



# Robust Bird’s Eye View Segmentation by Adapting DINOv2

Merve Rabia Barın<sup>1,2</sup>, Görkay Aydemir<sup>1</sup>, and Fatma Güney<sup>1,2</sup>

<sup>1</sup> Department of Computer Engineering, Koç University

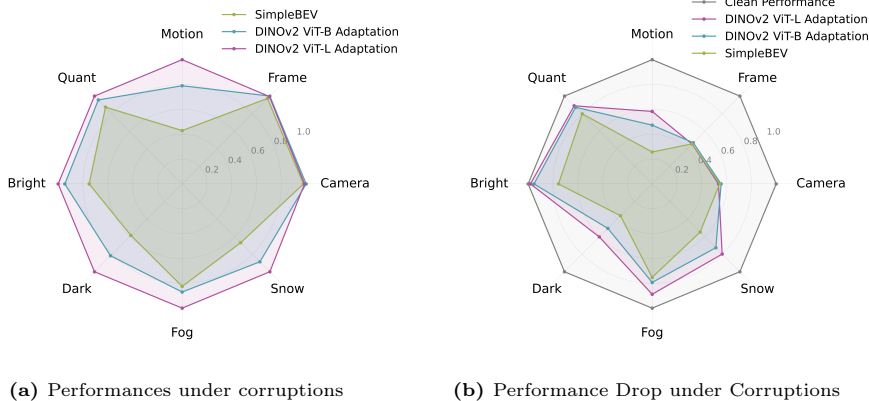
<sup>2</sup> KUIS AI Center

**Abstract.** Extracting a Bird’s Eye View (BEV) representation from multiple camera images offers a cost-effective, scalable alternative to LIDAR-based solutions in autonomous driving. However, the performance of the existing BEV methods drops significantly under various corruptions such as brightness and weather changes or camera failures. To improve the robustness of BEV perception, we propose to adapt a large vision foundational model, DINOv2, to BEV estimation using Low Rank Adaptation (LoRA). Our approach builds on the strong representation space of DINOv2 by adapting it to the BEV task in a state-of-the-art framework, SimpleBEV. Our experiments show increased robustness of BEV perception under various corruptions, with increasing gains from scaling up the model and the input resolution. We also showcase the effectiveness of the adapted representations in terms of fewer learnable parameters and faster convergence during training.

## 1 Introduction

Accurate perception of the surrounding scene in 3D is crucial for safe navigation in autonomous driving. LIDAR sensors can provide accurate 3D measurements, however, due to their high cost and consequently scalability issues, camera-based solutions are explored as an alternative. Specifically, substantial efforts have been directed toward extracting bird’s-eye view (BEV) representations from multi-camera images [6, 12, 17, 20, 29], providing a more cost-effective solution. BEV representations are assumed as input to motion prediction methods [5, 7] and are increasingly used as part of end-to-end driving systems [8, 9, 11, 14, 28].

While the robustness of these models under various conditions is a critical factor for ensuring safety, recent work [24] shows that BEV perception models suffer from performance degradation when exposed to different types of corruption such as brightness changes, adversarial weather conditions, motion blur, quantization, frame loss, and camera crash, highlighting a significant challenge. The accuracy of BEV perception plays a crucial role in both motion prediction [25] and end-to-end driving. While privileged agents that have access to ground truth BEV demonstrate an impressive driving performance, their student counterparts suffer from mistakes in predicted BEV maps [28]. Similarly, the performance of motion prediction methods drops notably while switching from ground truth BEV to the predicted BEV [25].



**Fig. 1: Robustness Analysis on nuScenes-C.** We compare the models under different types of corruptions in a, where each axis is normalized over the maximum performing model, *i.e.* ViT-L adaptation. We show the performance drop of models relative to their performance on clean data in b, where each axis is normalized to the clean data performance of the corresponding model.

The availability of large-scale datasets has facilitated the emergence of visual foundation models, renowned for their generalization capabilities. Notably, DINOv2 [19] stands out for its robust, general-purpose visual features, making it a suitable choice for various tasks such as zero-shot correspondence estimation [27], robotics [4], object-centric learning [1], point tracking [2], and object segmentation [18]. Despite the rich representation capacity of DINOv2, only a few works [22, 23] have explored its potential for BEV segmentation. In this work, we aim to explore the effectiveness of DINOv2 for robust BEV perception including performance, parameter efficiency, and convergence behavior.

We integrate DINOv2 into a state-of-the-art BEV segmentation model, SimpleBEV [12], by utilizing an efficient adaptation technique [13]. Specifically, we replace the backbone of SimpleBEV with DINOv2 for feature extraction and then efficiently update it using Low Rank Adaptation (LoRA). We systematically analyze the effectiveness of our approach by comparing the robustness of our adaptation to the original SimpleBEV with ResNet-101 backbone in terms of accuracy, parameter efficiency, and convergence behavior. Our experiments reveal that our adaptation improves the robustness of BEV perception under adversarial conditions with significantly lower learnable parameters and shorter training times.

## 2 Related Work

**Camera-based BEV Segmentation:** Bird’s-eye-view (BEV) representation is frequently used in autonomous driving to capture spatial information effectively. BEV summarizes the scene in a clear and compact representation by

extracting the necessary 3D information from 2D images. BEV methods first process camera images using a backbone, and then construct a discrete 3D representation of the surrounding scene by transforming 2D image features to 3D voxel grids. After constructing the 3D voxel grid, the features are decoded into a 2D BEV representation of the scene, after flattening the height dimension.

We can categorize the existing work on BEV into three based on how they extract 3D information from 2D images:

- i) *Depth-based* methods [16, 20], learn an explicit depth distribution at each pixel. They perform a weighted sum over the depth probabilities to back-project image features to 3D. However, learning an explicit depth distribution can be challenging due to the complexity of accurately modeling the depth for every pixel. Furthermore, voxel grids can only accumulate sparse features, which magnifies projection errors.
- ii) *Attention-based* models [3, 29] learn to align image features with voxel grids with an attention mechanism, leveraging camera-aware position embeddings to implicitly learn depth. However, learning projection via attention causes performance issues in attention-based approaches.
- iii) *Sampling-based* methods [12, 17] address these issues by sending rays from voxel grids to the images and bilinearly sampling the intersecting image features. Changing the direction of sampling increases the density of features in the voxel grid, relying on the quality of the features rather than the accuracy of projection.

Due to these reasons, we choose a sampling-based method, SimpleBEV [12] for our analysis.

**Vision Foundational Models in BEV:** Foundational models trained on large-scale data provide robust representations, improving performance in various downstream tasks [1, 4, 18, 27] due to their inherent semantic understanding and strong generalization capabilities. Although foundational models are trained solely on 2D data, they are shown to capture some 3D information from images, as tested on multi-view correspondence and depth estimation tasks [10, 26]. Among the foundational models evaluated, DINOv2 [19] performs the best for the 3D tasks considered, alongside Stable Diffusion [21], indicating its potential for 3D scene understanding. There is recent work building on foundational models for pre-training a BEV network for occupancy prediction [23] or sensor fusion [22] by benefiting from semantic capabilities of DINOv2 [19]. We directly target improving BEV estimation from camera images by adapting DINOv2.

### 3 Methodology

Our goal is to integrate visual foundation model DINOv2 [19] into BEV prediction model, SimpleBEV [12]. In this section, we first explain our adaptation strategy, Low Rank Adaptation (Section 3.1) and then SimpleBEV (Section 3.2) for completeness, and finally, present our approach to integrate DINOv2 into the SimpleBEV framework (Section 3.3).

### 3.1 Low Rank Adaptation (LoRA)

Low Rank Adaptation [13], *i.e.* LoRA, is widely used in Natural Language Processing to adapt pretrained large models to different tasks. LoRA is parameter efficient compared to fine-tuning, as only low-rank matrices are trained as residuals to frozen weights. Following the approach in MeLo [30], we update only the query and value projections in all attention layers in the ViT. Formally, given a pre-trained weight matrix  $\mathbf{W}_{\{Q,V\}} \in \mathbb{R}^{d \times d}$ , we obtain  $\mathbf{W}'$  as the result of adaptation:

$$\mathbf{W}'_{\{Q,V\}} = \mathbf{W}_{\{Q,V\}} + \mathbf{B}\mathbf{A} \quad (1)$$

where  $\mathbf{A} \in \mathbb{R}^{r \times d}$  and  $\mathbf{B} \in \mathbb{R}^{d \times r}$  are learned matrices,  $r$  is the rank, and  $d$  is the feature dimension. We train only the  $\mathbf{B}$  and  $\mathbf{A}$  matrices while keeping the original  $\mathbf{W}$  frozen for each attention layer.

### 3.2 SimpleBEV

Given images from  $N$  cameras, SimpleBEV [12] first extracts a feature map  $\mathbf{f}_i \in \mathbb{R}^{d \times H_f \times W_f}$  for each image  $i \in \{1, 2, \dots, N\}$ , where  $d$  represents the feature dimension, and  $H_f \times W_f$ , the size of the feature map. The method then employs a parameter-free lifting technique to transform image features into a 3D voxel grid  $\mathbf{V} \in \mathbb{R}^{d \times X \times Y \times Z}$ , where  $X$ ,  $Y$ , and  $Z$  correspond to the width, height and depth of the grid, respectively.

For each voxel  $\mathbf{V}_{\mathbf{p}} \in \mathbb{R}^d$  in the grid, represented by the 3D coordinate  $\mathbf{p} = [x, y, z]$ , the corresponding 2D pixel coordinates  $\mathbf{q}_i = (u_i, v_i)$  is calculated by projecting  $\mathbf{p}$  onto the 2D image plane of each camera  $i$ , using the intrinsic and extrinsic matrices  $\mathbf{K}_i$  and  $\mathbf{E}_i$ :

$$\mathbf{q}_i = \mathbf{K}_i \mathbf{E}_i \mathbf{p} \quad (2)$$

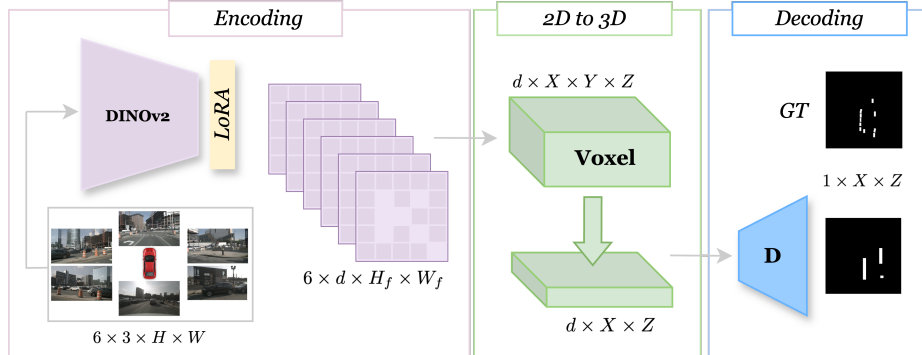
The feature values are then bilinearly sampled from the feature map  $\mathbf{f}_i$  at these projected coordinates,  $\mathbf{q}_i$ . Then, sampled features from all  $N$  cameras are aggregated and assigned to the corresponding voxel on the 3D grid:

$$\mathbf{V}_{\mathbf{p}} = \frac{1}{N} \sum_{i=1}^N \text{sample}(\mathbf{f}_i, \mathbf{q}_i) \quad (3)$$

After aggregating features for all  $\mathbf{V}_{\mathbf{p}} \in \mathbf{V}$  and constructing a 3D voxel grid  $\mathbf{V}$  that encapsulates the entire scene, a compressor reduces this voxel grid into a 2D BEV feature map  $\mathbf{B} \in \mathbb{R}^{d \times X \times Z}$ . The decoder then predicts the probability of occupancy for each grid cell. The model is trained using binary cross-entropy loss to optimize these predictions.

### 3.3 Adaptation of DINOv2 to BEV

**Integrating DINOv2:** We perform two main steps to integrate DINOv2 into SimpleBEV. First, we introduce additional layers in the query and key components of the ViT-based DINOv2 for adaptation, as described in Section 3.1.



**Fig. 2: Overview.** In this work, we propose to adapt DINOv2 to BEV segmentation using Low-Rank Adaptation (LoRA) for a robust BEV model. There are three main steps: i) We encode the camera images using DINOv2 to obtain tokens for each view, with attention weights updated through LoRA. ii) Transform image features from 2D to 3D using pull mechanism proposed by [12]. iii) Decode BEV features to 2D vehicle BEV masks.

Next, we replace the feature map  $\mathbf{f}_i \in \mathbb{R}^{d \times H_f \times W_f}$  with the output tokens of DINOv2 for each corresponding camera view, followed by pooling as shown in (3). During training, only the additional layers are updated, while the original DINOv2 model remains frozen. SimpleBEV, with its original ResNet-101 backbone, serves as the baseline for comparison.

**Evaluation Aspects:** We assess the quality of our adaptation by focusing on three key aspects in the evaluation:

- i) First, we consider **input resolution**, both image and feature map, where the ability to achieve high performance with lower-resolution inputs is advantageous, as it demonstrates efficiency in processing.
- ii) Second, we evaluate the **number of learnable parameters**, which not only indicates the resource efficiency of training but also reflects the robustness and general-purpose nature of the input features. Models with fewer parameters that still perform well suggest that the features are inherently strong, as they require minimal transformation.
- iii) Lastly, we examine **convergence speed**, favoring models that reach optimal performance quickly, as this reduces computational and time-related costs during training. Convergence is measured by the number of updates needed, with faster convergence indicating more effective use of general-purpose features.

As a result, we prioritize models that excel in using lower-resolution inputs, have fewer learnable parameters, and converge more quickly, as these traits highlight the efficiency and robustness of our adaptation.

**Table 1: Quantitative Results on nuScenes Validation Set.** This table compares the results of SimpleBEV and DINOv2 adaptations with different backbones. The ●●● represents the number of iterations in the full training of SimpleBEV, which corresponds to 25K steps, while ● represents one-third of it. See text for details.

Model	Backbone	Input Resolution	Feature Resolution	#Steps	mIoU	Exp.
SimpleBEV	ResNet-101	224 × 400	28 × 50	●●●	42.3	A
		448 × 800	56 × 100		47.4	B
DINOv2	ViT-B	224 × 392	16 × 28	●●●	42.3	C
	ViT-L				43.4	D
	ViT-L	224 × 392	16 × 28	●	43.2	E
		392 × 700	28 × 50		47.6	F
		448 × 784	32 × 56		47.7	G

## 4 Experiments

### 4.1 Experimental Setup

**Datasets and Metrics:** We conducted our experiments on the nuScenes [5] dataset, which is widely used for training and evaluating camera-based BEV segmentation methods. The dataset contains 28130 instances in the training set and 6019 instances in the validation set. For robustness analysis, we conducted experiments on the nuScenes-C benchmark [24], which is designed to measure the robustness of camera-based BEV perception models across eight types of corruptions under three levels of severity. The corruptions in this dataset are grouped into 8 categories: *brightness*, *darkness*, *fog*, *snow*, *motion blur*, *color quantization*, *camera crash*, and *frame loss*. For detailed descriptions of these augmentations, please refer to RoboBEV [24]. Following prior work, we use the mean Intersection-over-Union (mIoU) to evaluate model performance, measuring the overlap between our predictions and the ground truth boxes.

**Training Details:** We train the adaptation models using ViT-B and ViT-L architectures with the AdamW optimizer, a learning rate of  $1 \times 10^{-3}$ , and one cycle scheduler. For fine-tuning, we use a lower learning rate,  $1 \times 10^{-5}$ , and set the effective batch size to 16. For the ViT-B model, we followed the default training schedule of 25K updates as suggested in [12]. For the ViT-L model, we additionally explored shorter training time, using approximately one-third of the default steps to demonstrate the effect of adaptation on convergence with a large backbone.

### 4.2 Adaptation

We compare the performance of our adaptation models to the original SimpleBEV results in Table 1. Our comparisons focus on three key aspects: i) Same

input resolution, ensuring the models are evaluated under identical conditions; ii) Same feature map resolution to assess the efficiency of the spatial features encoded by the models; iii) Update iterations to understand the convergence speed and training efficiency of the models.

**Same Input Resolution:** Vision-based BEV perception models rely entirely on image input, making image resolution critically important. As demonstrated in [12], increasing the input resolution can improve performance up to a certain limit. Considering these, we choose to experiment on  $224 \times 400$  and  $448 \times 800$ . We compare the models with nearly identical input resolutions<sup>3</sup> to assess their effectiveness when processing a similar number of pixels for a fair evaluation. For the smaller resolution, SimpleBEV and DINOv2 ViT-B reach the same level of performance, 42.3 mIoU (A *vs.* C), while ViT-L surpasses SimpleBEV by 1.1 points (A *vs.* D). For the higher resolution, ViT-L outperforms SimpleBEV by a small margin, with mIoUs of 47.4 and 47.7, respectively (B *vs.* G). This demonstrates that the adapted DINOv2 backbones can reach the performance of SimpleBEV using similar resolutions.

**Same Feature Resolution:** In addition to input resolution, we also consider feature map resolution. By this, we refer to the resolution of  $\mathbf{f}_i$  as introduced in Section 3.2, which is the output of the backbone. Different backbones operate with different strides, *i.e.* the downsampling ratio of the final feature map. A backbone with a lower stride is expected to be advantageous [15], as it can better preserve the details of spatial information. Specifically, the SimpleBEV downsamples the image to  $1/8^{th}$  of the original input resolution, while DINOv2 downsamples to  $1/14^{th}$  due to patch size. Considering the same feature resolution of  $28 \times 50$ , the adapted ViT-L outperforms SimpleBEV by a significant margin of 5.3 IoU (A *vs.* F). This indicates that the adapted DINOv2 backbone can preserve spatial information more efficiently.

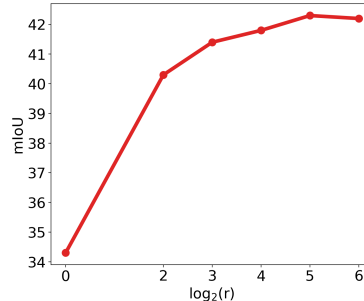
Moreover, increasing the feature resolution consistently improves performance, as expected, by providing more accurate and dense interpolation among features. For instance, raising the resolution from  $16 \times 28$  to  $28 \times 50$  results in a 4.2 mIoU increase (D *vs.* F). Additionally, a further resolution increase to  $32 \times 56$  yields a slight improvement, raising the mIoU from 47.6 to 47.7 (F *vs.* G).

**Number of Updates:** We chose the DINOv2 ViT-L model for convergence experiments due to its large scale, making it ideal for assessing convergence efficiency. The performance of DINOv2 ViT-L shows a drop of 0.2 IoU when trained only for one-third of the total iterations (D *vs.* E). In contrast, ViT-L not only surpasses SimpleBEV but does so even with fewer training iterations, given the same input resolution (A *vs.* E) and the same feature resolution (A *vs.* F). These results indicate that DINOv2 adaptation can converge quickly, and shorter training times have minimal impact on performance.

<sup>3</sup> The input resolution for ViTs should be divisible by the patch size, which is 14 in the case of DINOv2. Therefore, we use the closest multiples of 14 to match the target resolution.

**Table 2: Training Method and Parameter Efficiency.** This table shows the results of different weight update strategies with the corresponding number of learnable parameters. Note that there are additional 5M parameters for the decoder.

Backbone	Method	#Params	mIoU
ResNet-101	-	37M	42.3
ViT-B	Frozen	0	34.3
ViT-B	Fine-tune	86M	41.5
ViT-B	LoRA	1M	42.3
ViT-L	LoRA	3M	43.4



**Fig. 3: Varying the Rank of LoRA.** This plot illustrates the effect of increasing the LoRA rank (in log scale) on the performance, with rank 0 representing a frozen backbone.

### 4.3 Robustness Evaluation

In Figure 1, we present a robustness analysis comparing the SimpleBEV baseline with our two adaptation variants, ViT-B and ViT-L. We evaluate the models under various corruptions from the nuScenes-C dataset and report their performance separately to assess how each model handles different types of corruptions. For a fair comparison, all models are trained at similar resolutions:  $224 \times 400$  for SimpleBEV and  $224 \times 392$  for the ViT backbones, as discussed in Section 4.2.

In Figure 1a, the performance of all methods is normalized based on the best-performing model for each type of corruption. For color quantization, frame loss, and camera crash, all methods perform similarly. However, the ViT-L adaptation significantly outperforms other methods in most corruption types, surpassing SimpleBEV by at least 20%, with ViT-B also showing strong results. The difference is particularly pronounced in the case of motion blur, where SimpleBEV operates at only 40% of the ViT-L adaptation’s performance.

Figure 1b illustrate the performance drop of each model relative to its performance on clean data. This figure highlights the relative decline in performance for each model when exposed to various corruptions. As seen in the previous analysis, the ViT-B and ViT-L adaptations exhibit less performance degradation compared to SimpleBEV. Excluding the camera crash and frame loss scenarios, the ViT-L adaptation consistently maintains its performance, never dropping below 60% of its clean performance, with degradation of less than 20% for brightness, fog, and quantization. For the ViT-B adaptation, the threshold is around 45%. In contrast, SimpleBEV struggles to maintain its performance, with drops below 40% in motion blur conditions and below 30% in darkness. Overall, the DINOv2 adaptations demonstrate greater robustness than SimpleBEV in six out of eight corruption types and are comparable in the remaining two, camera crash and frame loss, highlighting the superior robustness of the adaptation approach. This



finding underscores the value of exploring foundational models like DINOv2 for enhancing robustness in BEV perception tasks.

#### 4.4 Ablation Study

**Training Method:** To highlight the impact of LoRA, we conducted experiments with the ViT-B DINOv2 adaptation across three configurations, as shown in Table 2: i) Frozen, where only the decoder (with 5M parameters) is trained and no additional learnable parameters are introduced, *i.e.* no updates to the backbone; ii) Fine-tuning, where all 86M parameters of the ViT-B DINOv2 backbone are updated; and iii) LoRA, where a small set of parameters is learned within the attention layers of the backbone.

The experiments reveal that, while the frozen model demonstrates a decent zero-shot representation performance, it significantly lags behind SimpleBEV (34.3 *vs.* 42.3). Fine-tuning the entire backbone offers improved results but demands significant computational resources due to the large number of learnable parameters (86M). In contrast, the LoRA configuration, which adds just 1M parameters to the decoder, achieves results that are on par with (ViT-B; 42.3) or even superior (ViT-L; 43.4) to SimpleBEV (42.3). This is achieved by updating only 1.12% of the parameters in ViT-B and 2.70% in ViT-L. Notably, LoRA outperforms the full fine-tuning by 0.8 points, showcasing its parameter efficiency and effectiveness in enhancing the performance.

**Rank of LoRA:** We experimented by varying the rank of adaptation in LoRA as shown in Figure 3. Rank essentially controls the capacity of the adaptation, with higher ranks providing more parameters for fine-tuning. Higher ranks can capture more complex relationships and lead to better performance, while potentially losing the information from pre-training. A rank of 0 corresponds to no updates to the backbone, *i.e.* frozen backbone. We found that increasing the rank leads to consistent improvements up to a certain point, specifically up to rank 32. However, increasing the rank from 32 to 64 results in a performance decrease of 0.1. A potential reason for this is the loss of the inductive bias of DINOv2, due to the larger updates on the attention weights. This indicates that rank 32 strikes an optimal balance between adapting the features to the task and preserving the valuable information acquired during pre-training.

## 5 Conclusion

We investigated the effectiveness of DINOv2 with Low Rank Adaptation (LoRA) for BEV segmentation on nuScenes. We first showed comparable results to SimpleBEV with a smaller backbone in the clean setting. By scaling up the model or the input and feature resolutions, we could obtain significant improvements compared to the baseline performance. Our experiments on nuScenes with corruptions (nuScenes-C) show increased robustness to various corruptions, even with

a shorter training time. These results justify our motivation to build on general-purpose features from large foundation models for improving BEV segmentation. Our approach requires updating significantly fewer parameters compared to full fine-tuning or the supervised baseline.

Our analysis is currently limited to DINOv2 and does not explore the performance of other promising 3D-aware foundation models, such as Stable Diffusion. Comparing different foundation models for BEV presents an interesting direction for further research, offering more insights about how to incorporate these models into BEV frameworks.

## Acknowledgements

This project is co-funded by KUIS AI and the European Union (ERC, ENSURE, 101116486). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them.

## References

1. Aydemir, G., Xie, W., Güney, F.: Self-supervised object-centric learning for videos. In: *Advances in Neural Information Processing Systems (NeurIPS)* (2023)
2. Aydemir, G., Xie, W., Güney, F.: Can visual foundation models achieve long-term point tracking? *arXiv preprint arXiv:2408.13575* (2024)
3. Bartoccioni, F., Zablocki, É., Bursuc, A., Pérez, P., Cord, M., Alahari, K.: Lara: Latents and rays for multi-camera bird’s-eye-view semantic segmentation. In: *Proc. Conf. on Robot Learning (CoRL)* (2023)
4. Blumenkamp, J., Morad, S., Gielis, J., Prorok, A.: Covis-net: A cooperative visual spatial foundation model for multi-robot applications. In: *Proc. Conf. on Robot Learning (CoRL)* (2024)
5. Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2020)
6. Chambon, L., Zablocki, E., Chen, M., Bartoccioni, F., Pérez, P., Cord, M.: Point-bev: A sparse approach for bev predictions. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2024)
7. Chang, M.F., Lambert, J., Sangkloy, P., Singh, J., Bak, S., Hartnett, A., Wang, D., Carr, P., Lucey, S., Ramanan, D., et al.: Argoverse: 3d tracking and forecasting with rich maps. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2019)
8. Chen, D., Krähenbühl, P.: Learning from all vehicles. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2022)
9. Chen, D., Zhou, B., Koltun, V., Krähenbühl, P.: Learning by cheating. In: *Proc. Conf. on Robot Learning (CoRL)* (2019)
10. El Banani, M., Raj, A., Maninis, K.K., Kar, A., Li, Y., Rubinstein, M., Sun, D., Guibas, L., Johnson, J., Jampani, V.: Probing the 3d awareness of visual foundation models. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2024)
11. Hanselmann, N., Renz, K., Chitta, K., Bhattacharyya, A., Geiger, A.: KING: generating safety-critical driving scenarios for robust imitation via kinematics gradients. In: *Proc. of the European Conf. on Computer Vision (ECCV)* (2022)
12. Harley, A.W., Fang, Z., Li, J., Ambrus, R., Fragkiadaki, K.: Simple-bev: What really matters for multi-sensor bev perception? In: *Proc. IEEE International Conf. on Robotics and Automation (ICRA)* (2023)
13. Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L., Chen, W.: LoRA: Low-rank adaptation of large language models. In: *Proc. of the International Conf. on Learning Representations (ICLR)* (2022)
14. Hu, Y., Yang, J., Chen, L., Li, K., Sima, C., Zhu, X., Chai, S., Du, S., Lin, T., Wang, W., Lu, L., Jia, X., Liu, Q., Dai, J., Qiao, Y., Li, H.: Planning-oriented autonomous driving. In: *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)* (2023)
15. Karaev, N., Rocco, I., Graham, B., Neverova, N., Vedaldi, A., Ruppel, C.: CoTracker: It is better to track together. In: *Proc. of the European Conf. on Computer Vision (ECCV)* (2024)
16. Li, Y., Ge, Z., Yu, G., Yang, J., Wang, Z., Shi, Y., Sun, J., Li, Z.: Bevdepth: Acquisition of reliable depth for multi-view 3d object detection. In: *Proc. of the Conf. on Artificial Intelligence (AAAI)* (2023)

17. Li, Z., Wang, W., Li, H., Xie, E., Sima, C., Lu, T., Qiao, Y., Dai, J.: Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In: Proc. of the European Conf. on Computer Vision (ECCV) (2022)
18. Nguyen, V.N., Groueix, T., Ponimatkin, G., Lepetit, V., Hodan, T.: Cnos: A strong baseline for cad-based novel object segmentation. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2023)
19. Oquab, M., Darcet, T., Moutakanni, T., Vo, H.V., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., Howes, R., Huang, P.Y., Xu, H., Sharma, V., Li, S.W., Galuba, W., Rabbat, M., Assran, M., Ballas, N., Synnaeve, G., Misra, I., Jegou, H., Mairal, J., Labatut, P., Joulin, A., Bojanowski, P.: DINOv2: Learning robust visual features without supervision. arXiv preprint arXiv:2304.07193 (2023)
20. Phillion, J., Fidler, S.: Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In: Proc. of the European Conf. on Computer Vision (ECCV) (2020)
21. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models (2021)
22. Schramm, J., Vödisch, N., Petek, K., Kiran, R.B., Yogamani, S., Burgard, W., Valada, A.: Bevcars: Camera-radar fusion for bev map and object segmentation. In: Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS) (2024)
23. Sirko-Galouchenko, S., Boulch, A., Gidaris, S., Bursuc, A., Vobecky, A., Pérez, P., Marlet, R.: Occfeat: Self-supervised occupancy feature prediction for pretraining bev segmentation networks. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2024)
24. Xie, S., Kong, L., Zhang, W., Ren, J., Pan, L., Chen, K., Liu, Z.: Robobev: Towards robust bird’s eye view perception under corruptions. arXiv preprint arXiv:2304.06719 (2023)
25. Xu, Y., Chambon, L., Zablocki, É., Chen, M., Alahi, A., Cord, M., Pérez, P.: Towards motion forecasting with real-world perception inputs: Are end-to-end approaches competitive? In: Proc. IEEE International Conf. on Robotics and Automation (ICRA) (2024)
26. Zhan, G., Zheng, C., Xie, W., Zisserman, A.: What does stable diffusion know about the 3D scene? arXiv preprint arXiv:2310.06836 (2023)
27. Zhang, J., Herrmann, C., Hur, J., Cabrera, L.P., Jampani, V., Sun, D., Yang, M.H.: A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence. In: Advances in Neural Information Processing Systems (NeurIPS) (2023)
28. Zhang, Z., Liniger, A., Dai, D., Yu, F., Van Gool, L.: End-to-end urban driving by imitating a reinforcement learning coach. In: Proc. of the IEEE International Conf. on Computer Vision (ICCV) (2021)
29. Zhou, B., Krähenbühl, P.: Cross-view transformers for real-time map-view semantic segmentation. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2022)
30. Zhu, Y., Shen, Z., Zhao, Z., Wang, S., Wang, X., Zhao, X., Shen, D., Wang, Q.: MeLo: Low-rank adaptation is better than fine-tuning for medical image diagnosis. arXiv preprint arXiv:2311.08236 (2023)