

ShelfGaussian: Shelf-Supervised Open-Vocabulary Gaussian-based 3D Scene Understanding

Supplementary Material

In this supplementary material, we first describe the robot hardware setup and our custom dataset collection process in Sec. 1. Then we conduct more ablation experiments in Sec. 2. More qualitative results are provided in Sec. 3. The implementation of our CUDA-accelerated Gaussian-to-voxel (G2V) splatting module is detailed in Sec. 4. Finally, we discuss the limitations and future directions in Sec. 5.

1. In-the-Wild Robot Experiment

1.1. UGV Hardware Setup

We illustrate the setup of our unmanned ground vehicle (UGV) in this section. A Clearpath Jackal robot is employed as our base UGV in this experiment. For collecting monocular RGB images, we use an on-board ZED 2i stereo camera. For collecting paired LiDAR point clouds, we utilize one Velodyne VLP16 LiDAR and one Hesai XT16 LiDAR. The layout of the coordinate frames of all sensors is given in Fig. 1. ROS2 [11] is adopted for the communication between the robot, sensors and personal computer. The intrinsics of ZED 2i camera is calibrated through built-in ZED SDK while the extrinsics of LiDARs are calibrated via open-source toolkit [5].

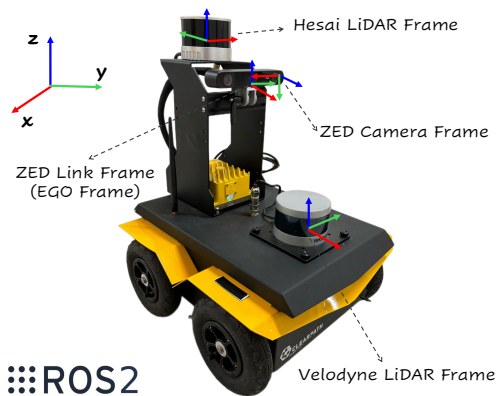


Figure 1. Visualization of the coordinate frames of different sensors and the ego vehicle. Red, green and blue arrows denote the x , y and z axes, respectively.

1.2. Custom Dataset Collection

By teleoperating our UGV, we collect a custom dataset in common urban scenarios. We choose four scenes: *street*, *park*, *grassland* and *garden*. We split our dataset into a 90% subset for training and a 10% subset for testing, resulting in 2638 and 294 samples, respectively. The frequency of the

key frames is set to 2Hz. We use our DINO-driven pseudo labeling engine to generate pseudo labels at key frames in a range of $[0\text{m}, 20\text{m}] \times [-10\text{m}, 10\text{m}] \times [-1\text{m}, 4\text{m}]$ with a voxel resolution of 0.2m. Each voxel label is decorated with 128-channel compressed DINO features. We visualize the scene reconstruction in Fig. 2.

2. Ablation Study

We conduct more ablation experiments in this section to study the effectiveness of individual modules and justify our architectural choices. Unless otherwise noted, all ablation studies are conducted on our model with camera-only input using the DINOv2 ViT-B/14 [12] model, 1000 Gaussian queries per view, and the Talk2DINO [2] language alignment method, ensuring consistency with our main paper.

2.1. Ablation on Number of Gaussians

We conduct an ablation study on the number of Gaussians per view to examine its effect on model accuracy. As shown in Tab. 1, with an increased number of Gaussian queries, the final model performance improves steadily, although some metrics exhibit minor fluctuations.

Mod.	Number of Gaussians			IoU	mIoU
	300	600	1000		
C	✓			61.42	17.85
		✓		61.07	18.26
			✓	63.25	19.07

Table 1. Ablation on the number of Gaussian queries.

2.2. Ablation on 2D and 3D Loss Weights

To study the effects of 2D and 3D losses on the final performance, we conduct experiments with different 3D loss weights. The weights for the BCE and feature losses are empirically selected based on their magnitudes relative to the 2D losses. As shown in Tab. 2, increasing 3D loss weights consistently boosts both IoU and mIoU.

2.3. Ablation on Language Alignment Methods

Talk2DINO [2] and DINO.txt [10] are two popular approaches for aligning natural language with DINO [12, 14] features. Talk2DINO [2] directly learns a mapping layer to transform CLIP [13] embeddings of text queries into the DINO feature space. DINO.txt [10] learns a vision block on top of the original DINO backbone and a text encoder for

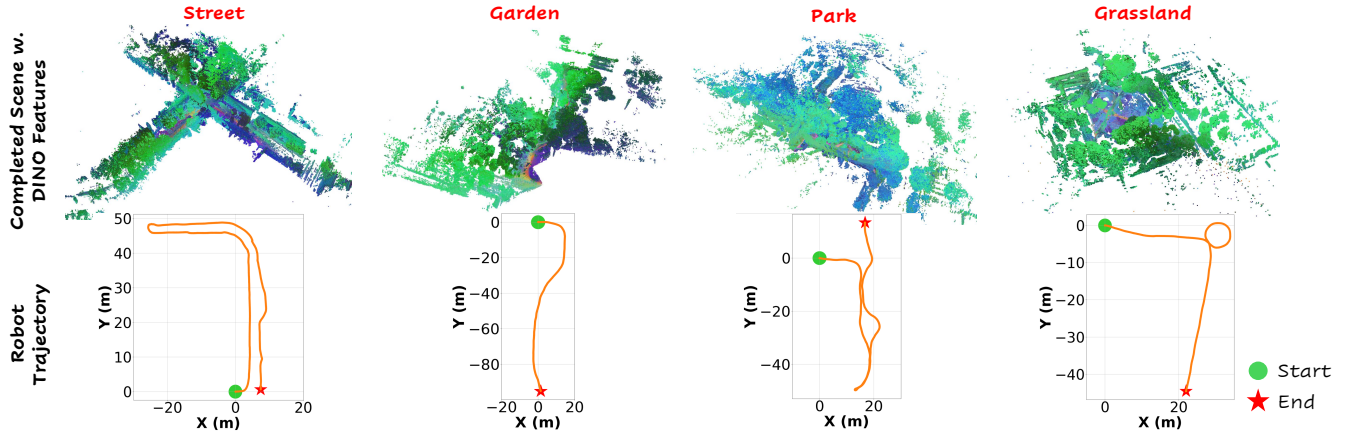


Figure 2. **Scene reconstruction results of four urban scenes.** The top row shows the completed scenes decorated with DINO features, visualized by mapping PCA components to RGB colors. The bottom row shows the robot trajectories within four urban scenes.

Mod.	2D Loss	3D Loss		IoU	mIoU
		BCE Loss	Feat. Loss		
C	1.0	1.0	1.0	58.66	17.52
	1.0	4.0	8.0	61.38	18.56
	1.0	8.0	16.0	63.25	19.07

Table 2. **Ablation on 2D and 3D loss weights.**

raw text queries, and then aligns the encoded vision and text features through contrastive learning. We test these two methods in our framework, and the results are shown in Tab. 3. With the same VFM backbone, our model using Talk2DINO [2] demonstrates better performance than using DINO.txt [10], especially in mIoU. This indicates that directly encoding Gaussians with DINO features leads to stronger semantic and physical representations.

Mod.	Sup.	VFM	Align.	IoU	mIoU
C	2D	DINOv3 ViT-L/16 [14]	Talk2DINO [2]	46.55	12.04
			DINO.txt [10]	47.08	5.71

Table 3. **Ablation on language alignment methods.** Since DINO.txt [10] only releases model weights based on DINOv3 ViT-L/16 [14], we conduct ablation experiments using this. PCA [1] is utilized to compress the feature dimension from 1024 to 128.

2.4. Ablation on Feature Compression Dimensions

To explore the influence of (PCA) [1] feature compression on the final model performance, we conduct ablation experiments with different compressed feature dimensions. As shown in Tab. 4, increasing compressed dimension generally improves the final model performance, particularly in the case of 2D and 3D joint supervision. However, the computational overhead increases as more feature channels are involved in the forward and backward passes. Considering

the tradeoff between computational efficiency and model performance, we eventually choose 128 as the final compressed feature dimension.

Mod.	Sup.		PCA Dim.		IoU	mIoU
	2D	3D	2D	3D		
C	✓		64	-	46.80	12.47
	✓	✓	64	64	58.96	18.18
	✓		128	-	47.26	12.39
	✓	✓	128	128	63.25	19.07

Table 4. **Ablation on PCA feature compression dimensions.**

2.5. Ablation on Trajectory Planning

We further conduct ablation study on the effects of ego status and past trajectory on vehicle trajectory planning. As shown in Tab. 5, with the addition of vehicle state information, the planning performance is steadily enhanced. Compared to other methods using vehicle status, our Gaussian-Planner demonstrates comparable L2 error and lower collision rate, which indicates the effectiveness of object-centric Gaussian representations in reliable, collision-aware vehicle trajectory planning. However, removing vehicle state input leads to suboptimal performance, suggesting that combining object-centric Gaussians with basic vehicle status remains beneficial and offers an avenue for further refinement.

Method	Ego Status	Past Traj.	L2 (m) ↓				CR (%) ↓			
			1s	2s	3s	Avg.	1s	2s	3s	Avg.
Gaussian-Planner	✓		1.05	1.76	2.50	1.77	0.31	2.38	5.31	2.67
	✓	✓	0.16	0.33	0.59	0.36	0.03	0.23	0.84	0.37
			0.16	0.32	0.58	0.35	0.00	0.18	0.78	0.32

Table 5. **Ablation on open-loop trajectory planning.** L2 (m): the mean Euclidean distance between the predicted and ground-truth trajectories. CR (%): the percentage of the planned trajectory intersecting an obstacle. Lower is better for both metrics.

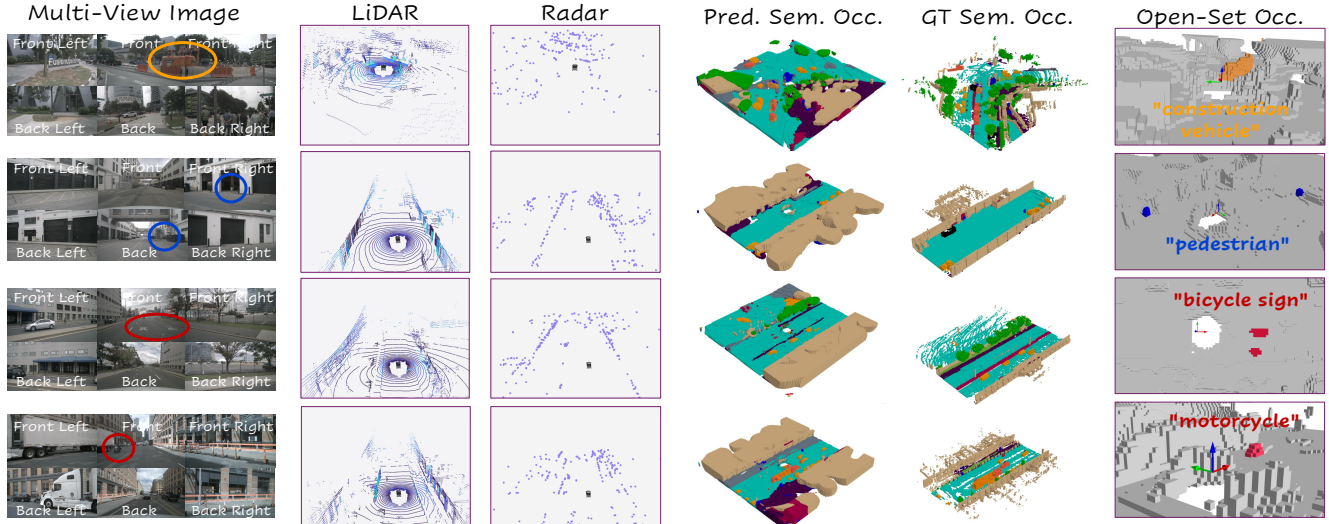


Figure 3. Qualitative results of ShelfGaussian on nuScenes [4] dataset validation split.

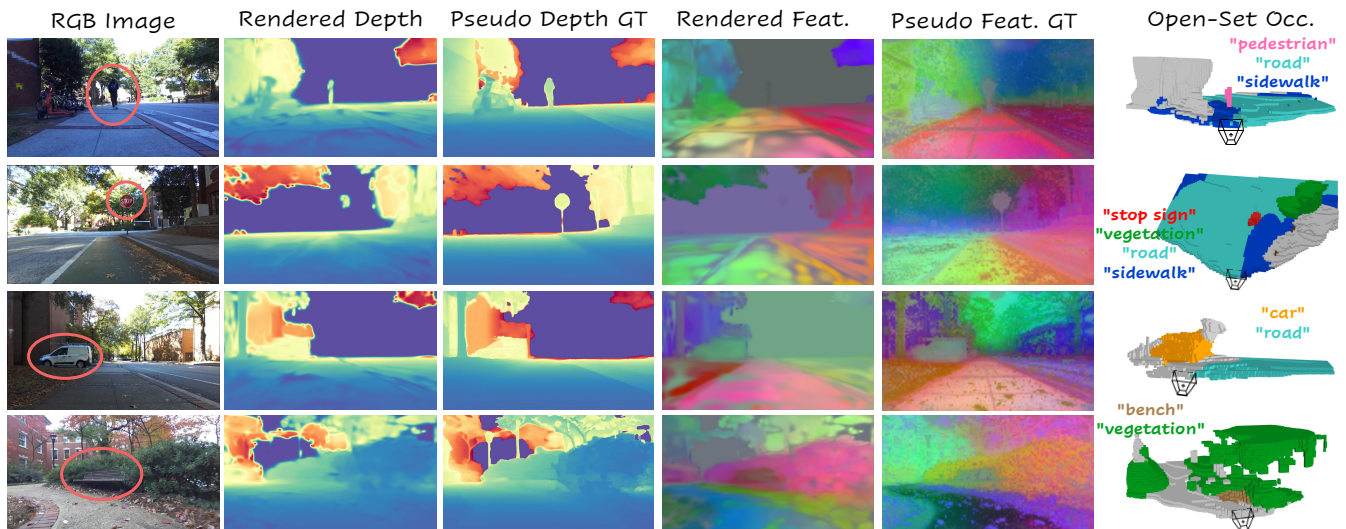


Figure 4. Qualitative results of ShelfGaussian on our custom dataset testing split.

2.6. Training Efficiency

Method	Mod.	Train. Time (h)	Memory (GB)	IoU	mIoU
GaussTR [9]	C	22	20	44.54	12.27
ShelfGaussian	C	25	15	63.25	19.07
	L	31	28	66.10	19.34
	C+R	31	28	62.84	19.42
	L+C	32	29	69.24	21.52
	L+C+R	32	29	69.45	21.78

Table 6. **Benchmark of training efficiency.** For a fair comparison, all experiments are conducted on 8 Nvidia L40S 48GB GPUs.

We compare the training time and training memory usage of our model with GaussTR [9], one of the most training-efficient self-supervised methods to date. As shown in Tab. 6, our method achieves better performance than GaussTR [9] with similar computational resources. As

more sensor modalities are incorporated, both model performance and training cost increase accordingly. Furthermore, we evaluate the inference efficiency of our method. The camera-only and LiDAR-camera-radar models achieve **3.3** and **2.3** FPS respectively, demonstrating the potential of real-time on-board deployment.

2.7. Gaussian-to-Voxel Splatting Efficiency

We further compare our designed G2V module with the ones used in Taichi [6] and GaussianFlowOcc [3], and the results are shown in Tab. 7. Under our challenging experimental setup, Taichi [6] fails to complete the operation due to exceeding the structural compiler limit, while GaussianFlowOcc [3] encounters out-of-memory issues. It is also worth mentioning that GaussianFlowOcc [3] models each Gaussian’s influence using an isotropic, fixed-size cube cen-

Method	18k Gauss. w. 1024D			9k Gauss. w. 768D		
	FW(s)	BW(s)	Mem.(GB)	FW(s)	BW(s)	Mem.(GB)
Taichi [6]	-	-	Fail	-	-	Fail
GaussianFlowOcc [3]	-	-	OOM	-	-	OOM
GaussTR [9]	14.1	596.7	22.3	12.4	446.4	14.8
GaussianFormer [7]	5.4	13.4	12.9	1.9	4.2	7.7
ShelfGaussian	0.5	2.7	4.9	0.2	1.0	3.7

Table 7. **Benchmark of G2V splatting module.** FW: forward time. BW: backward time. Mem.: memory. OOM: out of memory. Fail: exceed structural compiler limit. All the experiments are conducted on one Nvidia L40S 48GB GPU.

tered at its mean position. This approximation neglects the anisotropic covariance geometry of 3D Gaussians, leading to inaccurate computation.

3. Qualitative Results

3.1. nuScenes Dataset

We provide more qualitative results of ShelfGaussian on the nuScenes [4] dataset in Fig. 3. Our model is able to predict dynamic agents (e.g., *pedestrian*), long-tail objects (e.g., *construction vehicle*) and novel classes (e.g., *bicycle sign*) in urban driving scenarios.

3.2. Custom Dataset

We provide more qualitative results of ShelfGaussian on our custom dataset in Fig. 4. By querying the model with open-vocabulary categories such as “*pedestrian*” and “*stop sign*”, it successfully localizes and completes the objects.

4. Gaussian-to-Voxel Splatting Module

This section presents a comprehensive description of the gradient computation implemented in our CUDA-accelerated G2V splatting operator. We derive the gradients with respect to Gaussians’ means, covariance matrices, opacities, and high-dimensional features, respectively.

4.1. Preliminaries

For a Gaussian g evaluated at voxel center x_v , its unnormalized density contribution to voxel v is computed as:

$$w_{g,v} = \alpha_g \exp\left(-\frac{1}{2}(x_v - \mu_g)^\top \Sigma_g^{-1}(x_v - \mu_g)\right),$$

where $\mu_g \in \mathbb{R}^3$ is the mean, $\Sigma_g \in \mathbb{R}^{3 \times 3}$ is the covariance matrix, and α_g is the opacity. For convenience, we define:

$$d_{g,v} = x_v - \mu_g, \quad \phi_{g,v} = \exp\left(-\frac{1}{2}d_{g,v}^\top \Sigma_g^{-1}d_{g,v}\right).$$

The density contribution $w_{g,v}$ can be then expressed as:

$$w_{g,v} = \alpha_g \cdot \phi_{g,v}.$$

Each Gaussian g also encodes a high-dimensional feature vector $f_g \in \mathbb{R}^C$, where C is the feature dimension.

Voxel-wise Accumulators. For a fixed voxel v , we aggregate the contributions of all Gaussians to this voxel as:

$$F_v = \sum_g w_{g,v}, N_v = \sum_g w_{g,v} f_g, S_v = \max(F_v, \varepsilon).$$

Here, F_v is the accumulated density at voxel v , N_v is the density-weighted features, and S_v is the clamped denominator for numerical stability. ε is a constant, set as 1×10^{-6} .

Voxel Features. The final features of voxel v are defined as the normalized, density-weighted average of the Gaussian features:

$$G_v = \frac{N_v}{S_v} \in \mathbb{R}^C.$$

The upstream gradient with respect to the voxel features can be expressed as:

$$\bar{G}_v = \frac{\partial \mathcal{L}}{\partial G_v} \in \mathbb{R}^C.$$

In the following sections, we derive gradients of the loss \mathcal{L} with respect to different Gaussian parameters ($\mu_g, \Sigma_g, \alpha_g$) and features f_g using these voxel-wise accumulators.

4.2. Gradient w.r.t. Gaussian Features

N_v depends on the Gaussian features f_g , while S_v is independent of it. Thus, the gradients can be computed as:

$$\frac{\partial N_v}{\partial f_g} = w_{g,v} \cdot I_C, \quad \frac{\partial S_v}{\partial f_g} = 0,$$

and

$$\frac{\partial G_v}{\partial f_g} = \frac{1}{S_v} \frac{\partial N_v}{\partial f_g} = \frac{w_{g,v}}{S_v} \cdot I_C,$$

where I_C denotes the identity matrix. Applying the upstream gradient \bar{G}_v , we obtain the gradient with respect to the Gaussian feature f_g :

$$\frac{\partial \mathcal{L}}{\partial f_g} = \frac{\partial \mathcal{L}}{\partial G_v} \frac{\partial G_v}{\partial f_g} = \sum_v \left(\frac{w_{g,v}}{S_v} \right) \bar{G}_v.$$

4.3. Gradient w.r.t. Density Contribution

Based on $G_v = N_v/S_v$, by applying the quotient rule, we can obtain:

$$\frac{\partial G_v}{\partial w_{g,v}} = \begin{cases} \frac{f_g - G_v}{S_v}, & F_v > \varepsilon. \\ \frac{f_g}{\varepsilon}, & F_v \leq \varepsilon. \end{cases}$$

Therefore the gradient with respect to the density contribution $w_{g,v}$ can be calculated as:

$$s_{g,v} = \frac{\partial \mathcal{L}}{\partial G_v} \frac{\partial G_v}{\partial w_{g,v}} = \begin{cases} \bar{G}_v \cdot \frac{f_g - G_v}{S_v}, & F_v > \varepsilon. \\ \bar{G}_v \cdot \frac{f_g}{\varepsilon}, & F_v \leq \varepsilon. \end{cases}$$

4.4. Gradient w.r.t. Gaussian Mean

The density contribution of Gaussian g at voxel v is expressed as:

$$w_{g,v} = \alpha_g \exp\left(-\frac{1}{2}d_{g,v}^\top \Sigma_g^{-1} d_{g,v}\right),$$

where $d_{g,v} = x_v - \mu_g$. By taking the derivative of $w_{g,v}$ with respect to the Gaussian mean μ_g , we have:

$$\frac{\partial w_{g,v}}{\partial \mu_g} = w_{g,v} \Sigma_g^{-1} (x_v - \mu_g).$$

Thus, the gradient of the loss \mathcal{L} with respect to μ_g can be formed as:

$$\frac{\partial \mathcal{L}}{\partial \mu_g} = \sum_v s_{g,v} w_{g,v} \Sigma_g^{-1} (x_v - \mu_g).$$

4.5. Gradient w.r.t. Gaussian Opacity

Since we have $w_{g,v} = \alpha_g \cdot \phi_{g,v}$, we can compute the derivative of $w_{g,v}$ with respect to Gaussian opacity α_g as:

$$\frac{\partial w_{g,v}}{\partial \alpha_g} = \phi_{g,v} = \frac{w_{g,v}}{\alpha_g}.$$

Therefore, the gradient of the loss \mathcal{L} with respect to α_g can be computed as:

$$\frac{\partial \mathcal{L}}{\partial \alpha_g} = \sum_v s_{g,v} \phi_{g,v} = \sum_v s_{g,v} \frac{w_{g,v}}{\alpha_g}.$$

4.6. Gradient w.r.t. Gaussian Covariance

By taking the derivative of density contribution with respect to Gaussian covariance Σ_g , we obtain:

$$\frac{\partial w_{g,v}}{\partial \Sigma_g} = \frac{1}{2} w_{g,v} \Sigma_g^{-1} d_{g,v} d_{g,v}^\top \Sigma_g^{-1}.$$

Thus, the gradient of the loss \mathcal{L} with respect to Σ_g can be formulated as:

$$\frac{\partial \mathcal{L}}{\partial \Sigma_g} = \frac{1}{2} \sum_v s_{g,v} w_{g,v} \Sigma_g^{-1} (x_v - \mu_g) (x_v - \mu_g)^\top \Sigma_g^{-1}.$$

5. Limitations and Future Work

Our current 2D image-level alignment and 3D LiDAR-based pseudo labeling approach is prone to occlusion, causing ShelfGaussian to perform suboptimally on flat surfaces (e.g., *drivable surface*) and small objects (e.g., *barrier*). These findings suggest future directions in exploring occlusion-aware alignment and more reliable, robust 3D supervision signals. Additionally, ShelfGaussian only supports learning 3D Gaussians on a per-view basis, which overlooks the interactions of Gaussians across different camera views. Future work could explore learning Gaussians directly in world space to eliminate view-specific restrictions. Lastly, a more extensive evaluation of our Gaussian-based trajectory planner in more challenging driving scenarios such as WOD-E2E [15] and on closed-loop benchmarks such as Bench2Drive [8] is an essential step toward real-world deployment.

References

- [1] Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4):433–459, 2010. 2
- [2] Luca Barsellotti, Lorenzo Bianchi, Nicola Messina, Fabio Carrara, Marcella Cornia, Lorenzo Baraldi, Fabrizio Falchi, and Rita Cucchiara. Talking to dino: Bridging self-supervised vision backbones with language for open-vocabulary segmentation. In *ICCV*, pages 22025–22035, 2025. 1, 2
- [3] Simon Boeder, Fabian Gigengack, and Benjamin Risse. Gaussianflowocc: Sparse and weakly supervised occupancy estimation using gaussian splatting and temporal flow. *arXiv preprint arXiv:2502.17288*, 2025. 3, 4
- [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 multimodal dataset for autonomous driving. In *CVPR*, pages 11621–11631, 2020. 3, 4
- [5] Clemente Donoso (CDonosok). ros2_camera_lidar_fusion: Ros2 package to calculate the intrinsic and extrinsic camera calibration and fuse camera & lidar. https://github.com/CDonosok/ros2_camera_lidar_fusion, 2025. 1
- [6] Yuanming Hu, Jiafeng Liu, Xuanda Yang, Mingkuan Xu, Ye Kuang, Weiwei Xu, Qiang Dai, William T Freeman, and Frédo Durand. Quantaichi: a compiler for quantized simulations. *ACM Transactions on Graphics (TOG)*, 40(4):1–16, 2021. 3, 4
- [7] Yuanhui Huang, Wenzhao Zheng, Yunpeng Zhang, Jie Zhou, and Jiwen Lu. Gaussianformer: Scene as gaussians for vision-based 3d semantic occupancy prediction. In *ECCV*, pages 376–393. Springer, 2024. 4
- [8] Xiaosong Jia, Zhenjie Yang, Qifeng Li, Zhiyuan Zhang, and Junchi Yan. Bench2drive: Towards multi-ability benchmarking of closed-loop end-to-end autonomous driving. In *NeurIPS 2024 Datasets and Benchmarks Track*, 2024. 5
- [9] Haoyi Jiang, Liu Liu, Tianheng Cheng, Xinjie Wang, Tianwei Lin, Zhizhong Su, Wenyu Liu, and Xinggang Wang. Gausstr: Foundation model-aligned gaussian transformer for self-supervised 3d spatial understanding. In *CVPR*, pages 11960–11970, 2025. 3, 4
- [10] Cijo Jose, Théo Moutakanni, Dahyun Kang, Federico Baldassarre, Timothée Darcet, Hu Xu, Daniel Li, Marc Szafraniec, Michaël Ramamonjisoa, Maxime Oquab, et al. Dinov2 meets text: A unified framework for image-and pixel-level vision-language alignment. In *CVPR*, pages 24905–24916, 2025. 1, 2
- [11] Steven Macenski, Tully Foote, Brian Gerkey, Chris Lalancette, and William Woodall. Robot operating system 2: Design, architecture, and uses in the wild. *Science robotics*, 7(66):eabm6074, 2022. 1
- [12] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 1

- [13] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021. [1](#)
- [14] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. [1](#), [2](#)
- [15] Runsheng Xu, Hubert Lin, Wonseok Jeon, Hao Feng, Yuliang Zou, Liting Sun, John Gorman, Kate Tolstaya, Sarah Tang, Brandyn White, et al. Wod-e2e: Waymo open dataset for end-to-end driving in challenging long-tail scenarios. *arXiv preprint arXiv:2510.26125*, 2025. [5](#)