

Improving Millimeter Wave Radar Perception with Deep Learning

Junfeng Guan

Sohrab Madani

ECE 544 Project Report I

jguan8@illinois.edu

smadani2@illinois.edu

1. Introduction and Project Descriptions

Since the past few years, AI-Powered autonomy revolution in the automotive industry has attracted great attention worldwide. It is believed that in the not-too-distant future, fully autonomous vehicles will be the norm rather than the exception, redefining mobility in our daily lives. With deep learning widely applied on sensor data, self-driving cars are able to localize and map objects, understand the environment, and make correct decisions. As the most fundamental task, previous works have demonstrated accurate object detection and classification, but they are limited to data obtained from LiDARs and cameras. These optical sensors have high imaging resolution, but they naturally fail in low visibility conditions such as fog, rain, and snow, because light beams are narrower than water droplets and snowflakes. [4] This fundamental limitation of optical sensors is one of the major roadblocks to achieving the 5th SAE level of full automation. [1] On the contrary, Radar wave propagates through small particles and can provide an alternate imaging solution in such inclement weather. Besides, radar can also directly measure the velocity of vehicles based on the doppler shift of the reflected signal instead of going through cluster tracking among frames like LiDAR. With the additional realtime speed information, AI can potentially have a better perception than human drivers.

Although the low resolution of traditional automotive radar overshadows its advantages, the advent of Millimeter-wave (mmWave) technology makes it possible to have a reliable 3D imaging system in inclement weather with relatively higher resolution. Along with good propagation characteristics, mmWave also provides much wider bandwidth and enables miniature-sized large-aperture antenna arrays, which improve distance and angle measurement respectively. Previous work has demonstrates sub-centimeter level imaging resolution for short range objects [5], but the fundamental difference in wavelength between radar (5 mm) and LiDAR (905 nm) cannot be easily overcome. Moreover, radar images are not as readable to non-experts. Therefore, in order to convince the automotive industry and general public that mmWave is a reliable imaging solution for autonomous driving, we should try to generate radar images

targeting the resolution comparable to LiDAR, and make them more perceptually intuitive to people.

In this project, we propose to develop high-resolution and reliable imaging techniques for self-driving cars with mmWave radar. Specifically training neural networks to enhance the low-resolution and unreadable radar images to be similar to the LiDAR point cloud, which has been extensively and successfully used for self-driving perception. Eventually enable various crucial vision applications for autonomous vehicles like lane detection, image mapping, localization, and object identification with mmWave radar data only. To do this, we have to overcome challenges analyzed in the following section 2.

2. Challenges

2.1. Radar Imaging Primer

Radar generates images of objects by localizing the cluster of point reflectors that forms the object. 3D space imaging requires mapping point reflectors to voxels in a spherical coordinates with its distance, azimuth angle, and elevation angle. Distance is measured through round-trip Time-of-Flight (ToF) of reflected radar signal, while azimuth and elevation angles can be either obtained by the beam steering angle of phased array antennas or be estimated with Direction-of-Arrival (DoA) estimation algorithms such as beam-forming or Multiple Signal Classification (MUSIC). Notice that as distance becomes further, the voxel size corresponds to the same angle increases, so that long-range objects contain much fewer voxels and appear to be more blurry. Besides, unlike the extremely narrow width of light beams, the cone-shape radar beam with sidelobes cause interference and leakage among nearby voxels and even the environment, which smears generated images. Last but not least, with a few centimeter resolution of distance measurement, a continuum of a large number of point reflector sums up and makes it an under-determined problem to localize them. Besides, radar reflection tends to be more specular than the mostly scattering reflection at optical frequencies. In other words, reflection off smooth objects might mainly towards an angle away from the radar receiver and disap-

pears in the image. Hence, "edge detection" in radar images is very different from that of pictures. Firstly, it needs to predict high frequency information with only low frequency data. Secondly, it needs to learn to fill up missing parts of the object.

2.2. Related Work

iiiiii HEAD Previous works have attempted to adapt the optical camera-oriented Convolutional Neural Network (CNN) to its microwave counterpart to classify single-object images from high-resolution 2D synthetic aperture radar (SAR) images. [2] [7] [3] There has also been successful application of CNNs on recreating short range human body skeletons in 2D images and 3D spaces from radar images by tracking 14 key points. [8] [9] However, the scope of these application of neural networks on radar images are restricted to feature and single object classification. Besides, their radar image are either of super high resolution generated with airborne geographical sensing synthetic aperture radar or short-range where voxel resolution does degrade too much. On the contrary, in our application of autonomous cars, we not only don't have as high resolution for long range, but required precisely traced boundary which contains the information of size, shape, and orientation of cars, bikes, and pedestrians. Therefore if we rely on fine grained feature detection that consists boundaries, we need very deep neural nets. Since CNN structures haven't been well investigated for radar signal and image, this task becomes extremely difficult.

===== Previous works have attempted to adapt the optical camera-oriented Convolutional Neural Network (CNN) to its microwave counterpart to classify single-object images from high-resolution 2D synthetic aperture radar (SAR) images. [2] [7] [3] There has also been successful application of CNNs on recreating short range human body skeletons in 2D images and 3D spaces from radar images by tracking 14 key points. [8] [9] However, the scope of these application of neural networks on radar images are restricted to feature and single object classification. Besides, their radar image are either of super high resolution generated with airborne geographical sensing synthetic aperture radar or short-range where voxel resolution does not degrade too much. On the contrary, in our application of autonomous cars, we not only do not have as high resolution for long range, but required precisely traced boundary which contains the information of size, shape, and orientation of cars, bikes, and pedestrians. Therefore if we rely on fine grained feature detection that consists boundaries, we need very deep neural nets. Since CNN structures haven't be well investigated for radar signal and image, this task becomes extremely difficult. 9d9e575ff3d6c11b45fcb77359b58858254bb413

2.3. Dataset Availability

Once the network is setup, it needs training data -i.e., it needs many labeled exp Dataset Variation between systems Experiment Processing

2.4. 3D Complexity

3.3D 3D CNN 3D GAN size complexity

2.5. Evaluation Metrics

3. Method

Describe the overall method on how you solve the proposed problem, and a bit of original derivation that has some relevance to what youre trying to accomplish

Problem statement: low blurred images no boundary can be seen, specularity causes missing parts, no availbale dataset, 4D CNN NN complex.

Overview 3D, in the method, we say that we start with 2D version, and based on that build 3D. We are trying to use cGAN to generate higher resolution images from low resolution radar images, which should have a sharp and accurate boundary of the object. Also, the missing parts due to specularity of reflection need to be feel up.

Generative Adversarial Networks (or GANs) have been widely used to generate images [cite some stuff here] and fill in missing parts of data. In our case, we are looking to generate images with accurate boundaries using low resolution images with missing parts as input. Conditional GANs have already proven successful in similar settings, such as in [6], where the authors have used thermal images under low light conditions where there some parts of the image are missing to retrieve the human face boundaries and estimate its orientation. Another motivation behind using conditional GANs is that loss functions such as the L_2 and L_1 (i.e. loss functions that are equal the Euclidean or L_1 distance between the input and the output) which are the de facto standard loss function for restoring images render blurry images, which are not suitable for our application as they do not emphasize on boundaries. On the other hand, using conditional GAN, we were be able to motivate the loss function in GAN to learn to focus on the boundaries, by designing the ground to contain information mostly about the boundary of objects, as discussed in more detail the dataset section.

We have the CGAN from [cite pix2pix] to train and test our model. Denoting the input, ground truth, and noise by x , y , and z respectively, we can write the objective for a conditional GAN where the Generator and discriminator are G and D as

$$\min_G \max_D \mathcal{V}_{CGAN}(D, G) = \mathbb{E}_x(\log(D(x|y))) + \mathbb{E}_z(\log(1 - D(x, G(x|z)))) \quad (1)$$

Here, $D()$ is the score that the discriminator gives for a certain input, and $G(x, z)$ is what the generator tries to generate given the input x and noise vector z . The point of having the noise vector, of course, is to avoid the problem of over-fitting. The difference between the standard GAN and the conditional version is that in the latter everything is conditioned to y . This y could be anything we want, that adds extra information about what the output G should look like. In our setting, we call y the ground truth, as it contains information of real boundaries of the objects. Given y , D tries to maximize its output value when the input is real, and at the same time minimize it when the input is artificially generated by the G , the generator. At the same time, the generator tries to generate an image that gets a high score from D , motivating it to generate images that are similar to real ones.

It is suggested in [pix2pix] that we can combine generator's task of trying to get a high score from discriminator with a pre-determined loss function, such as L_1 . That is, one could change the objective to be

$$\begin{aligned} \min_G \max_D \mathcal{V}_{CGAN}(D, G) = \\ \mathbb{E}_x(\log(D(x|y))) + \mathbb{E}_z(\log(1 - D(x, G(x|z)))) \\ + \lambda \mathcal{L}_{L_1}(G(x, z)). \end{aligned} \quad (2)$$

The motivation behind this is to capture the low-frequency information using the L_1 loss, and motivate GAN to model high-frequency information, which in our case, will translate to more precision in identifying boundaries.

For the generator, we have adopted the architecture of [unet], which is an auto-encoder. Similar to most such architectures, U-net consists of a contracting path and an expansive path. The first layer contracting path is made of two successive 3×3 convolutions followed by a ReLu (rectified linear unit) and a 2×2 max-pooling layer without overlapping, which shrinks the size of its input by a factor of four. Let us call such a layer a down-sampling layer. This layer is then repeated multiple times, each time doubling the feature channels. The expansive path reverses this process: at each layer, the data is first up-sampled by a factor of two, then filtered with using a 2×2 window (2 by 2 convolution) halving the number of feature channels. At this point, a cropped version feature maps from the corresponding layer in the contracting path is forwarded to this layer (i.e. for layer k , the feature maps from layer $n - k$ are forwarded, assuming a total number of n layers) and are concatenated with the current feature maps. There are two points to be mentioned here; first, the cropping should happen because of the loss of edges that has happened due to successive convolutions. Second, and more importantly, this idea is used as a way to detour the bottleneck layer (the middle layer through which all the information should pass) by creating these skip connections between layers. One reason behind this choice of architecture is that this forwarding

encourages the similar structure between the input and the output. In other use cases, [U-net] evaluated this architecture in a biomedical image segmentation problem, and in [pix2pix], it was adopted to generate a general translator of high-resolution output images to high resolution input images, which is quite similar to our setting, except that in our problem the input images are not high-resolution.

As for the discriminator, what we need is for the network to be able to identify local structures, in order to capture the properties of local objects (e.g. cars). Since the network is relying on the L_1 loss to guarantee the correctness of low-frequency information, it is possible to restrict the discriminator to only penalize structures that occur within a patch window of the image. In other words, the discriminator slides over the image looking at patches of size $n \times n$, and scoring each of them, and finally averaging over all patches to derive the final score. For our implementation, we chose n to be 70 where the image size was 256×256 . This structure has been dubbed patchGAN [pix2pix].

Input: Low resolution radar images, because there is no available public dataset of mmWave radar images, and collecting a big enough dataset by ourselves is not possible. We synthesized radar images. This heat-map is fed into the network. Groundtruth: 1. Size of 3D which makes the training phase very slow even when using GPUs. 2. 2D: 3. 3D

The input to our problem is pre-processed radar data. First, using the raw data from the antenna array, a coarse heat-map of objects are generated. After some further processing.

Why not raw data: (Goes to detail of Dataset jayden) Radar imaging processing algorithms is a well established field for 70 years. and there is fruitless to try and learn them using machine learning. So instead of The reason why we did not choose the raw data as input is twofold.

3.1. Conditional GAN

3.2. Dataset Generation

We were not able to find any public dataset of 3D or even 2D radar images especially street views for the self-driving car application. Also unlike cameras and LiDARs that generate standard images, different radar imaging system and image processing algorithms can lead to non-negligible variation in radar images. Hence we were initially forced to build our own dataset of 3D radar images with our custom-built 60GHz imaging system. However, since our system design compensates frame rate for cost and resolution, continuously collecting 500 frames will take at least 41.67 hours. We also have to wisely pick a good number of locations to mimic various scenes autonomous cars would actually face. Last but not in addition to the radar images as the input to cGAN, we also have to obtain another perfectly synced camera or LiDAR dataset as the ground truth of object boundary.

at 500 well selected locations , which becomes a huge overhead. Moreover,

Once we have designed our Unlike the digital camera, imaging radar less common, no dataset available . Since there is no

Challenge of unavailable radar dataset.

Experiment to collect radar images size time umbiguity

Simulation with EM

3D

2D Mask R-CNN input: radar image groundtruth: mask
pros: large dataset with car truck human cons: No 3D info,
no specularity 3D CAD input: contour groundtruth: pros:
small dataset, single element cons: 3D shape info, specular-
ity

Evaluate the simulation with EM simulation Feko and
experimental results.

4. Experimental Results

Describe the setup of the experiments you ran, e.g., what
evaluation metrics, datasets are used. Present the results,
preferably in the form of tables and/or figures

4.1. Dataset

Training dataset 118 Testing dataset 509

4.2. Results

5. Discussion and Conclusion

Analyze the results, summarize the findings and point
out possible future directions

References

- [1] Automated vehicles for safety. 1
- [2] S. Chen and H. Wang. Sar target recognition based on deep learning. *International Conference on Data Science and Advanced Analytics, DSAA*, 2014. 2
- [3] M. Gong, J. Zhao, J. Liu, Q. Miao, and L. Jiao. Change detection in synthetic aperture radar images based on deep neural networks. *IEEE Transactions on Neural Networks and Learning Systems* , *DSAA*, 2016. 2
- [4] W. Jones. Keeping cars from crashing. *IEEE Spectrum*, 2001. 1
- [5] B. Mamandipoor, G. Malysa, A. Arbabian, U. Madhow, and K. Noujeim. 60 ghz synthetic aperture radar for short-range imaging: Theory and experiments. *Asilomar Conference on Signals, Systems and Computers*, 2014. 1
- [6] A. Nambi, S. Bannur, I. Mehta, H. Kalra, A. Virmani, V. Padmandabhan, R. Bhandari, and B. Raman. Demo: Hams: Driver and driving monitoring using a smartphone. *ACM MobiCom*, 2018. 2
- [7] C. Schwegmann, W. Kleynhans, B. Salmon, L. Mdakane, and R. Meyer. Very deep learning for ship discrimination in synthetic aperture radar imagery. *IEEE International Geoscience and Remote Sensing Symposium, IGARSS*, 2016. 2

[8] M. Zhao, T. Li, M. Alsheikh, Y. Tian, H. Zhao, A. Torralba, and D. Katabi. Throug-wall human pose estimation using radio signals. *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2018. 2

[9] M. Zhao, Y. Tian, H. Zhao, M. Alsheikh, T. Li, R. Hristov, Z. Kabelac, D. Katabi, and A. Torralba. Rf-based 3d skeletons. *SIGCOMM*, 2018. 2