

SLPC: a VRNN-based approach for stochastic lidar prediction and completion in autonomous driving

George Eskandar¹, Alexander Braun², Martin Meinke², Karim Armanious¹ and Bin Yang¹

¹University of Stuttgart, Institute of Signal Processing and System Theory

Stuttgart, Germany

²Robert Bosch GmbH, Lidar Perception

Stuttgart, Germany

Abstract—Predicting future 3D LiDAR pointclouds is a challenging task that is useful in many applications in autonomous driving such as trajectory prediction, pose forecasting and decision making. In this work, we propose a new LiDAR prediction framework that is based on generative models namely Variational Recurrent Neural Networks (VRNNs), titled Stochastic LiDAR Prediction and Completion (SLPC). Our algorithm is able to address the limitations of previous video prediction frameworks when dealing with sparse data by spatially inpainting the depth maps in the upcoming frames. Our contributions can thus be summarized as follows: we introduce the new task of predicting and completing depth maps from spatially sparse data, we present a sparse version of VRNNs and an effective self-supervised training method that does not require any labels. Experimental results illustrate the effectiveness of our framework in comparison to the state of the art methods in video prediction.

Index Terms—Video Prediction, Depth completion, LiDAR pointcloud, Autonomous Driving, Deep Learning

I. INTRODUCTION

Autonomous cars continuously use depth information to capture the geometry of their environment. Many vision tasks such as object detection, scene flow, path planning and segmentation can be built on this input. Light detection and ranging technology (LiDAR) is heavily integrated in the autonomous driving industry as a remote sensing method by virtue of its high accuracy and long range. However, in bad weather conditions such as heavy rain, snow, fog, low hanging clouds or high sun angles, LiDAR can suffer from reduced precision due to reflection and refraction of light beams, unlike radar sensors which are more robust in these situations [1, 2]. Additionally, many LiDAR sensors currently available on the market do have a limited number of horizontal scanning lines (e.g. 16, 64, 128) which provides sparse depth measurements, thus capturing the environmental geometry with a precision that can still be improved. As such, predicting LiDAR measurement of future frames is desirable not just in academic research but also in industrial applications as in autonomous driving. For instance, on a short timescale it can work as a fallback real-time system in case of critical weather conditions and/or sensor failure.

Despite its potential, LiDAR pointcloud prediction poses several technical challenges. First, unlike images, LiDAR scans are sparse and do not present strong local features like edges and corners. Second, it's hard to simultaneously learn

the depth information, the distribution of the points projected on a camera frame (or sparsity shape) and the distribution of possible outcomes. Thirdly, we want to predict LiDAR scans at a high resolution, which usually makes the prediction task more challenging because a larger pixel-space leads to more uncertainty in the prediction. Thus, our research lies at the intersection of video prediction and depth inpainting, also known as depth completion.

The video prediction task is currently a hot topic in computer vision. Deep learning techniques have been recently used to predict future frames in videos. There have been two big approaches to video prediction: the first is to directly generate the pixels of the new frame like [3–6]. The second is to predict motion vectors and learn a transformation of the past pixels, like in [7–9]. What these frameworks have in common is that they use a combination of convolutional (CNNs and 3D-CNNs) and recurrent neural networks (LSTMs, GRUs, convLSTMs) to predict subsequent measurements. More recent frameworks use generative approaches like VAEs and Generative Adversarial Networks (GANs) [10–14] and learn a time-dependent multimodal probability distribution in pixel-space. In [15], the Variational Recurrent Neural Network (VRNN) was introduced where a prior probability distribution is defined over latent random variables in the hidden state of an RNN. Sampling from the prior leads to different outcomes. State-of-the-art video prediction consists of variations of VRNNs like [16–18], as they model uncertainty in the prediction in contrast to the deterministic approaches.

On the other hand, Depth Completion refers to the task of spatially inpainting sparse depth maps. Traditionally, pointcloud is projected onto the camera frame resulting in an irregular depth map. Research in this field consists in developing network architectures which produce accurate dense depth maps [19]. Recently, Patil et al [20] demonstrated that considering temporal correlation improves depth completion. In autonomous driving and robotics, the prediction of the future dynamic environment using deep learning has been investigated from two main aspects. Some methods involve predicting a dynamic occupancy grid (a pixel is occupied if there is a rigid object or obstacle) like [21] and [22]. On the other hand, some recent frameworks estimate absolute depth values like [23] which predict subsequent LiDAR measurements for indoor robotic navigation. The most recent

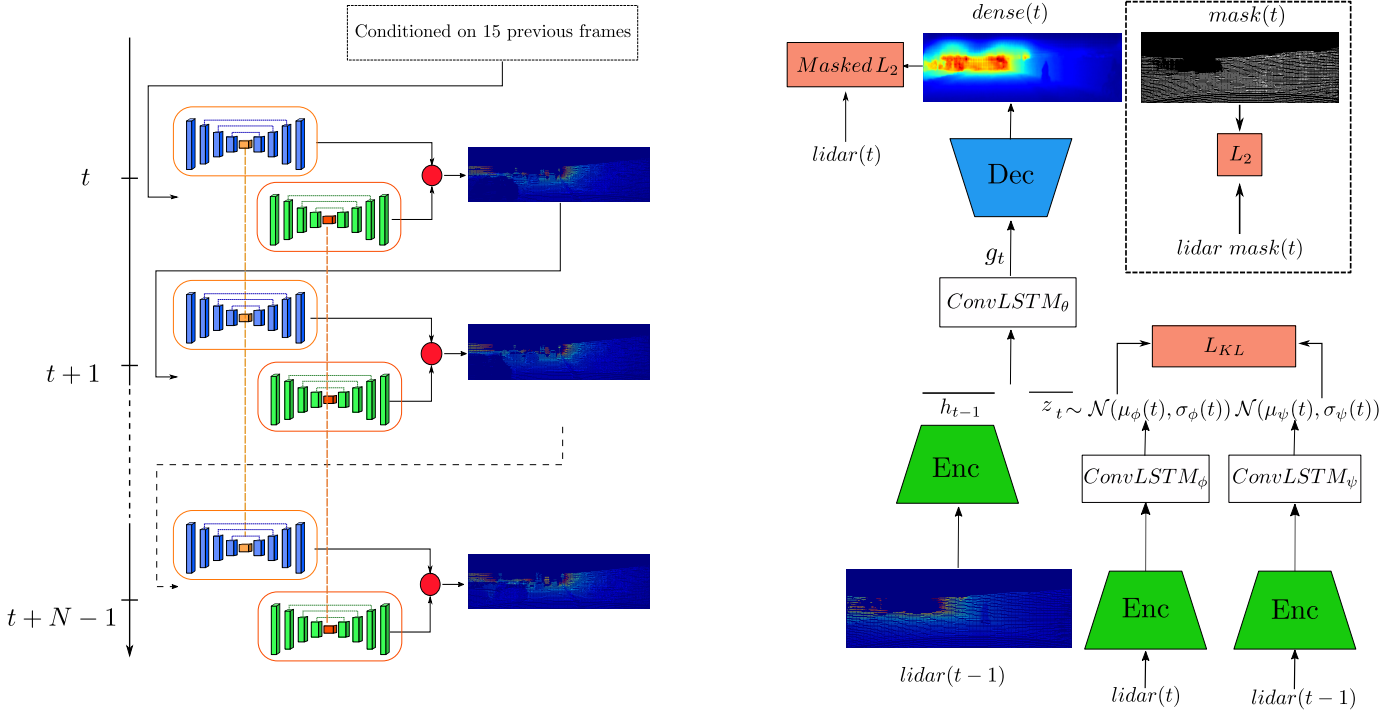


Fig. 1: Left: In inference, the Depth Network predicts a more complete Depth map while the Mask network predicts a sparsity mask that downsamples the dense depth map. Right: Network architecture and losses. The two parallel networks, depth and mask, share the same architecture but use different loss functions as shown in the box with dashed boundaries.

framework, Scene PointCloud Sequence Forecasting (SPCSF) [24], generates future LiDAR pointclouds in a determined way using the Chamfer Distance loss function [25], by first predicting future range map frames and projecting them in 3D.

In this work, we introduce the task of predicting dense depth maps from previous sparse LiDAR measurements. This would enable the generation of more reliable depth information in bad weather condition, by leveraging temporal correlation between past frames. We use the depth map representation, which projects the pointcloud in the front view of the car on the camera frame of reference, in contrast to the the range map representation in [24]. This is better suited to autonomous driving pipelines which perform sensor fusion between RGB monoculars and LiDAR sensors. To tackle the challenge, our model extends the state-of-the-art of video prediction with a modified architecture, losses and a novel training setup, tailored to the sparse LiDAR maps. Not only do we predict future pointclouds, but we also inpaint them spatially to have dense depth maps that are richer in geometric information. A thorough quantitative and qualitative analysis is conducted between the proposed framework and the state-of-the-art baselines in both video and LiDAR prediction frameworks for single and multiple frame prediction.

II. METHODS

In this section, we introduce the proposed SLPC framework for Stochastic LiDAR Prediction and Completion. This framework utilizes as baseline the Stochastic Video Prediction with

Learned Prior SVG-LP model. However, we expand on the prior framework with a different architecture, loss functions and training methods that are tailored to the sparse nature of our data. This should enable the prediction of subsequent LiDAR frames in a more robust manner.

A. VRNN with learned prior

A VRNN is a variational RNN that generates future frames Q_t based on the encoded previous frames $lidar_{1..t-1}$ and a latent variable z_t sampled from a prior probability distribution $p(z)$. The latter encodes stochastic spatio-temporal information not captured by the deterministic encoder. In our baseline SVG-LP, this prior distribution, $q_\psi(z_t|x_{1..t-1})$, is chosen to be a conditional time-varying Gaussian distribution with fixed variance $\mathcal{N}(\mu_\psi(t), \sigma_\psi(t))$. During training, the prior is forced to be similar to the posterior distribution $q_\phi(z_t|x_{1..t}) = \mathcal{N}(\mu_\phi(t), \sigma_\phi(t))$ by minimizing the KL-divergence between them. As depicted in Fig. 1, the whole model thus consists of four parts:

- a frame encoder (Enc) that encodes the previous frame $lidar_{t-1}$ into a variable in a latent space h_{t-1} .
- a prior/posterior model ($ConvLSTM_\psi$, $ConvLSTM_\phi$) maps the latent variable into a probability distribution $\mathcal{N}(\mu_\psi(t), \sigma_\psi(t)) / \mathcal{N}(\mu_\phi(t), \sigma_\phi(t))$.
- a frame predictor ($ConvLSTM_\theta$) encodes the combined h_{t-1} and $z_t \sim \mathcal{N}(\mu(t), \sigma(t))$ into a state g_t .
- a decoder (Dec) that generates $dense_t / mask_t$ from g_t and skip connections from encoder layers.

B. The sparsity problem

The sparsity problem of projected pointclouds is addressed in three ways: the architecture, losses and the training method. Since our goal is to predict and complete a sequence of sparse LiDAR frames $\hat{Q}(t)$, we formulate our problem as $\hat{Q}(t) = dense(t) \odot mask(t)$, where $dense(t)$ is a dense depth map and $mask(t)$ is a mask that downsamples the dense map to obtain a sparse pointcloud. We use a divide-and-conquer approach and use 2 separate VRNNs to predict each of $dense(t+1)$ and $mask(t+1)$. The choice of this formulation allows for more accurate and complete depth maps all while having a sparse pointcloud that can then be fed back as input of the model to generate more frames in an autoregressive manner.

The architecture and input of the **Depth** and **Mask** networks are the same. For more robustness against the sparse LiDAR, the encoder has three sparse convolution layers [19]. The operation is a masked convolution that normalizes the output over the number of pixels having valid depth values, and it has been shown in [19] to lead to less noisy results. The operation can be mathematically described in Eq. (1) as:

$$f_{u,v}(\mathbf{y}, \mathbf{o}) = \frac{\sum_{i,j=-k}^k y_{u+i,v+j} o_{u+i,v+j} w_{i,j}}{\epsilon + \sum_{i,j=-k}^k o_{u+i,v+j}} + b \quad (1)$$

where o is a mask indicating the pixels that contain a valid depth value, y is the input feature to the convolution, w is the learned kernel, k is the kernel size, b is the bias, and $\epsilon > 0$ prevents division by 0. We choose a large kernel size of $k=11$ to have a big field of view before the encoder. The sparse convolution layers are followed by a ResNet-34 [26] network, with group normalization layers replacing the usual batch normalization layers, because it has been shown in [27] that they lead to a better generalization with smaller batch sizes, necessary for high resolution video training. Skip connections between encoder and decoder allow to transfer features from the previous frame to the new predicted frame in an implicit manner as shown in [19]. Contrary to our baseline, we use convLSTMs [28] instead of LSTM to leverage more granular spatio-temporal features. The latent variables z_t become tensors of size $3 \times 10 \times 512$, instead of $1 \times 1 \times 512$. Each 512-vector in z_t encodes a complex stochastic probability for a patch in the original 192×640 image.

The loss function (Eq. (2)) for the **Depth** Network is a combination of a masked L_2 loss and a KL-divergence loss to optimize the variational lower bound as in [16].

$$L_{\text{Depth}} = \sum_t \left(\sum_{\text{pixels}} \|\mathbb{1}_{lidar(t)>0} \cdot (dense(t) - lidar(t))\|_2^2 + \lambda_1 L_{KL}(q_\phi(z_t|x_{1:t}), q_\psi(z_t|x_{1:t-1})) \right) \quad (2)$$

Rather than utilizing the masked L_2 loss in the **Mask** network, instead a standard L_2 loss between the predicted mask and the groundtruth mask (which consists of all pixels that have a non-zero depth value) is employed. Although in recent works like SPCSF [24], the mask prediction is treated as a

TABLE I: Quantitative comparisons

Framework	(a) Next frame prediction		
	MAE	RMSE	CD
(a) SVG-LP	5.36	11.01	0.246
(b) Ours SLPC	0.74	2.371	0.071
Framework	(b) Prediction of 15 frames		
	MAE	RMSE	CD
(a) SPCSF	16.01	20.76	1.135
(b) SVG-LP	11.87	17.30	0.963
(c) Ours SLPC	2.68	5.38	0.700

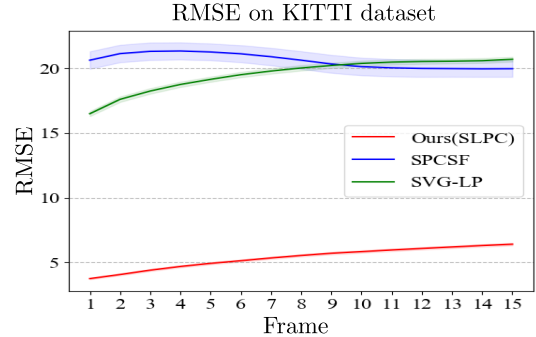


Fig. 2: Quantitative analysis of our depth maps accuracy in comparison to the baselines. RMSE is plotted on 15 frames with 95% confidence interval. All models are warmed up by 15 previous frames.

pixel classification problem, we have found that the L_2 loss is able to capture better details than the binary cross entropy loss. The latter leads to extremely sparse depth maps and a deterioration of the subsequent predicted frames.

To enable training on sequences consisting of 30 high resolution frames, we follow [20] by training with Truncated Back Propagation Through Time (TBPTT) [29]. By updating the weights after each frame, less computations are required in comparison to Back Propagation Through Time and we were able to experiment with high capacity encoders. We trained with a batch size of 8 to prevent overfitting. Our training strategy consists in training the **Depth** Network and the **Mask** network each in an isolated way to predict the next frame. Then we trained the two networks together in an autoregressive manner by teacher forcing to predict 15 frames after being warmed up by 15 previous frames.

III. DATASETS AND EXPERIMENTS

The proposed framework was evaluated on the KITTI dataset for depth completion [30]. It contains 85898 frames from 138 drives for training and around 7000 frames from 13 drives for validation. The LiDAR data was reprojected on a downscaled camera frame with a resolution of 192×640 using the given calibration data. We limited the videos to subsequences of 30 frames for training, and obtained 2554 training sequences in total. We trained on single NVIDIA TITAN RTX GPU.

Since there is no direct baseline, we compared with the SPCSF and the SVG-LP frameworks for Lidar Prediction. To

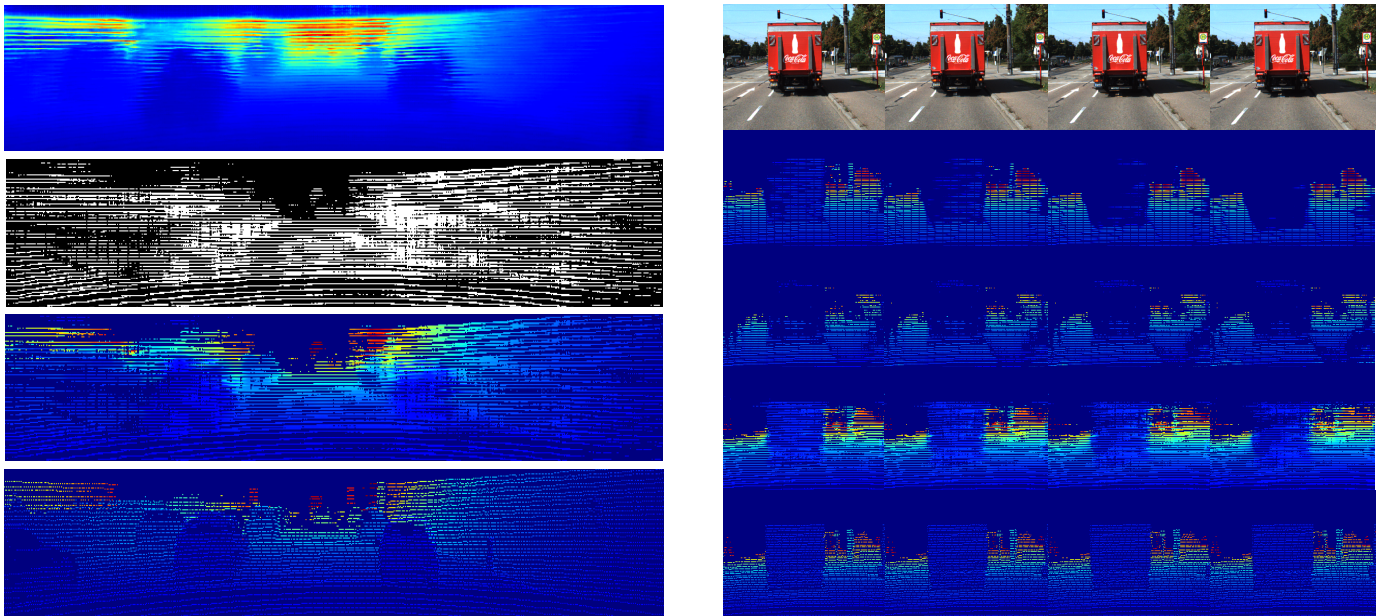


Fig. 3: Qualitative Comparison. From top to bottom: Left: the predicted dense map, predicted mask, predicted LiDAR frame, groundtruth. Right: RGB images, SVG-LP, SPCSF, Ours, Groundtruth. We illustrate 4 frames (17-20 in the sequence).

achieve a fair comparison, we used the verified open-source implementation of SVG-LP while reproducing a faithful implementation of the architecture of SPCSF and using the original hyperparameters. We have performed two experiments: first, we have trained our **Depth** and **Mask** networks on the KITTI dataset for the task of next frame prediction and compared with SVG-LP (since SPCSF is trained to produce multiple frames at once). Then we trained the two networks together to predict 15 frames after being conditioned by 15 frames (which means 1.5 seconds of driving on the KITTI dataset). Comparisons were made with both SVG-LP and SPCSF. Since our model is stochastic, we draw 20 samples from the model for each testing sequence as in [16] and report the mean RMSE for each predicted frame.

To evaluate the quantitative performance of our framework, the utilized metrics are the Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) in meters against the projected lidar frames, and the Chamfer Distance (CD) between the generated pointclouds and the original LiDAR pointclouds, projected in 3D using the camera calibration matrix in the KITTI dataset.

IV. RESULTS AND DISCUSSION

The quantitative results in Table I and Fig. 2 indicate that our model outperforms the baselines by a large margin in all three metrics. Interestingly, the SPCSF model, although optimized to LiDAR prediction on range maps, performed worse than the SVG-LP baseline. We hypothesize, this is due to the nature of the depth maps: they are sparser than range maps and lie on the camera frame of reference as opposed to the reference LiDAR frame. We also offer the advantage of generating dense depth maps as seen in Fig. 3, which results in better accuracy. Moreover, SPCSF is deterministic

as opposed to both SVG-LP and our model, which might lead to overfitting in some driving scenarios.

The qualitative analysis in Fig. 3 shows the robustness of our model against the sparse LiDAR. Both SVG-LP and SPCSF cannot simultaneously predict the depth and sparsity information accurately, and this resulted in depth maps with few points (or a lot of black holes in the depth estimation) especially in the latter frames. The mask network in our case prevents this degradation scenario from happening as it almost keeps the number of predicted points constant in each frame. In the case of SVG-LP, the resulting depth maps suffer from grid artifacts because the model was originally designed to deal with images not pointclouds. SPCSF and our Model do not suffer from this visual problem, because the loss functions of both models are more adapted to the 3D nature of LiDAR. However, the chamfer distance loss in SPCSF leads to pointclouds that are overpopulated in some locations and sparse in others, as was shown in [31].

Our model however has some limitations. Most importantly is the blurry nature of our depth maps. This is attributed to the use of VAE as noted in [17, 18]. We will work in the near future to address this issue: visual quality can be enhanced by either adversarial methods [17] or VAEs with hierarchical latent codes [18]. Simple photometric losses from the depth completion literature such as in [32, 33] can be used for our self-supervised framework. Another limitation lies in our mask predictions, which do not have the same level of sparsity as the original LiDAR frames and lead to a deterioration in the prediction quality over time, as it can be seen in Fig. 2. We plan in the near future to address these shortcomings.

V. CONCLUSION

In this work, we introduced the new task of depth map prediction and completion. To accomplish this, we presented

a new framework and training algorithm. Our model SLPC provides robust depth prediction against sparse LiDAR. This was achieved by two VRNNs working in parallel, one to complete the depth information and the other to predict a sparsity mask that downsamples the predicted depth map. The architecture and losses are tailored to the nature of the data. A comparative analysis with state-of-the-art frameworks has illustrated the superior performance of our model for LiDAR prediction. Our method is orthogonal to previous frameworks in video prediction, meaning that any video prediction model can be combined with our method to produce state-of-the-art predicted depth maps.

REFERENCES

- [1] Sherif Abdulatif et al. “Person Identification and Body Mass Index: A Deep Learning-Based Study on Micro-Dopplers”. In: *IEEE Radar Conference (RadarConf)*. 2019, pp. 1–6.
- [2] Sherif Abdulatif et al. “Stairs detection for enhancing wheelchair capabilities based on radar sensors”. In: *2017 IEEE 6th Global Conference on Consumer Electronics (GCCE)*. 2017, pp. 1–5.
- [3] William Lotter, Gabriel Kreiman, and David D. Cox. “Deep Predictive Coding Networks for Video Prediction and Unsupervised Learning”. In: *International Conference on Learning Representations, (ICLR)*. 2017.
- [4] Nal Kalchbrenner et al. “Video Pixel Networks”. In: *International Conference on Machine Learning, (ICML)*. Vol. 70. 2017, pp. 1771–1779.
- [5] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. “Unsupervised Learning of Video Representations using LSTMs”. In: *International Conference on Machine Learning, (ICML)*. Vol. 37. JMLR Workshop and Conference Proceedings. 2015, pp. 843–852.
- [6] Junhyuk Oh et al. “Action-Conditional Video Prediction using Deep Networks in Atari Games”. In: *Conference on Neural Information Processing Systems (NIPS)*. 2015.
- [7] Pulkit Agrawal, João Carreira, and Jitendra Malik. “Learning to See by Moving”. In: *IEEE International Conference on Computer Vision (ICCV)*. 2015, pp. 37–45. URL: <https://doi.org/10.1109/ICCV.2015.13>.
- [8] Chelsea Finn, Ian J. Goodfellow, and Sergey Levine. “Unsupervised Learning for Physical Interaction through Video Prediction”. In: *Conference on Neural Information Processing Systems (NIPS)*. 2016.
- [9] Fitsum A. Reda et al. “SDCNet: Video Prediction Using Spatially-Displaced Convolution”. In: *European Conference on Computer Vision (ECCV)* (2018).
- [10] Diederik P. Kingma and Max Welling. “Auto-Encoding Variational Bayes”. In: *International Conference on Learning Representations, (ICLR)*. 2014.
- [11] Danilo Jimenez Rezende et al. “Stochastic Backpropagation and Approximate Inference in Deep Generative Models”. In: *International Conference on Machine Learning, (ICML)*. Vol. 32. 2014, pp. 1278–1286.
- [12] Ian Goodfellow et al. “Generative adversarial nets”. In: *Advances in neural information processing systems*. 2014, pp. 2672–2680.
- [13] Karim Armanious et al. “An Adversarial Super Resolution Remedy for Radar Design Trade-offs”. In: *27th European Signal Processing Conference (EUSIPCO)*. 2019.
- [14] Sherif Abdulatif et al. “Towards Adversarial Denoising of Radar Micro-Doppler Signatures”. In: *International Radar Conference (RADAR)*. 2019, pp. 1–6.
- [15] Junyoung Chung et al. “A Recurrent Latent Variable Model for Sequential Data”. In: *Conference on Neural Information Processing Systems (NIPS)*. 2015.
- [16] Emily Denton and Rob Fergus. “Stochastic Video Generation with a Learned Prior”. In: *International Conference on Machine Learning, (ICML)*. 2018.
- [17] Alex X. Lee et al. “Stochastic Adversarial Video Prediction”. In: *arXiv preprint* (2018). URL: <http://arxiv.org/abs/1804.01523>.
- [18] Lluís Castrejón et al. “Improved Conditional VRNNs for Video Prediction”. In: *International Conference on Computer Vision (ICCV)*. 2019.
- [19] Jonas Uhrig et al. “Sparsity Invariant CNNs”. In: *International Conference on 3D Vision*. 2017, pp. 11–20.
- [20] Vaishakh Patil et al. “Don’t Forget The Past: Recurrent Depth Estimation from Monocular Video”. In: *IEEE Robotics and Automation Letters* 5 (2020), 6813–6820.
- [21] Marcel Schreiber et al. “Long-Term Occupancy Grid Prediction Using Recurrent Neural Networks”. In: *International Conference on Robotics and Automation (ICRA)*. 2019, pp. 9299–9305.
- [22] Masha Itkina, Katherine Rose Driggs-Campbell, and Mykel J. Kochenderfer. “Dynamic Environment Prediction in Urban Scenes using Recurrent Representation Learning”. In: *IEEE Intelligent Transportation Systems Conference (ITSC)*. 2019.
- [23] Yafei Song et al. “2D LiDAR Map Prediction via Estimating Motion Flow with GRU”. In: *International Conference on Robotics and Automation, (ICRA)*. 2019, pp. 6617–6623.
- [24] Xinshuo Weng et al. *Unsupervised Sequence Forecasting of 100,000 Points for Unsupervised Trajectory Forecasting*. 2020. arXiv: 2003.08376.
- [25] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. “A Point Set Generation Network for 3D Object Reconstruction from a Single Image”. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2017, pp. 2463–2471.
- [26] Kaiming He et al. “Deep Residual Learning for Image Recognition”. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 770–778.
- [27] Yuxin Wu and Kaiming He. “Group Normalization”. In: *International Journal of Computer Vision* 128.3 (2020).
- [28] Xingjian Shi et al. “Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting”. In: *Conference on Neural Information Processing Systems (NIPS)*. 2015, 802–810.
- [29] Ronald Williams and Jing Peng. “An Efficient Gradient-Based Algorithm for On-Line Training of Recurrent Network Trajectories”. In: *Neural Computation* (Sept. 1998).
- [30] Andreas Geiger, Philip Lenz, and Raquel Urtasun. “Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite”. In: *Conference on Computer Vision and Pattern Recognition (CVPR)*. 2012.
- [31] Panos Achlioptas et al. “Learning Representations and Generative Models for 3D Point Clouds”. In: *ICML*. 2018.
- [32] Clément Godard et al. “Digging Into Self-Supervised Monocular Depth Estimation”. In: *International Conference on Computer Vision (ICCV)*. 2019, pp. 3827–3837.
- [33] Fangchang Ma, Guilherme Cavalheiro, and Sertac Karaman. “Self-Supervised Sparse-to-Dense: Self-Supervised Depth Completion from LiDAR and Monocular Camera”. In: *International Conference on Robotics and Automation (ICRA)*. 2019, pp. 3288–3295.