

# Radar-Conditioned 3D Bounding Box Diffusion for Indoor Human Perception

Ryoma Yataka<sup>1,2\*</sup>, Pu (Perry) Wang<sup>2</sup>, Petros Boufounos<sup>2</sup>, and Ryuhei Takahashi<sup>1</sup>

<sup>1</sup>Information Technology R&D Center, Mitsubishi Electric Corporation, Kanagawa 247-8501, Japan

<sup>2</sup>Mitsubishi Electric Research Laboratories (MERL), Cambridge, MA 02139, USA

## Abstract

Privacy-preserving and cost-effective indoor sensing is vital for embodied agents to collaborate safely with people in dynamic scenes. Multi-view millimeter-wave radar shows great potential for this purpose. However, prevailing methods rely on implicit cross-view association, which this reliance often results in ambiguous feature matches and degraded performance in cluttered environments. To address these limitations, we propose **REXO** (multi-view Radar object dEtection with 3D bounding boX diffusiOn), which lifts DiffusionDet’s 2D box denoising to the full 3D radar space. Noisy 3D boxes are projected onto all radar views to enable **explicit** association and radar-conditioned denoising. Evaluated on two open indoor radar datasets, our approach outperforms state-of-the-art methods by +11.02 AP on MMVR and +4.22 AP on HIBER.

## 1. Introduction

Reliable perception is crucial for embodied agents in indoor settings (homes, factories, clinics), where scene understanding, motion capture, and human-robot collaboration are required. Radar is increasingly used for navigation, manipulation, and safer human-robot interaction because it provides robust awareness in low light, smoke, dust, and even through cardboard and plastic [23, 31]. For example, FuseGrasp [9] fuses radar and camera to grasp transparent objects, exploiting millimeter-wave (mmWave) radar’s ability to render transparent materials opaque and robotic-arm radar imaging to recover shapes invisible to RGB-D. On the other hand, radar-only perception (see Appendix A) remains challenging. Multi-view radar methods either pair horizontal proposals with fixed-height vertical ones [37] (Fig. 1 (a)) or use query-based transformers to regress 3D bounding boxes (BBoxes) from both views [40] (Fig. 1 (b)). Image-based object detection has been redefined as a generative denoising process, where a random noisy 2D BBox is iteratively refined through a diffusion denoising process

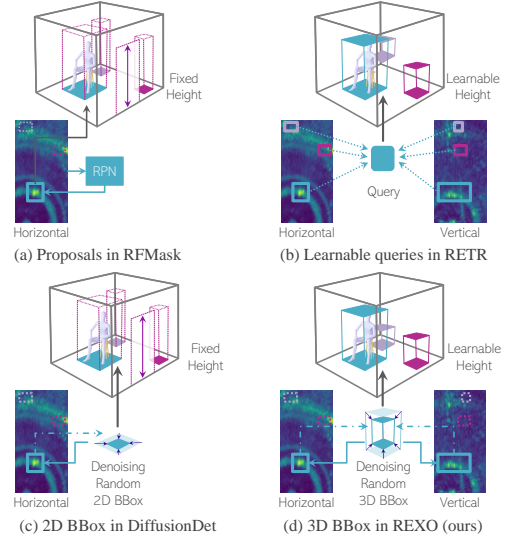


Figure 1. (a) RFMask [37] generates horizontal-view proposals with fixed height; (b) RETR [40] implicitly links queries to cross-view features; (c) DiffusionDet [7] needs pairing with fixed-height vertical BBoxes; (d) REXO (ours) performs diffusion directly in 3D radar space for simple, explicit cross-view association.

to yield a final clean BBox [7] and this approach generally surpasses query-based detectors. When ported to horizontal radar heatmaps (Fig. 1 (c)), it denoises 2D BBoxes but still requires the fixed-height vertical pairing used by RFMask.

We therefore *lift* the diffusion procedure from a 2D plane (image or horizontal radar view) in DiffusionDet to the full 3D radar space, as illustrated in Fig. 1 (d). This simple lifting facilitates cross-view radar feature association and radar-conditioned BBox denoising, while enabling the integration of geometry-aware loss functions and prior constraints on the 3D BBox. Consequently, we introduce the proposed framework as **Radar object dEtection with 3D bounding boX diffusiOn (REXO)** with the following contributions:

1. **2D-to-3D Lifting with Explicit Cross-View Association:** At each diffusion timestep, a noisy 3D BBox is projected onto every radar view, and RoI-aligned crops supply view-specific features. This BBox-guided association grows *linearly* with the number of views, whereas proposal- or query-based schemes grow quadratically.

\*The work was done at MERL as a visiting scientist.

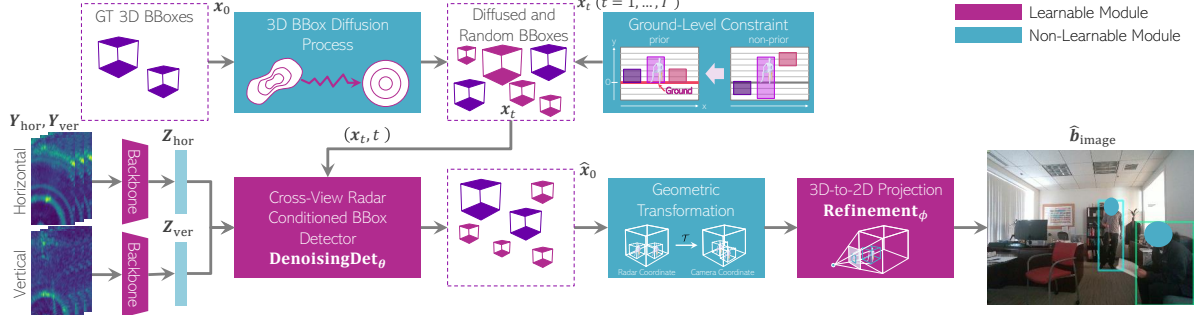


Figure 2. **REXO training:** 1) A backbone extracts horizontal/vertical radar features; 2) Ground-truth 3D BBox  $x_0$  are diffused to noisy  $x_t$ ; 3)  $x_t$  is grounded using a ground-level constraint; 4)  $\text{DenoisingDet}_\theta$  projects  $x_t$  onto both views and uses the aligned features to recover  $\hat{x}_0$ ; 5) A radar-to-camera transform and 3D-to-2D projection yield image BBox  $\hat{b}_{\text{image}}$ .

2. **3D BBox Denoising:** While the cross-view feature association is simplified due to the 2D-to-3D lifting, the denoising process may be more challenging. In turn, the associated radar features are used as conditioning to alleviate the more challenging 3D BBox denoising. To the best of our knowledge, REXO is the first diffusion model in the multi-view radar perception field.

We demonstrate the effectiveness of our contributions through evaluations on two open radar datasets.

## 2. REXO: BBox Diffusion in 3D Radar Space

DiffusionDet [7] reformulates object detection as a denoising diffusion process [17, 32], treating  $x_t$  as 2D BBox parameters instead of image pixels. We extend this to multi-view radar by lifting  $x_t$  to 3D BBox in radar coordinates system. Conditioned on radar heatmaps (see in Fig. 2), REXO performs 3D BBox diffusion in two phases: 1) a **forward process** that adds noise to ground-truth (GT) BBox  $x_0$  to produce random  $x_T$  during training, and 2) a **reverse process** that denoises random  $x_T$  to estimate noise-free  $\hat{x}_0$  during inference. The denoised BBox is also projected to the 2D image plane for supervision in both radar and image domains. We describe REXO in two parts: training and inference.

### 2.1. Training

**Backbone:** As illustrated in Fig. 2, we first generate two radar heatmaps (horizontal  $Y_{\text{hor}} \in \mathcal{R}^{M \times W \times D}$  and vertical  $Y_{\text{ver}} \in \mathcal{R}^{M \times H \times D}$  where  $M$ ,  $W$ ,  $H$  and  $D$  denote the number of consecutive frames, width, height and depth, respectively. More details are described in Appendix B) derived from raw data captured by horizontal and vertical radar arrays. Taking the two radar heatmaps  $Y_{\text{hor}} \in \mathcal{R}^{M \times W \times D}$  and  $Y_{\text{ver}} \in \mathcal{R}^{M \times H \times D}$  as inputs, a shared backbone network (e.g., ResNet [14]) generates horizontal-view and vertical-view radar feature maps:  $Z_{\text{hor}} = \text{backbone}(Y_{\text{hor}})$  and  $Z_{\text{ver}} = \text{backbone}(Y_{\text{ver}})$ .

**Initialization of  $x_0$  and Forward Process to  $x_t$  with Ground-Level Constraint:** For a given number of BBoxes  $N_{\text{train}}$  to be detected,  $x_0$  is simply initialized by the 3D BBox GT in the radar space  $x_{\text{radar}} = \{c_x, c_y, c_z, w, h, d\}^\top \in \mathbb{R}^6$  and padded with random 3D BBox  $x_{\text{rand}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_6)$  if  $N_{\text{train}} > N_{\text{GT}}$ . The diffused 3D BBox  $x_t$  at time  $t$  can be generated as

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \quad (1)$$

where  $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_6)$ , and  $\bar{\alpha}_t$  denotes the noise variance schedule. Since the BBox is now explicitly defined in the 3D radar coordinate system, it is natural to incorporate prior knowledge as a constraint into the diffusion process. Unlike DiffusionDet and RETR, we enforce the reduced five 3D parameters by grounding with  $h^t/2$ , allowing 3D and 2D gradients to flow jointly and guiding the denoising process under strict geometric constraints. This ensures that objects are correctly positioned on the floor, reflecting realistic spatial relationships:  $x_t = \{c_x^t, h^t/2, c_z^t, w^t, h^t, d^t\}^\top$  (see the Ground-Level Constraint in Fig. 2).

**Cross-View Radar-Conditioned BBox Detector:**  $\text{DenoisingDet}_\theta$  includes explicit cross-view feature association and radar-conditioned 3D BBox detector. Given the noisy 3D BBox  $x_t$  in (1), the  $x_t$ -guided cross-view feature association first projects  $x_t$  onto the two radar views, resulting in two 2D BBoxes,

$$x_{t,\text{hor}} = \{c_x^t, c_z^t, w^t, d^t\}^\top, x_{t,\text{ver}} = \{c_y^t, c_z^t, h^t, d^t\}^\top, \quad (2)$$

and then crops out the cross-view 2D radar features:  $Z_{\text{hor/ver}}^{\text{crop}} = \text{RoIAlign}(Z_{\text{hor/ver}}, x_{t,\text{hor/ver}}) \in \mathbb{R}^{C \times r \times r}$  via a standard ROIAlign operation [15], where  $r$  denotes a fixed spatial resolution, e.g.,  $r = 7$ . At time  $t$ , this process yields  $N_{\text{train}}$  pairs of associated radar features

$$Z_{\text{radar}}^{\text{crop}} = \{Z_{\text{hor}}^{\text{crop}}, Z_{\text{ver}}^{\text{crop}}\} \in \mathbb{R}^{C \times r \times 2r}, \quad (3)$$

each corresponding to a noisy 3D BBox  $x_t$ . Conditioned on  $Z_{\text{radar}}^{\text{crop}}$ , a  $\text{DenoisingDet}_\theta$  with learnable weights  $\theta$  is

Table 1. Evaluation on 4 data splits of the MMVR dataset and WALK of the HIBER dataset.

Method	MMVR:P1S1			MMVR:P1S2			MMVR:P2S1			MMVR:P2S2			HIBER:WALK		
	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>	AP	AP <sub>50</sub>	AP <sub>75</sub>
RFMask	25.53	67.30	15.86	24.46	66.82	11.22	31.37	61.50	27.48	6.03	22.77	0.88	17.77	52.46	6.78
RFMask3D	34.84	69.57	31.74	30.75	76.48	16.23	39.89	80.38	35.35	12.26	37.01	4.34	16.58	48.10	6.53
DETR	35.64	77.59	28.00	28.51	75.90	13.42	29.53	63.08	25.35	9.29	34.69	2.49	14.45	47.33	4.25
RETR	39.62	80.55	33.84	30.16	78.95	15.17	46.75	83.80	46.06	12.45	41.30	4.96	22.09	59.83	10.99
<b>REXO</b>	39.23	73.46	37.83	36.48	87.02	20.51	48.35	85.89	48.38	23.47	64.41	10.44	25.33	62.55	15.83

trained to estimate the BBox  $\hat{x}_0$  and the class scores  $\hat{p}$  as

$$\{\hat{x}_0, \hat{p}\} = \text{DenoisingDet}_\theta(x_t, t, Z_{\text{radar}}^{\text{crop}}), \quad (4)$$

where  $t$  specifies the timestep embedding. In our indoor setting, we use a two-class softmax over  $\{\text{person}, \text{background}\}$ . The class-head can extend to  $C$  classes (including background) by using a  $C$ -way softmax with cross-entropy.

### 3D-to-2D Projection with Learnable Refinement:

REXO further projects  $\hat{x}_0$  in (4) into the 2D image plane. By setting  $\hat{x}_{\text{radar}} = \hat{x}_0$ , we convert each of the 8 corners of the corresponding 3D BBox  $\hat{x}_{\text{radar}}$  using  $x_{\text{camera}}^i = R\hat{x}_{\text{radar}}^i + v$ , where  $\hat{x}_{\text{radar}}^i$  is the  $i$ -th corner of  $\hat{x}_{\text{radar}}$ ,  $R$  and  $v$  are the calibrated 3D rotation matrix and translation vector: Each 3D corner  $x_{\text{camera}}^i$  is projected to the image plane through the calibrated pinhole model:

$$b_{\text{init}} = \{\bar{c}_x, \bar{c}_y, \bar{w}, \bar{h}\}^\top = \text{proj}_{\text{init}}(x_{\text{camera}}). \quad (5)$$

Since  $b_{\text{init}}$  systematically overshoots the ground-truth extent (see Appendix C), we attach a refinement module with learnable parameter  $\phi$  to obtain the offset:

$$\Delta b = \{\Delta \bar{x}, \Delta \bar{y}, \Delta \bar{w}, \Delta \bar{h}\}^\top = \text{Refinement}_\phi(f), \quad (6)$$

where  $f = \text{Predictor}(e_t, Z_{\text{radar}}^{\text{crop}})$  is the time-dependent feature.  $e_t$  denotes the timestep embedding [17] and  $\text{Predictor}$  denotes the time-dependent predictor [7] with the radar feature and the embedding. Applying these offsets produces the final image-plane box  $\hat{b}_{\text{image}}$ , achieving tighter alignment without sacrificing geometric consistency.

$$\hat{b}_{\text{image}} = \{\bar{c}_x + \bar{w}\Delta \bar{x}, \bar{c}_y + \bar{h}\Delta \bar{y}, e^{\Delta \bar{w}}\bar{w}, e^{\Delta \bar{h}}\bar{h}\}^\top. \quad (7)$$

**Loss:** To ensure consistency between the radar and image plane representations, we adopt a simplified scheme of the Tri-plane loss [40] that directly calculates the loss of 3D BBox. REXO employs the Hungarian match cost [21] with a loss function computed in both the 3D and 2D spaces:  $\mathcal{L}_{\text{box}}^{\text{GA}} = \lambda_{3D}\mathcal{L}_{\text{box}}^{3D}(x_{\text{radar}}, \hat{x}_{\text{radar}}) + \lambda_{2D}\mathcal{L}_{\text{box}}^{2D}(b_{\text{image}}, \hat{b}_{\text{image}})$ , where the 3D/2D BBox loss is defined as  $\mathcal{L}_{\text{box}}^*(x, \hat{x}) = \lambda_{\text{GIoU}}\mathcal{L}_{\text{GIoU}}(x, \hat{x}) + \lambda_{L_1}\mathcal{L}_{L_1}(x, \hat{x})$  representing a weighted combination of the generalized intersection over union (GIoU) loss  $\mathcal{L}_{\text{GIoU}}$  [27] and the  $\ell_1$  loss  $\mathcal{L}_{L_1}$ , and the coefficients  $\lambda$  balance the relative contribution of each loss term.

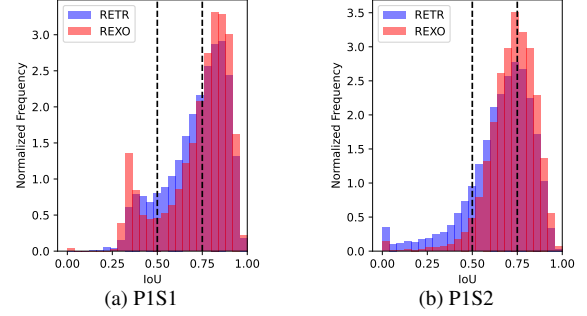


Figure 3. AP breakdowns with IoU histograms on MMVR.

## 2.2. Inference

REXO infers objects by reversing the diffusion process. Given a target count  $N$ , we sample random 3D boxes  $x_T \sim \mathcal{N}(0, I_6)$  in the radar coordinate system at  $t = T$  and denoise them down to  $t = 1$ . With  $x_t$  and radar features  $\{Z_{\text{hor}}, Z_{\text{ver}}\}$ , the trained  $\text{DenoisingDet}_\theta$  in (4) predicts  $\hat{x}_0$ , giving

$$p_\theta(x_{t-1} | x_t, Z_{\text{hor}}, Z_{\text{ver}}) = \mathcal{N}(\sqrt{\alpha_{t-1}}x_0 + \gamma\epsilon_\theta^{(t)}, \sigma_t^2 I_6),$$

$$x_{t-1} = \sqrt{\alpha_{t-1}}\hat{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta^{(t)} + \sigma_t \epsilon_t, \quad (8)$$

where  $\epsilon_\theta^{(t)} = (x_t - \sqrt{\alpha_t}\hat{x}_0)/\sqrt{1 - \alpha_t}$  specifies the direction pointing to the noisy BBox  $x_t$  at time  $t$ , and  $\epsilon_t \sim \mathcal{N}(0, I_6)$  represents a random BBox. Note that the denoising step is inherently conditioned on the cross-view radar feature maps via the estimated  $\hat{x}_0$  from the  $\text{DenoisingDet}_\theta$  module. After the final step,  $x_0 (= \hat{x}_{\text{radar}})$  is converted to image plane boxes  $\hat{b}_{\text{image}}$  via the radar-to-camera transform and the 3D-to-2D projection. Boxes whose class scores exceed a threshold are output as detections.

## 3. Experiments

We demonstrate the effectiveness of REXO through evaluations on two open high-resolution radar datasets.

**Datasets:** *MMVR* [26] includes multi-view radar heatmaps collected from over 25 human subjects across 6 rooms over a span of 9 days. It consists of 345K data frames collected in 2 protocols: P1: Open Foreground) with 107.9K frames in an open-foreground room with a single subject; and P2: Cluttered Space with 237.9K frames in 5 cluttered rooms with single and multiple subjects.

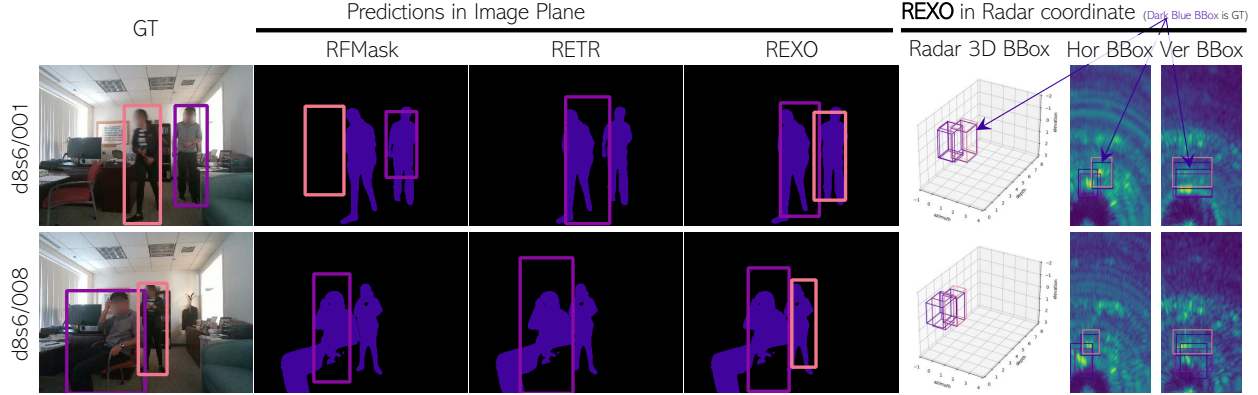


Figure 4. Visualization of unseen frames in P2S2 of MMVR: The left column shows the radar heatmaps, followed by the second column displaying predicted/GT 3D BBoxes in the radar space. Corresponding image-plane 2D BBox predictions are shown in the middle column for two baselines (RFMask and RETR) and REXO, with purple segmentation masks overlaid to illustrate the alignment with human GT. The right column presents the RGB images with GT 2D BBoxes for qualitative check.

Under each protocol, two data splits are defined to evaluate radar perception performance: S1: a random data split and S2: a cross-session, unseen split. *HIBER* [37], partially released, includes multi-view radar heatmaps from 10 human subjects in a single room but from different angles with multiple data splits. In our evaluation, we used the “WALK” data split, consisting of 73.5K data frames with one subject walking in the room.

**Implementation:** We consider RFMask [37], DETR [5] and RETR [40] as baseline methods. Additionally, we evaluate a 3D extension of RFMask (RFMask3D; see Appendix D), that takes the two radar views as inputs for BBox prediction. Hyperparameter settings are provided in Appendix E.

**Metrics:** We evaluate performance using average precision (AP) at two IoU thresholds of 0.5 ( $AP_{50}$ ) and 0.75 ( $AP_{75}$ ), along with the mean AP (AP) computed over thresholds in the range of  $[0.5 : 0.05 : 0.95]$ . For detailed metric definitions, refer to Appendix F.

**Result of MMVR:** Table 1 presents the results under the four combinations of two protocols and two data splits of the MMVR dataset. REXO demonstrates significant performance improvements in P1S2, P2S1, and P2S2. Notably, in P2S2 where the test radar frames contain an entirely unseen environment during training, REXO outperforms the best baseline RETR by a large margin, boosting AP from 12.45 to 23.47, highlighting its strong generalization capabilities. Surprisingly, under the simplest combination P1S1 where a single subject is recorded in the same room with a random data split, REXO’s performance is slightly lower than that of RETR, particularly on the metric  $AP_{50}$ . To understand these differences, we break down the AP into IoU histograms for (a) P1S1 and (b) P1S2, as illustrated in

Fig. 3, where blue and red histograms represent the IoU distributions for RETR and REXO, respectively, and the left and right dotted lines mark the two IoU thresholds at 0.5 and 0.75. It is seen that in Fig. 3a, the excess of RETR over REXO (blue areas) over the IoU interval  $[0.5, 0.75]$  is greater than that of REXO over RETR (pink areas) over the interval  $[0.75, 1.0]$ , explaining RETR’s higher  $AP_{50}$  under P1S1. Meanwhile, REXO has better  $AP_{75}$  as it provides more high-quality predictions with IoU above 0.75.

**Result of HIBER:** Table 1 presents the results evaluated on the “WALK” data split of the HIBER dataset. As well as MMVR cases, REXO outperforms RETR across all evaluation metrics with an AP of 25.33, surpassing RETR’s AP at 22.09. REXO attains  $AP_{50}$  of 62.55 and  $AP_{75}$  of 15.83, demonstrating strong performance in both low- and high-IoU BBox performance evaluations.

**Visualization:** Fig. 4 visualizes selected “Unseen” frames from a room never encountered during training in P2S2. It is seen that 2D BBox predictions by REXO align more closely with human segmentation masks (purple pixels) than those of RETR and RFMask. This improvement is potentially due to the explicit cross-view feature association, which strengthens consistency across radar views even in new environments, yielding better generalization. More challenging examples are provided in Appendix G.

## 4. Conclusion

We proposed REXO, a multi-view radar object detection method that refines the 3D BBox through a diffusion process. By explicitly guiding cross-view radar feature association, REXO achieves consistent performance improvements on two open indoor radar datasets over a list of strong baselines.



## References

- [1] Tomer Amit, Tal Shaharbandy, Eliya Nachmani, and Lior Wolf. SegDiff: Image segmentation with diffusion probabilistic models, 2022. [7](#)
- [2] Dan Barnes, Matthew Gadd, Paul Murcutt, Paul Newman, and Ingmar Posner. The Oxford radar robotcar dataset: A radar extension to the Oxford robotcar dataset. In *International Conference on Robotics and Automation*, pages 6433–6438, 2020. [7](#)
- [3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 22563–22575, 2023. [7](#)
- [4] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11621–11631, 2020. [7](#)
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision (ECCV)*, page 213–229, 2020. [4](#), [8](#)
- [6] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2416–2425, 2023. [7](#)
- [7] Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. DiffusionDet: Diffusion model for object detection. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 19773–19786, 2023. [1](#), [2](#), [3](#), [7](#)
- [8] Guoxuan Chi, Zheng Yang, Chenshu Wu, Jingao Xu, Yuchong Gao, Yunhao Liu, and Tony Xiao Han. RF-diffusion: Radio signal generation via time-frequency diffusion, 2024. [7](#)
- [9] Hongyu Deng, Tianfan Xue, and He Chen. Fusegrasp: Radar-camera fusion for robotic grasping of transparent objects. *IEEE Transactions on Mobile Computing*, 24(8):7028–7041, 2025. [1](#)
- [10] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The PASCAL Visual Object Classes (VOC) Challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010. [9](#)
- [11] Junqiao Fan, Jianfei Yang, Yuecong Xu, and Lihua Xie. Diffusion model is a good pose estimator from 3D RF-vision. In *Computer Vision – ECCV 2024*, pages 1–18, Cham, 2025. [7](#)
- [12] Xiangyu Gao, Guanbin Xing, Sumit Roy, and Hui Liu. RAMP-CNN: A novel neural network for enhanced automotive radar object recognition. *IEEE Sensors Journal*, 21(4): 5119–5132, 2021. [7](#)
- [13] Zhangxuan Gu, Haoxing Chen, and Zhuoer Xu. Diffusion-Inst: Diffusion model for instance segmentation. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2730–2734, 2024. [7](#)
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. [2](#), [7](#)
- [15] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 2980–2988, 2017. [2](#)
- [16] Cheng-Ju Ho, Chen-Hsuan Tai, Yen-Yu Lin, Ming-Hsuan Yang, and Yi-Hsuan Tsai. Diffusion-SS3D: Diffusion model for semi-supervised 3D object detection. In *Advances in Neural Information Processing Systems*, pages 49100–49112, 2023. [7](#)
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, pages 6840–6851, 2020. [2](#), [3](#)
- [18] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *Advances in Neural Information Processing Systems*, pages 8633–8646, 2022. [7](#)
- [19] Jan Hosang, Rodrigo Benenson, Piotr Dollár, and Bernt Schiele. What makes for effective detection proposals? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(4):814–830, 2016. [9](#)
- [20] S. M. Kay. *Fundamentals of Statistical Signal Processing: Detection Theory*. Prentice Hall, 1998. [7](#)
- [21] Harold W. Kuhn. The Hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. [3](#)
- [22] Wuyang Li, Xinyu Liu, Jiayi Ma, and Yixuan Yuan. Cliff: Continual latent diffusion for open-vocabulary object detection. In *ECCV*, 2024. [7](#)
- [23] Chris Xiaoxuan Lu, Stefano Rosa, Peijun Zhao, Bing Wang, Changhao Chen, John A. Stankovic, Niki Trigoni, and Andrew Markham. See through smoke: robust indoor mapping with low-cost mmWave radar. In *The 18th International Conference on Mobile Systems, Applications, and Services (MobiSys)*, page 14–27, 2020. [1](#)
- [24] Kai Luan, Chenghao Shi, Neng Wang, Yuwei Cheng, Huimin Lu, and Xieyuanli Chen. Diffusion-based point cloud super-resolution for mmwave radar data. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11171–11177, 2024. [7](#)
- [25] Dong-Hee Paek et al. K-Radar: 4D radar object detection for autonomous driving in various weather conditions. In *NeurIPS*, pages 3819–3829, 2022. [7](#)
- [26] M. Mahbubur Rahman, Ryoma Yataka, Sorachi Kato, Pu Wang, Peizhao Li, Adriano Cardace, and Petros Boufounos. MMVR: Millimeter-wave multi-view radar dataset and benchmark for indoor perception. In *European Conference on Computer Vision (ECCV)*, pages 306–322, 2024. [3](#), [7](#)
- [27] Hamid Reza Tofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 658–666, 2019. [3](#)

- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 7
- [29] Marcel Sheeny, Emanuele De Pellegrin, Saptarshi Mukherjee, Alireza Ahrabian, Sen Wang, and Andrew Wallace. RA-DIATE: A radar dataset for automotive perception in bad weather. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1–7, 2021. 7
- [30] Guangsheng Shi, Ruifeng Li, and Chao Ma. PillarNet: Real-time and high-performance pillar-based 3D object detection. In *European Conference on Computer Vision (ECCV)*, page 35–52, 2022. 7
- [31] Mikael Skog, Oleksandr Kotlyar, Vladimír Kubelka, and Martin Magnusson. Human detection from 4D radar data in low-visibility field conditions. *arXiv:2404.05307*, 2024. 1
- [32] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 2, 7
- [33] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019.
- [34] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 32211–32252. PMLR, 2023. 7
- [35] Dayi Tan, Hansheng Chen, Wei Tian, and Lu Xiong. DiffusionRegPose: Enhancing multi-person pose estimation using a diffusion-based end-to-end regression approach. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2230–2239, 2024. 7
- [36] Jincheng Wu, Ruixu Geng, Yadong Li, Dongheng Zhang, Zhi Lu, Yang Hu, and Yan Chen. Diffradar: high-quality mmWave radar perception with diffusion probabilistic model. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8291–8295, 2024. 7
- [37] Zhi Wu, Dongheng Zhang, Chunyang Xie, Cong Yu, Jinbo Chen, Yang Hu, and Yan Chen. RFMask: A simple baseline for human silhouette segmentation with radio signals. *IEEE Transactions on Multimedia*, 25:4730–4741, 2023. 1, 4, 7, 8, 9
- [38] Xinhao Xiang, Simon Dräger, and Jiawei Zhang. 3DiffusionDet: Diffusion model for 3D object detection with robust lidar-camera fusion. *ArXiv*, abs/2311.03742, 2023. 7
- [39] Bo Yang, Ishan Khatri, Michael Happold, and Chulong Chen. ADCNet: Learning from raw radar data via distillation. *arXiv:2303.11420*, 2023. 7
- [40] Ryoma Yataka, Adriano Cardace, Perry Wang, Petros Boufounos, and Ryuhei Takahashi. RETR: Multi-view radar detection transformer for indoor perception. In *Advances in Neural Information Processing Systems*, pages 19839–19869, 2024. 1, 3, 4, 7, 9
- [41] Jason J. Yu, Fereshteh Forghani, Konstantinos G. Derpanis, and Marcus A. Brubaker. Long-term photometric consistent novel view synthesis with diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7094–7104, 2023. 7
- [42] Ruibin Zhang, Donglai Xue, Yuhang Wang, Ruixu Geng, and Fei Gao. Towards dense and accurate radar perception via efficient cross-modal diffusion model. *IEEE Robotics and Automation Letters*, 9(9):7429–7436, 2024. 7
- [43] Mingmin Zhao, Tianhong Li, Mohammad Abu Alsheikh, Yonglong Tian, Hang Zhao, Antonio Torralba, and Dina Katabi. Through-wall human pose estimation using radio signals. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7356–7365, 2018. 7
- [44] Peijun Zhao, Chris Xiaoxuan Lu, Bing Wang, Niki Trigoni, and Andrew Markham. Cubelearn: End-to-end learning for human motion recognition from raw mmWave radar signals. *IEEE Internet of Things Journal*, 10(12):10236–10249, 2023. 7