

2D-RADAR IMAGING WITH DEEP CONVOLUTIONAL NEURAL NETWORKS

Mumin Jin^{1,2}, Atulya Yellepeddi², Gregory Wornell¹

¹Dept. EECS, MIT, Cambridge, MA 02139,

²Analog Devices, Inc. (Analog Garage) Boston, MA 02110

ABSTRACT

Millimeter-wave (mmWave) frequency-modulated continuous-waveform (FMCW) radar technology has become widely used for advanced driver assistance systems (ADAS) because of its ability to operate in harsh environmental conditions and provide direct measurements of range and velocity. However, the spatial resolution of an FMCW radar system is limited by the number of individual radar elements in it. While many algorithms have been developed to increase sensor array resolution for sparsely populated scenes with simplistic priors, many real-world scenes have neither the required level of sparsity nor easily described priors. In this work, we propose a system that uses deep convolutional neural networks (DCNN) to produce high-resolution radar images of realistic driving scenes. Our proposed system is able to generate radar point clouds that are five times as dense as traditional algorithms such as MUSIC and Orthogonal Matching Pursuit (OMP) from simulated radar data, enabling downstream tasks such as object detection and target classification by either a human or another neural network.

Index Terms— FMCW radar, deep learning, convolutional neural networks, ADAS

1. INTRODUCTION

In Advanced Driver Assistance Systems (ADAS), mmWave FMCW radar systems complement various other sensors on the vehicle by providing direct measurements of object ranges and velocities [1]. In recent years, due to their low cost and their ability to operate effectively in rain, snow, and extreme ambient lighting conditions, there is growing interest in using FMCW radar systems as an alternative to Lidars and provide increased redundancy [2]. However, since the azimuth resolution of a radar system is proportional to the number of radar sensors in an antenna array, implementing a large enough antenna array to achieve the required azimuth resolution of 1° for Level 4/Level 5 (L4/L5) autonomous driving can be very

expensive and technically challenging even with the help of multiple-input-multiple-output (MIMO) technologies [3, 4].

While many algorithms, such as multiple signal classification (MUSIC) and the various compressed sensing algorithms, exist to improve the resolution of sensor arrays, they have limited applications in automotive settings because they are designed for reconstructing signals with known sparsity or tractable priors such as Gaussian and Gaussian-Bernoulli distributions [5–15]. In automotive settings, however, it is difficult to find a basis in which a realistic driving scene is sparse, especially in high-density cluttered environments with guardrails and buildings [16]. Moreover, the prior distribution of the scenes we wish to image with automotive radar is difficult to describe mathematically.

In recent years, the proliferation of deep learning algorithms has demonstrated the representation power of neural networks, and there has been growing interest in training neural networks to aid in radar data processing [16–19]. To use neural networks for improving spatial resolution of radar arrays, several challenges must be overcome. First, radar data are complex-valued. Second, radar data have high dynamic range due to signal strength decreasing rapidly with range. Finally, in the range-azimuth domain, object shapes warp depending on their locations. As a result, convolution layers perform poorly on range-azimuth radar signals. The authors of [16, 19] use neural networks to process radar data in the range-Doppler dimensions without dealing with the challenge of object shapes warping in azimuth-range dimensions; therefore, they do not improve the spatial resolution of radar arrays.

In this work, we propose data processing steps and a deep convolutional neural network (DCNN) to produce high-quality two dimensional radar images of realistic driving scenes using a small uniform linear array (ULA) of FMCW radar sensors. Our proposed neural network works by applying its learned knowledge of the scene prior to refine the blurry input radar images into high-resolution occupancy grids, which can be treated as high quality radar point-clouds and serve as a useful intermediate representation of driving scenes for a variety of downstream tasks such as object detection, classification, segmentation, view generation, and other tasks that involve sensor fusion. To the best of our knowl-

This work was supported, in part, by Analog Devices, Inc., NSF under Grant No. CCF-1816209, and AFRL and the USAF AI Accelerator under Cooperative Agreement No. FA8750-19-2-1000.

edge, this work is the first to upsample radar point clouds using deep learning. In Section 2, we introduce the signal model and notations. We describe our proposed method in Section 3 and show proof-of-concept results using simulated data in Section 4.

2. BACKGROUND

Consider an FMCW system with one transmitter and a uniform linear array (ULA) of M receivers with spacing d imaging a scene containing K isotropic point reflectors¹, the received signal at m^{th} receiver, according to [1,3], can be modeled as

$$y_m(t) = \sum_{k=1}^K \alpha_k \exp \left\{ j \left(\omega_c \tau_{k,m} + \gamma \tau_{k,m} t - \frac{\gamma}{2} \tau_{k,m}^2 \right) \right\} \quad (1)$$

where α_k is the attenuation factor, which depends on the reflectivity and distance of the k^{th} reflector from the sensor. ω_c is the center frequency of the transmitted signal; γ is the sweep slope; and $\tau_{k,m}$ is the round-trip time delay from the transmitter to the k^{th} reflector and back to the m^{th} receiver. If the transmitted chirp lasts for T seconds, then the maximum frequency of $y_m(t)$ is $B = \gamma T$. The maximum detectable distance is $R_{\text{MAX}} = \frac{Tc}{2}$, where c is the speed of light. After sampling in time at the minimum Nyquist rate corresponding to B , we get

$$y_m[n] = \sum_{k=1}^K \alpha_k \exp \left\{ j \left(\omega_c \tau_{k,m} + \frac{2\pi \tau_{k,m}}{T} n - \frac{\gamma}{2} \tau_{k,m}^2 \right) \right\} \quad (2)$$

Then, for points far away enough from the receiver array, we have

$$\tau_{k,m} \approx \frac{2r_k + md \cos \theta_k}{c} \quad (3)$$

where d is the spacing between the receivers, r_k the radial distance of the reflector to the transmitter and the first receiver, and θ_k the angle of the reflector with respect to the receiver array as in Fig. 1.

Typically, the sensor spacing is chosen to be $d = \lambda/2$ to avoid aliasing in azimuth. By plugging in the results from (3) to (2) we get

$$\mathbf{Y}[m, n] \approx \sum_{k=1}^K \alpha_k \exp \left\{ j \left(\pi \cos \theta_k m + 2\pi \frac{r_k}{R_{\text{MAX}}} n + \psi_k \right) \right\} \quad (4)$$

where ψ_k is the phase shift corresponding to each reflector. Equation 4 suggests that we can retrieve the range and azimuth information, (r_k, θ_k) , by performing a Fast Fourier Transform (FFT) to each row then each column of \mathbf{Y} . With FFT processing alone, hundreds of radar elements or virtual antennas would be required to achieve azimuth resolution $< 1^\circ$ for L4/L5 autonomous driving, resulting in high cost in hardware implementation and energy usage.

¹Our single-transmitter ULA model captures the salient features of MIMO radar, in particular the resolution of the resulting virtual array

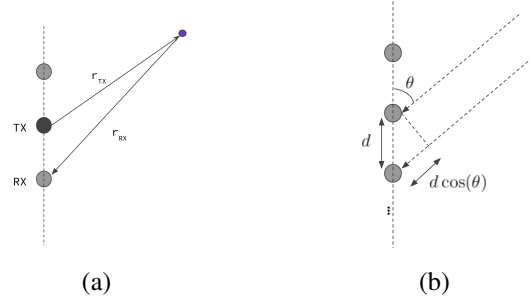


Fig. 1. Far-field approximation. When $r_k = r_{TX} \gg d$, the transmitted signal path and the received signal path shown in (a) are approximately parallel as in (b), $r_{RX} \approx r_k + d \cos(\theta)$.

3. PROPOSED METHOD

We propose a method that uses a neural network to produce high resolution radar point clouds using data from a small ULA during a single chirp. Our proposed system consists of signal pre-processing steps and a DCNN optimized using a supervised learning framework.

3.1. Pre-Processing

Given collected data, $\mathbf{Y} \in \mathbb{C}^{M \times N}$, from an array of M sensors over N time steps, we first use FFT processing to get a low-resolution proxy image of the scene. The result of FFT processing is a complex valued matrix of size $S \times D$, where $SD \gg MN$. One option to feed complex valued data to a neural network is to simply separate the real and imaginary parts into two channels. Since the attenuation α_k decays with $\frac{1}{r_k^s}$, with s depending on specularity of the object, we process the magnitude of the data in the log scale so that the neural network can more easily handle the high dynamic range of radar data when the scene contains objects very close to the sensors and objects farther away. We then normalize the log magnitude to be between 0 and 1. Finally, we create two channels of input by adding back the phase information as described in Algorithm 1. The result of pre-processing is a tensor of size $2 \times S \times D$ in the range $[-1, 1]$.

Algorithm 1 Data Pre-Processing

- 1: **Input:** $\mathbf{Y} \in \mathbb{C}^{M \times N}$
- 2: Range compression: compute D -point FFT for every row of \mathbf{Y} to get $\mathbf{Y}_r \in \mathbb{C}^{M \times D}$
- 3: Azimuth compression: compute S -point FFT for every column of \mathbf{Y}_r to get $\tilde{\mathbf{X}} \in \mathbb{C}^{S \times D}$, $\tilde{\mathbf{X}} = |\tilde{\mathbf{X}}| e^{j\angle \tilde{\mathbf{X}}}$
- 4: $\mathbf{L} \leftarrow \log(|\tilde{\mathbf{X}}|)$
- 5: $\mathbf{L} \leftarrow \frac{\mathbf{L} - \min \mathbf{L}}{\max \mathbf{L} - \min \mathbf{L}}$
- 6: **Output:** $[\mathbf{L} \odot \cos \angle \tilde{\mathbf{X}}, \mathbf{L} \odot \sin \angle \tilde{\mathbf{X}}]$, where \odot denotes element-wise product.

3.2. Neural Network

Inspired by the success of recent works which use convolution layers to process radar data for various tasks such as object detection and localization, we also adopt a network consisting of convolution layers [16, 18, 19]. In contrast to these works, which produce radar point clouds by detecting active range-Doppler cells without surpassing the FFT azimuth resolution limitation, we train our DCNN on a data from a single chirp without the Doppler dimension. Our DCNN increases the resolution of the bird's eye view image of the scene by applying its learned prior of the two-dimensional shapes in the scene to refine the low-resolution inputs.

Since the shape of an object warps based on its location in polar coordinates, convolution layers need to be able to model the location dependent distortion of object shape. This is accomplished in our system by concatenating the two channels from the proxy image with 32-dimensional hard-coded positional encodings proposed in [20] before feeding the result into the convolution layers. The details of the neural network architecture are shown in Fig. 2.

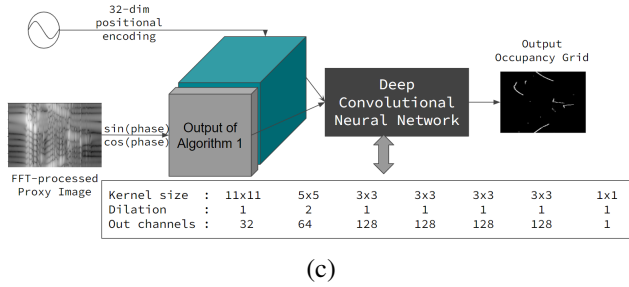


Fig. 2. Neural Network architecture. The activation function GELU is used for the intermediate layers, and Sigmoid is used for the final layer so that the NN output has values between 0 and 1. A wide kernel of 11×11 in the first layer is designed to capture object shapes that span both range and azimuth dimensions. We use circular padding in all layers because of the periodic nature of FFT and to keep the output dimensions the same as the input dimensions.

The DCNN was trained to minimize the “grid”-wise binary cross entropy loss between the output and the ground truth occupancy grid over a data set $\mathcal{D} = \{(\mathbf{Y}^{(i)}, \mathbf{G}^{(i)})\}$,

$$\mathcal{L} = -\frac{1}{S \cdot D} \mathbb{E}_{p_{\mathcal{D}}} \left[\left\| \mathbf{G} \log \hat{\mathbf{G}}(\mathbf{Y}) + (1 - \mathbf{G}) \log(1 - \hat{\mathbf{G}}(\mathbf{Y})) \right\|_1 \right] \quad (5)$$

where $\hat{\mathbf{G}}(\mathbf{Y})$ is the output of the neural network given inputs processed from \mathbf{Y} , and $\mathbf{G} \in \{0, 1\}^{S \times D}$ the ground truth grid indicating which azimuth-range cells contain reflectors. Similar to [21], which uses a neural network to recover active azimuthal sectors, the output of our neural network can be interpreted as a map of probabilities of the grid cells being active. By using two dimensional convolution kernels, our

proposed DCNN also exploits the global relations between grid-cells, thus learning a prior over two-dimensional scenes.

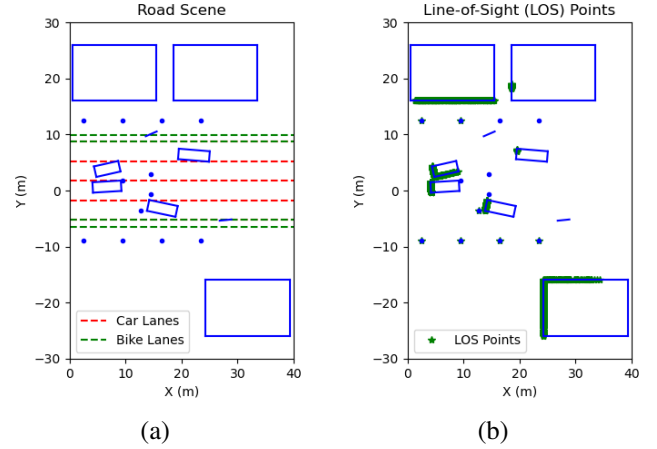


Fig. 3. Sample scene. (a) shows a scene with various shapes (solid blue lines and points) representing buildings, cars, and pedestrians that populate a four-lane road driving environment. The radar sensor is placed at the origin along the Y axis. (b) shows the final scene, which consists of point reflectors obtained by randomly tracing 1000 rays from the origin to locate the line-of-sight (LOS) points in the driving environment (right).

4. EXPERIMENTS AND RESULTS

4.1. Data Simulation

We trained our model on a simulated data set of driving scenes in a four-lane road driving environment. Fig. 3 shows an example of a simulated driving scene and the locations of point reflectors. For each scene, we simulate radar data using the signal model in equation (2) for a $M = 12$ -element ULA. We set the maximum detectable distance $R_{\text{MAX}} = 40\text{m}$ corresponding to a mid-range radar system, and we take $N = 128$ time samples over a single chirp duration $T = \frac{2R_{\text{max}}}{c}$. This radar system is small enough to be implemented on a single chip with time-division multiplex (TDM) MIMO [22].

We choose $S = D = 128$ for the dimension of the ground truth occupancy grid. The grid azimuth resolution of $\sim 1.4^\circ$ is close to meeting the L4/L5 autonomous driving standard. A grid cell has value 1 if a point reflector falls inside it. As seen in Fig. 3, a realistic driving scene, despite containing many reflecting points, is structured into a small union of geometric shapes, such as points, line segments, and corners of various sizes and orientations. Though the prior is not easily described mathematically, our neural network trained herein is able to represent it.

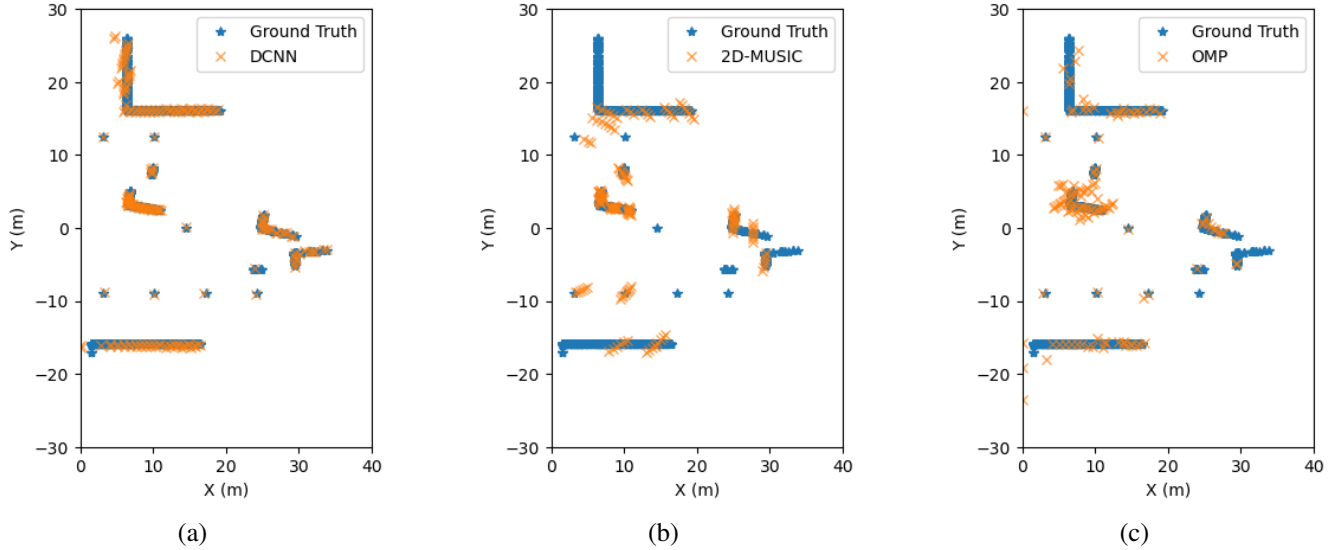


Fig. 4. Comparison of point clouds generated using (a) our proposed method, (b) 2D-MUSIC, and (c) OMP on a test sample.

4.2. Training

Using the simulation method described in 4.1, we generated a total of 10000 scenes. 8540 of these scenes were used for training, and 460 were used for validation. The remaining 1000 scenes were set aside for testing and evaluation. The neural network was trained with batch-size 16, using the ADAM optimizer with learning rate = 0.001, $(\beta_1, \beta_2) = (0.9, 0.99)$, decreasing the learning rate by a factor of 0.95 every 7 epochs. The final model check point was selected from when the validation loss was minimized after training for 20 epochs.

4.3. Results

Many of the recent deep-learning solutions to closely related problems in the radar domain inspired our proposed system architecture [16, 18, 19]. However, since these solutions use different kinds of inputs and processing from our work and are designed for different purposes other than generating high-quality spatial representations of driving scenes from low-resolution radar data, adapting the recent deep-learning models for our problem would lead to significant architecture changes. Therefore, we compare the outputs of our proposed method to the results of OMP and 2D-MUSIC on the 128×128 azimuth-range grid [7, 9].

Table 1 shows the average grid-wise probabilities of detection (p_D), false alarm (p_{FA}), precision and F1-Scores over the test data set. p_D , which is equivalent to the recall value, in this context measures the density of predicted points around the objects. Table 1 shows that, on average, the point clouds produced by our proposed DCNN are five times as dense as those by 2D-MUSIC and OMP given similar low-levels of

Table 1. Average Grid-wise p_D (Recall), p_{FA} , Precision, and F1-Scores Evaluated on Test Data Set

Method	p_D	p_{FA}	Precision	F1-Score
DCNN	79.52%	0.18%	84.07%	81.65%
OMP	16.47%	0.60%	23.33%	18.87%
2d MUSIC	16.40%	0.25%	44.07%	22.61%

p_{FA} . When OMP and 2D-MUSIC were allowed to predict equal number of data points as DCNN, their false positive rates increased faster than detection rates.

Fig. 4 shows a comparison of the results of using our proposed system, 2D-MUSIC, and OMP on a sample scene from the test data set. The locations of active range-azimuth cells extracted from each algorithm are overlaid on top of the ground-truth locations of the point reflectors. For a scene generated from the same process as the training data, the point clouds generated with our proposed system has greater overlap with the ground truth points than using 2D-MUSIC or OMP. Fig. 5 compares the results on a sample that is outside of the training data distribution but uses the same set of shapes. In this case, our proposed method is still able to generalize to scenes containing the same set of shapes, while OMP produces many false alarms around the points that are closer to the origin and miss the points that are farther away. Finally, Fig. 6 shows an example scene where super-resolution in azimuth is required because all the “cars” have the same distance to the origin. Our proposed system is able to resolve the separate shapes while 2D-MUSIC and OMP fail to recover the distinct shapes occupying the same range bins. Qualitatively, even when OMP and 2D-MUSIC are able to predict points clustered around an object, the point clusters

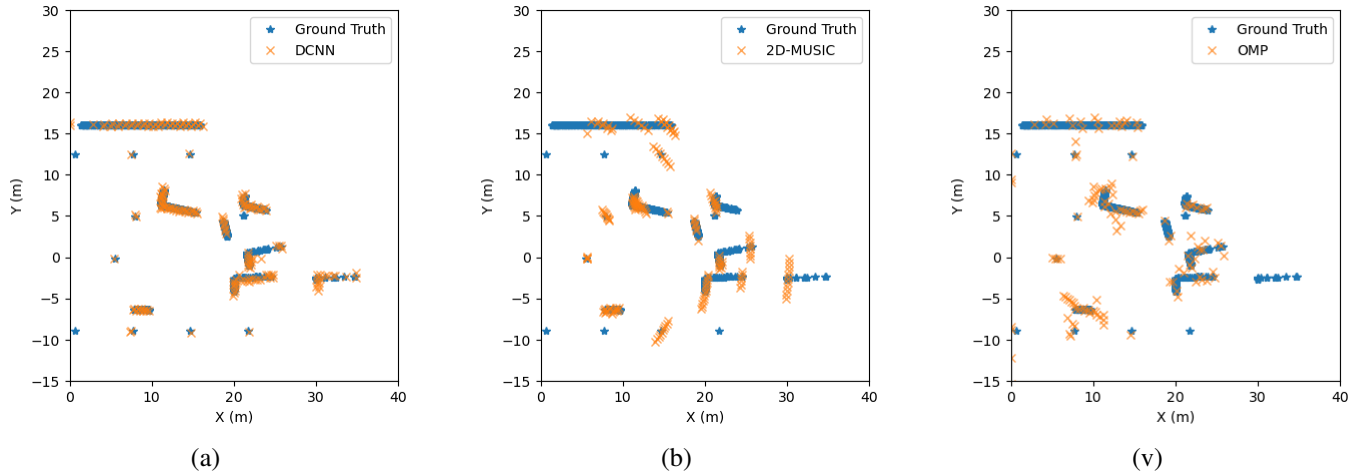


Fig. 5. Comparison of point clouds generated using (a) our proposed method, (b) 2D-MUSIC, and (c) OMP on an out-of-distribution sample. In this scene, the top and bottom lane contains two “cars”, while each lane only contains a single “car” in scenes in the training data set. Our DCNN generalizes to scenes with the same set of shapes.

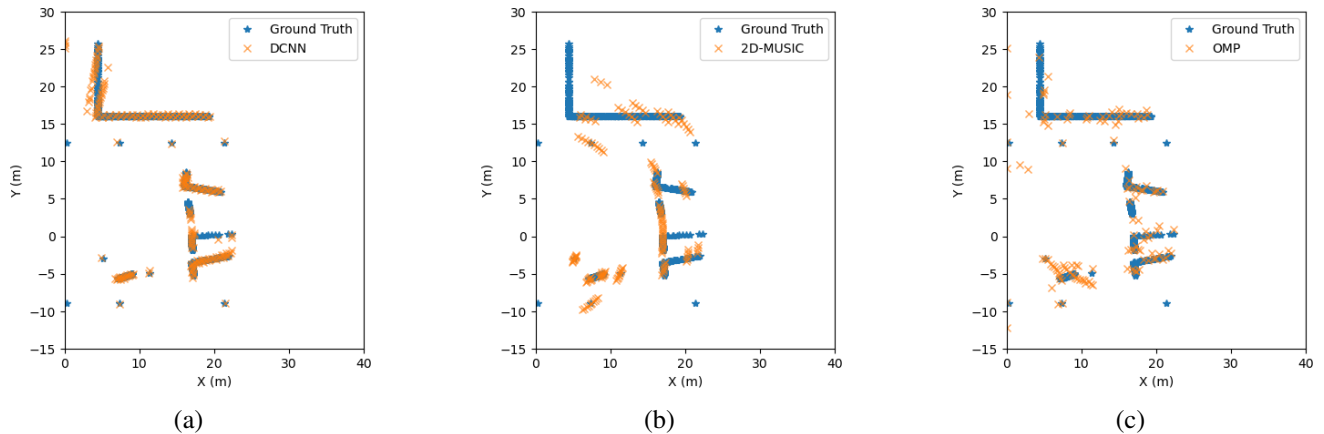


Fig. 6. Comparison of point clouds generated using (a) our proposed method, (b) 2D-MUSIC, and (c) OMP on a sample scene that requires super-resolution in azimuth.

do not necessarily form distinctive shapes. In comparison, the DCNN densely generates points that closely conform to the outlines of the objects in the scene, making a suitable representation of the scene for various downstream applications such as object detection and classification.

Our proposed system does have two main limitations. First, it requires supervised training using ground-truth occupancy grids. In practice, to apply our system to real-world data, we can simultaneously collect radar and Lidar data and use Lidar point clouds to produce the ground truth labels. Second, our system requires specialized training to tailor to each sensor configuration and driving environment. However, once our the model is trained, our system works much faster during inference than 2D-MUSIC or OMP.

5. CONCLUSION

In this paper, we proposed a deep learning system to produce high-resolution, dense radar point clouds using data collected from a small ULA. Our method produces point clouds that are five times as dense as traditional methods such as OMP and 2D-MUSIC and achieves F1-Scores four times as high as OMP and 2D-MUSIC. We focused on two dimensional imaging by using data from a single chirp, but extension to 3D, 4D imaging is straightforward. The results also promise potential applications of our methods beyond automotive radar imaging. In future work, we plan to apply our proposed system to real-world data and evaluate the performance of various downstream tasks using the outputs of our method as intermediate representations of driving scenes.

6. REFERENCES

- [1] S. M. Patole, M. Torlak, D. Wang, and M. Ali, "Automotive Radars: A Review of Signal Processing techniques," *IEEE Signal Processing Magazine*, vol. 34, no. 2, pp. 22–35, 2017.
- [2] T. Lim, A. Ansari, B. Major, D. Fontijne, M. Hamilton, R. Gowaikar, and S. Subramanian, "Radar and Camera Early Fusion for Vehicle Detection in Advanced Driver Assistance Systems," in *Machine Learning for Autonomous Driving Workshop at the 33rd Conf. on Neural Information Processing Systems*, vol. 2, 2019, p. 7.
- [3] S. Sun, A. P. Petropulu, and H. V. Poor, "MIMO Radar for Advanced Driver-Assistance Systems and Autonomous Driving: Advantages and Challenges," *IEEE Signal Processing Magazine*, vol. 37, no. 4, pp. 98–117, 2020.
- [4] M. Schoor, G. Kuehnle, K. Rambach, and B. Loesch, "Method for Operating A MIMO Radar," U.S. Patent 2014/0347211, Nov. 27, 2014.
- [5] R. Schmidt, "Multiple Emitter Location and Signal Parameter Estimation," *IEEE Trans. on Antennas and Propagation*, vol. 34, no. 3, pp. 276–280, 1986.
- [6] J. Odendaal, E. Barnard, and C. W. I. Pistorius, "Two-Dimensional Superresolution Radar Imaging Using the MUSIC algorithm," *IEEE Trans. on Antennas and Propagation*, vol. 42, no. 10, pp. 1386–1391, 1994.
- [7] F. Belfiori, W. van Rossum, and P. Hoogeboom, "Application of 2D MUSIC algorithm to Range-Azimuth FMCW Radar Data," in *2012 9th European Radar Conf.* IEEE, 2012, pp. 242–245.
- [8] S. G. Mallat and Z. Zhang, "Matching Pursuits with Time-Frequency Dictionaries," *IEEE Trans. on Signal Processing*, vol. 41, no. 12, pp. 3397–3415, 1993.
- [9] T. Cai and L. Wang, "Orthogonal Matching Pursuit for Sparse Signal Recovery with Noise," *IEEE Trans. on Information theory*, vol. 57, no. 7, pp. 4680–4688, 2011.
- [10] B. Mamandipoor, D. Ramasamy, and U. Madhow, "Newtonized Orthogonal Matching Pursuit: Frequency Estimation Over the Continuum," *IEEE Trans. on Signal Processing*, vol. 64, no. 19, pp. 5066–5081, 2016.
- [11] D. L. Donoho, A. Maleki, and A. Montanari, "Message-Passing Algorithms for Compressed Sensing," *Proceedings of the National Academy of Sciences*, vol. 106, no. 45, pp. 18 914–18 919, 2009.
- [12] A. Maleki, L. Anitori, Z. Yang, and R. G. Baraniuk, "Asymptotic Analysis of Complex LASSO via Complex Approximate Message Passing (CAMP)," *IEEE Trans. on Information Theory*, vol. 59, no. 7, pp. 4290–4308, 2013.
- [13] Z. Yang, X. Chen, and X. Wu, "A Robust and Statistically Efficient Maximum-Likelihood Method for DOA Estimation Using Sparse Linear Arrays," *arXiv preprint arXiv:2203.13433*, 2022.
- [14] F. Biondi, "Compressed Sensing Radar-New Concepts of Incoherent Continuous Wave Transmissions," in *2015 3rd Int. Workshop on Compressed Sensing Theory and Its Applications to Radar, Sonar and Remote Sensing (CoSeRa)*. IEEE, 2015, pp. 204–208.
- [15] W. Rowe, J. Li, and P. Stoica, "Sparse Iterative Adaptive Approach with Application to Source Localization," in *2013 5th IEEE Int. Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE, 2013, pp. 196–199.
- [16] Y. Cheng, J. Su, M. Jiang, and Y. Liu, "A Novel Radar Point Cloud Generation Method for Robot Environment Perception," *IEEE Trans on Robotics*, vol. 38, no. 6, pp. 3754–3773, 2022.
- [17] Y. Su, K. Fan, N. Bach, C. Kuo, and F. Huang, "Unsupervised Multi-Modal Neural Machine Translation," in *Proceedings of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2019, pp. 10 482–10 491.
- [18] N. Kumchaiseemak, I. Chatnuntaweck, S. Teerapitayanon, P. Kotchapansompote, T. Kaewlee, M. Piriya-jitakonkij, T. Wilaiprasitporn, and S. Suwajanakorn, "Toward Ant-Sized Moving Object Localization Using Deep Learning in FMCW Radar: A Pilot Study," *IEEE Trans. on Geoscience and Remote Sensing*, vol. 60, pp. 1–10, 2022.
- [19] D. Brodeski, I. Bilik, and R. Giryes, "Deep Radar Detector," in *2019 IEEE Radar Conf. (RadarConf)*. IEEE, 2019, pp. 1–6.
- [20] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, A. G. L. Jones, L. Kaiser, and I. Polosukhin, "Attention Is All You Need," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [21] M. Gall, M. Gardill, T. Horn, and J. Fuchs, "Spectrum-Based Single-Snapshot Super-Resolution Direction-of-Arrival Estimation Using Deep Learning," in *2020 German Microwave Conf. (GeMiC)*. IEEE, 2020, pp. 184–187.
- [22] "AWR1843AOP Single-chip 77- and 79-GHz FMCW mmWave Sensor Antennas-OnPackage (AOP)," AWR1843BOOST datasheet.