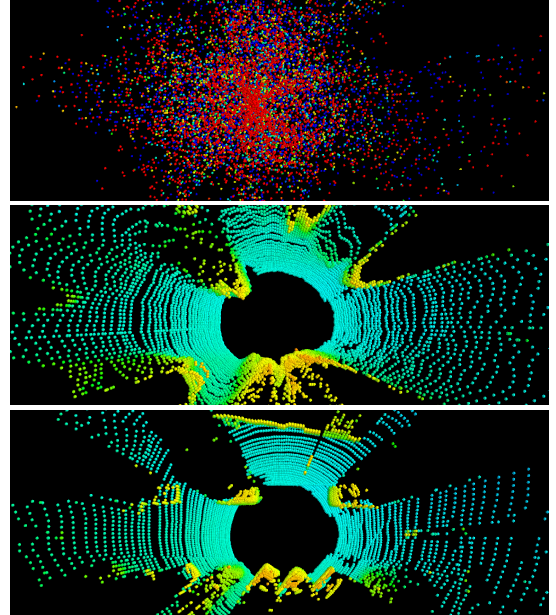


Fig. 4: Top : corrupted lidar from the test set, where we added noise drawn from $\mathcal{N}(0, 0.8)$ on the preprocessed scan. Middle : reconstructed point cloud given corrupted input. Bottom : original lidar scan

*1) Reconstructing clean data:* while the baseline models are able to reconstruct the global structure of the lidar scan, they are unable to recover the more fine grained detail of the input (see Fig.1). This suggests that leveraging the known structure of the lidar plays a key role in obtaining high quality reconstructions. Quantitative results are shown in Table I.

| Model | EMD | Chamfer |
|---|---|---|
| Random | 4331.9 | 253.6 |
| AtlasNet | 1571.2 | 2.85 |
| Ach. et al | 1103.1 | 2.16 |
| Ours(xyz) | 137.2 | 1.23 |
| **Ours(pol)** | **127.0** | **1.04** |

TABLE I: EMD and Chamfer distance measured on test set reconstructions (in both cases lower is better)

*2) Reconstructing corrupted data:* Next, we evaluate the proposed models on their ability to extract important in-
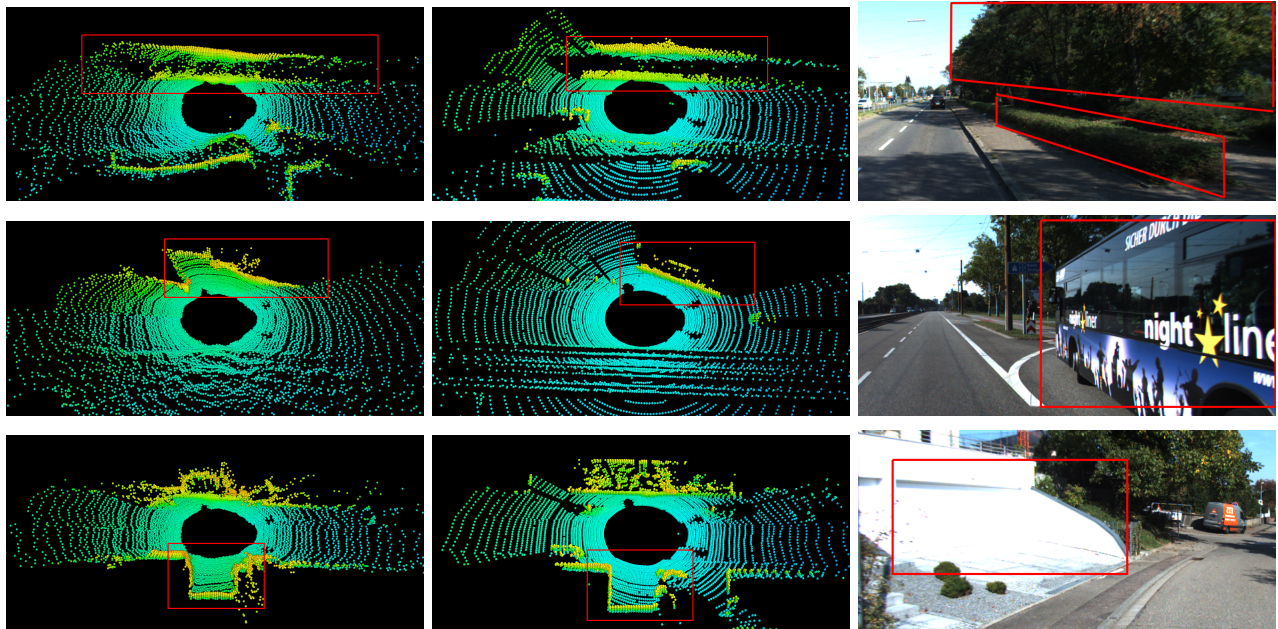
Fig. 5: We compare generated GAN samples (left) with their nearest neighbor in feature space (middle) from the test set. We show the corresponding RGB image (right). Regions of interest are highlighted in red.

formation from corrupted lidar scans. As shown in Fig. 4, the proposed VAE correctly reconstructs the defining components of the original cloud, even if the given input is seemingly uninformative. We emphasize that **our model was not trained with such corrupted data**, therefore these results are quite surprising. Animations and additional reconstructions can be found here .

Moreover, we observe that as soon as the input is moderately noisy, the proposed Cartesian representation yields better performance. As seen in Fig. 3, this representation performs better than its Polar alternative over the majority of the graph. In addition, we observe a similar trend when points are randomly removed from the input, as shown in Fig. 3; when more than 15% of the points are missing, using $(x, y, z)$ coordinates performs favorably according to EMD. This result suggests that in this corruption regime, having access to absolute positional information provides a better signal to the model. Interesting future work would be to leverage the best of the two representations.

We note that the suboptimal performance of the baselines is mainly due to two factors. First, since points are encoded independently, only information about the global structure is kept, and local fine-grained details are neglected. Second, the Chamfer Distance used for training assumes that the point cloud has a uniform density, which is not the case for lidar scans.

### B. Unconditional

We perform a visual inspection of generated samples, located in the leftmost column of Figure 5 (more samples are available here). We see that our model generates realistic samples. First, the scans have a well-defined global structure: an aerial view of the samples show points correctly aligned

to model the structure of the road. Second, the samples share local characteristics of real data: the model correctly generates road obstacles, such as cars, or cyclists. This amounts to having locations with a dense aggregation of points, followed by a trailing area with almost no points, similar to the shadow of an object. Third, model respects the point density of the data, where the density is roughly inversely proportional to the distance from the origin. Lastly, our models show good sample diversity.

*1) What is the GAN generating?:* In order to better interpret samples from the unconditional generator, we try to match them to real data examples. We perform the following procedure: we encode every sample to a latent representation, given by the output of the third layer of our discriminator. We similarly encode random datapoints from the test set, and match the generated sample to the real datapoint yielding the smallest latent L2 loss. We show three examples of this matching in Figure 5. In the first row, we see the model generating a two layer roadside to the right, consisting of a long shrub, followed by a line of trees. On the second row, we find a large tilted object to the right, which matches a bus turning right. Finally, on the last row we see a sharp enclosing, corresponding to a driveway leading to a garage door.

### VII. DISCUSSION AND FUTURE WORK

In this work we introduced two generative models for *raw* lidars, a GAN and a VAE. We have shown that the proposed adversarial network can generate highly realistic data, and captures both local and global features of real lidar scans. The LiDAR-VAE successfully encodes and reconstructs lidar samples, and is highly robust to missing or inputed data.

We demonstrate that when adding enough noise to render the scan uninformative to the human eye, the proposed VAE still extracts relevant information and generates the missing data. Our work in deep generative modeling of lidar enables concrete advancements in real life applications; the former model can help reduce the discrepancy between synthetic and real lidars in driving simulators, while the latter can be leveraged in deployed vehicles for reconstruction, compression, or prediction of the data stream.

Moreover, we proposed a simple way to encode absolute positional information in the lidar representation, and showed that this leads to better reconstructions when the input is noisy or incomplete. Interesting future work would be to see if this can also lead to improvements in standard lidar processing tasks.

## REFERENCES

[1] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *Proceedings of the 2nd International Conference on Learning Representations.*, 2013.

[2] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.

[3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[4] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3D object detection," *arXiv preprint arXiv:1711.06396*, 2017.

[5] M. Engelcke, D. Rao, D. Z. Wang, C. H. Tong, and I. Posner, "Vote3deep: Fast object detection in 3d point clouds using efficient convolutional neural networks," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 1355–1361.

[6] B. Li, T. Zhang, and T. Xia, "Vehicle detection from 3D lidar using fully convolutional network," *arXiv preprint arXiv:1608.07916*, 2016.

[7] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.

[8] A. Dewan, G. L. Oliveira, and W. Burgard, "Deep semantic classification for 3D lidar data," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 3544–3549.

[9] L. Caltagirone, S. Scheidegger, L. Svensson, and M. Wahde, "Fast lidar-based road detection using fully convolutional neural networks," *arXiv preprint arXiv:1703.03613*, 2017.

[10] B. Li, "3D fully convolutional network for vehicle detection in point cloud," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 1513–1518.

[11] M. Velas, M. Spanel, M. Hradis, and A. Herout, "Cnn for very fast ground segmentation in velodyne lidar data," in *Autonomous Robot Systems and Competitions (ICARSC), 2018 IEEE International Conference on*. IEEE, 2018, pp. 97–103.

[12] V. Vaquero, I. del Pino, F. Moreno-Noguer, J. Solà, A. Sanfeliu, and J. Andrade-Cetto, "Deconvolutional networks for point-cloud vehicle detection and tracking in driving scenarios," in *Mobile Robots (ECMR), 2017 European Conference on*. IEEE, 2017, pp. 1–7.

[13] V. Vaquero, A. Sanfeliu, and F. Moreno-Noguer, "Deep lidar cnn to understand the dynamics of moving vehicles," *arXiv preprint arXiv:1808.09526*, 2018.

[14] P. Ondruska, J. Dequaire, D. Z. Wang, and I. Posner, "End-to-end tracking and semantic segmentation using recurrent neural networks," *arXiv preprint arXiv:1604.05091*, 2016.

[15] P. Tomasello, S. Sidhu, A. Shen, M. W. Moskewicz, N. Redmon, G. Joshi, R. Phadte, P. Jain, and F. Iandola, "Dscnet: Replicating lidar point clouds with deep sensor cloning," *arXiv preprint arXiv:1811.07070*, 2018.

[16] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas, "Representation learning and adversarial generation of 3D point clouds," *arXiv preprint arXiv:1707.02392*, 2017.

[17] T. Groueix, M. Fisher, V. G. Kim, B. C. Russell, and M. Aubry, "Atlasnet: A papier-mâché approach to learning 3D surface generation," *arXiv preprint arXiv:1802.05384*, 2018.

[18] Y. Yang, C. Feng, Y. Shen, and D. Tian, "Foldingnet: Point cloud autoencoder via deep grid deformation," in *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, vol. 3, 2018.

[19] H. Fan, H. Su, and L. J. Guibas, "A point set generation network for 3D object reconstruction from a single image." in *CVPR*, vol. 2, no. 4, 2017, p. 6.

[20] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3D classification and segmentation," *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, vol. 1, no. 2, p. 4, 2017.

[21] R. Liu, J. Lehman, P. Molino, F. P. Such, E. Frank, A. Sergeev, and J. Yosinski, "An intriguing failing of convolutional neural networks and the coordconv solution," in *Advances in Neural Information Processing Systems*, 2018, pp. 9628–9639.

[22] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," *arXiv preprint arXiv:1701.07875*, 2017.

[23] S. Kolouri, G. K. Rohde, and H. Hoffmann, "Sliced wasserstein distance for learning gaussian mixture models," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3427–3436.

[24] A. v. d. Oord, N. Kalchbrenner, and K. Kavukcuoglu, "Pixel recurrent neural networks," *arXiv preprint arXiv:1601.06759*, 2016.

[25] C. Doersch, "Tutorial on variational autoencoders," *arXiv preprint arXiv:1606.05908*, 2016.

[26] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," *arXiv preprint arXiv:1710.10196*, 2017.

[27] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *Advances in Neural Information Processing Systems*, 2016, pp. 2234–2242.

[28] A. Jolicoeur-Martineau, "The relativistic discriminator: a key element missing from standard gan," *arXiv preprint arXiv:1807.00734*, 2018.

[29] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[30] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

[31] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," *arXiv preprint arXiv:1805.08318*, 2018.

[32] D. P. Kingma, T. Salimans, R. Jozefowicz, X. Chen, I. Sutskever, and M. Welling, "Improved variational inference with inverse autoregressive flow," in *Advances in neural information processing systems*, 2016, pp. 4743–4751.

[33] W. Lotter, G. Kreiman, and D. Cox, "Deep predictive coding networks for video prediction and unsupervised learning," *arXiv preprint arXiv:1605.08104*, 2016.

[34] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su *et al.*, "Shapenet: An information-rich 3D model repository," *arXiv preprint arXiv:1512.03012*, 2015.

[35] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "Carla: An open urban driving simulator," *arXiv preprint arXiv:1711.03938*, 2017.

[36] A. B. L. Larsen, S. K. Sønderby, H. Larochelle, and O. Winther, "Autoencoding beyond pixels using a learned similarity metric," *arXiv preprint arXiv:1512.09300*, 2015.

[37] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "Gans trained by a two time-scale update rule converge to a local nash equilibrium," in *Advances in Neural Information Processing Systems*, 2017, pp. 6626–6637.

[38] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.