

PseudoGD: Enhancing Spatial Reasoning in Vision-Language Models through Pseudo Geometric Knowledge Distillation

Anonymous ACL submission

Abstract

Recent Large Vision-Language Models (LVLMs) have shown remarkable success in general semantic understanding. However, they still struggle with 3D spatial reasoning tasks, such as estimating metric distances or understanding precise relative positions. Previous works, like SpatialVLM, tried to address this by using synthesized spatial VQA dataset. However, they are fundamentally limited because their vision encoders are biased toward 2D patterns learned from image-text pairs. In this paper, we argue that this lack of 3D awareness is a critical bottleneck that cannot be solved by data scaling alone. To address this, we propose Pseudo Geometric Distillation (PseudoGD), a framework designed to help vision encoders internalize 3D geometric information using only standard 2D images. PseudoGD explicitly injects metric scale and structural context into the encoder through a Joint Training strategy. This approach optimizes geometric learning and spatial VQA tasks together, ensuring that the Large Language Model (LLM) aligns well with the improved visual features in real-time. Extensive experiments on the OmniSpatial benchmark demonstrate that PseudoGD achieves enhanced performance across various model architectures. Notably, significant improvements in Hypothetical Perspective Taking and Locate tasks prove that our model has effectively learned a physical sense of space.

1 Introduction

The scope of Vision-Language Models (VLMs) now extends beyond basic tasks like image captioning and VQA to Embodied and Physical AI, where agents interact with the real world (Radford et al., 2021; Alayrac et al., 2022; Driess et al., 2023; Li et al., 2022; Goyal et al., 2017; Zitkovich et al., 2023; Hudson and Manning, 2019). Robust spatial understanding, such as accurate perception of the

locations, distances, and scales of objects, is essential for these systems (Chen et al., 2024; Wang et al., 2024; Yu et al., 2025). However, spatial reasoning remains a persistent bottleneck for current VLMs, limiting their effectiveness in tasks that require precise geometric and structural understanding (Lin et al., 2024; Liu et al., 2024; Kamath et al., 2023; Majumdar et al., 2024; Nikolov et al., 2025; Song et al., 2025).

To address this, SpatialVLM (Chen et al., 2024) improved spatial understanding by using large synthetic Spatial VQA datasets without explicit 3D training. Following this trajectory, subsequent studies have explored various avenues to increase spatial awareness. Based on this, OmniSpatial (Mengdi Jia et al., 2025) introduced a spatial reasoning benchmark grounded in cognitive psychology, VLM-3R (Fan et al., 2025) used 3D reconstruction, and SpatialRGPT (Cheng et al., 2024) improved spatial reasoning with region-based prompting.

Despite these advances, vision encoders of most VLMs are trained on 2D image-text pairs, limiting their ability to comprehend 3D spatial structures. To bridge this gap, we introduce Pseudo Geometric Knowledge Distillation (PseudoGD). This approach empowers vision encoders to internalize 3D geometric cues from monocular inputs by leveraging depth and segmentation models as "Teachers," effectively mitigating the inherent 2D bias. Unlike previous works (Huang et al., 2023; Li et al., 2025; Hong et al., 2023) that depend on 3D point clouds or multi-view data, our approach ensures high scalability and generalization performance through Pseudo knowledge distillation using only 2D images. Consequently, PseudoGD achieves a remarkable 59.6% accuracy, nearly doubling the performance of SpaceLLaVA (32.1%) and establishing a new baseline for comprehensive spatial understanding.

2 Related Works

Early VLM research primarily focused on visual grounding, establishing correspondences between visual objects and textual descriptions (Mao et al., 2016; Nagaraja et al., 2016; Yu et al., 2016). While datasets such as Visual Genome (Krishna et al., 2017) contributed to training 2D positional information at the bounding box level, they remained limited to planar recognition, excluding depth and 3D structural contexts. However, the advancement of Embodied AI and robotics has necessitated that VLMs possess 3D spatial understanding capabilities, such as metric distance estimation and relative spatial relations, beyond simple localization (Zhu et al., 2024; Sun et al., 2025).

The most dominant trend involves solutions through datasets synthesis. SpatialVLM (Chen et al., 2024) demonstrated that quantitative data expansion can enhance qualitative reasoning by constructing a massive VQA dataset that synthesizes 3D geometric information onto internet-scale 2D images. Similarly, RoboSpatial (Song et al., 2025) combined 3D scan data from robotics environments with 2D images to train the spatial awareness required for robotic manipulation. Meanwhile, SpatialRGPT (Cheng et al., 2024) and SR-3D (Cheng et al., 2025) were proposed to improve inference without modifying the model architecture. These methods enhance spatial reasoning performance by using region-based prompting to guide the model’s focus toward specific pixel areas. However, these studies share a common limitation: they rely on pre-trained encoders (e.g., CLIP (Radford et al., 2021)) fundamentally biased toward 2D semantic matching. Since adjustments at the text or prompt level do not fundamentally alter the encoder’s internal representations, reasoning in the absence of 3D structural information remains superficial (Hu et al., 2025; Qin et al., 2025).

Recently, model-centric approaches have emerged to geometrically tune vision encoders (Radford et al., 2021; Oquab et al., 2023; Dosovitskiy, 2020) themselves. VLM-3R (Fan et al., 2025) introduced an auxiliary module for 3D reconstruction from monocular video to assist visual perception, while 3D VLM-GD (Lee et al., 2025) proposed a geometric knowledge distillation method that extracts geometric cues from 3D foundation models and injects them into the vision encoder. Although 3D VLM-GD (Lee et al., 2025) aligns with our technical trajectory,

it faces a decisive constraint stemming from its strict dependency on fine-tuning with specific datasets paired with multi-view images or 3D point clouds. The construction of such datasets necessitates specialized capture equipment and strictly controlled environments. This dependency imposes a critical bottleneck on data scalability, fundamentally contradicting the philosophy of data abundance advocated by prior works like SpatialVLM. Consequently, it structurally precludes the utilization of vast web-scale 2D data, thereby isolating the model from the rich, diverse visual distributions required for universal spatial reasoning.

3 Methodology

We introduce PseudoGD, a framework designed to empower vision encoders to internalize 3D geometric reasoning directly from monocular 2D inputs. In this section, we define the fundamental cognitive bottlenecks inherent in existing VLM training paradigms and detail how our Pseudo Geometric Distillation and Joint Training strategies effectively bridge this gap.

3.1 Pseudo Geometric Knowledge Distillation (PseudoGD)

As illustrated in Figure 1, the core principle of this technique is to transfer the geometric reasoning capabilities of two teacher models, which are Depth Pro (Bochkovskii et al., 2024) and Segment Anything Model (SAM) (Kirillov et al., 2023; Ravi et al., 2024; Carion et al., 2025), to the vision encoder without requiring explicit 3D ground-truth data. By leveraging the knowledge distillation, this method enables the encoder to internalize spatial cues that are often absent in standard vision-language pre-training.

We integrate two complementary geometric properties to enrich visual representations. First, we employ Depth Pro as the Metric Depth Teacher to inject precise metric scale information. Unlike relative depth estimation, Depth Pro accounts for focal length and physical dimensions, providing the encoder with a physical sense of scale essential for quantitative reasoning (e.g., "5-meter distance"). Second, we utilize SAM as the Structural Segmentation Teacher to instill structural context. By encapsulating sophisticated object boundaries and part-whole relationships, SAM embeddings enhance the encoder’s perception of complex spa-

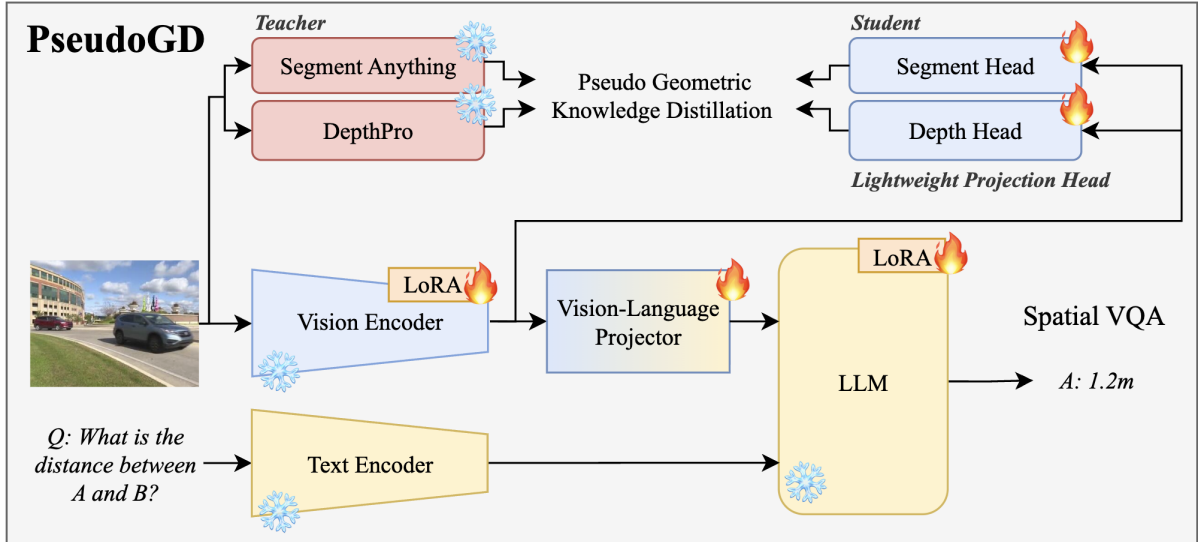


Figure 1: The training framework of PseudoGD. Our approach performs PseudoGD and Spatial VQA training simultaneously to enhance 3D spatial understanding.

tial arrangements, such as occlusions and relative positioning.

Through lightweight projection heads, we map the student encoder’s features into the teachers’ respective spaces. The distillation objective combines metric depth error and structural similarity loss, ensuring the model is not confined by the biases of specific 3D datasets. This facilitates the learning of universal geometric features, thereby securing a robust generalized capacity for comprehensive 3D spatial understanding.

3.2 Joint Training Strategy

Sequential training, where an encoder is pre-trained on geometric tasks before being connected to an LLM, often leads to catastrophic forgetting of semantic knowledge and feature misalignment, where the LLM fails to adapt to the shifted visual distribution. To circumvent these issues, we adopt a Joint Training strategy that co-optimizes geometric distillation and spatial VQA objectives within a unified loop. Formally, the total objective function is defined as Equation (1).

$$\mathcal{L}_{\text{Total}} = \mathcal{L}_{\text{VQA}} + \mathcal{L}_{\text{Depth}} + \mathcal{L}_{\text{Seg}} \quad (1)$$

Where \mathcal{L}_{VQA} denotes the autoregressive language modeling loss, while $\mathcal{L}_{\text{Depth}}$ and \mathcal{L}_{Seg} represent the distillation losses derived from the Metric Depth and Structural Segmentation teachers, respectively.

This integrated optimization process enables the vision encoder to acquire geometric inductive bi-

ases while simultaneously allowing the LLM to learn, in real-time, how to interpret these evolving representations. Consequently, PseudoGD achieves the capability to immediately leverage visual depth and structural cues for linguistic reasoning, ensuring a seamless alignment between visual perception and language generation.

4 Experiments

4.1 Experimental Setup

Datasets. For training, we utilize the VQASynth dataset collection for SpaceLLaVA, SpaceQwen, and SpaceThinker (Chen et al., 2024), which is publicly available on Hugging Face, to establish a baseline of fundamental spatial comprehension. In the evaluation phase, we adopt the OmniSpatial benchmark (Mengdi Jia et al., 2025) as the held-out test set to rigorously validate comprehensive spatial intelligence across diverse cognitive domains.

Evaluations. We conduct a comparative evaluation against four representative VLMs that demonstrate strengths in spatial perception or general multimodal capabilities, specifically SpaceLLaVA (Chen et al., 2024; Liu et al., 2023), SpaceMantis (Chen et al., 2024), SpaceQwen2.5-VL, and SpaceThinker-Qwen2.5 (Bai et al., 2023; Yang et al., 2025; Bai et al., 2025; Chen et al., 2024). These models represent the state-of-the-art in spatial interaction and general visual reasoning, thereby providing a rigorous standard.

To rigorously evaluate comprehensive spatial

Method	Overall	Dynamic Reasoning		Spatial Interaction			Complex Logic		Perspective Taking		
		Manipulate	Motion	Traffic	Locate	Geospatial Strategy	Pattern Recog.	Geometric Reasoning	Ego	Allo	Hypo.
SpaceLLaVA-13B	36.14	52.70	21.39	43.53	38.10	44.55	23.71	32.90	58.82	38.03	45.78
+ PseudoGD	38.88	56.76	30.35	50.59	41.90	47.27	28.87	25.81	53.92	39.10	45.78
SpaceMantis-8B	36.01	52.70	35.55	36.47	34.29	33.64	35.05	21.94	52.94	36.44	32.53
+ PseudoGD	37.18	54.05	33.24	42.35	35.24	36.36	34.02	21.94	51.96	38.83	43.37
SpaceQwen2.5VL-3B†	40.25	58.11	39.88	41.18	40.95	40.91	29.90	25.81	63.73	38.83	39.76
+ PseudoGD	39.73	50.00	42.20	48.24	43.81	37.27	28.87	27.10	50.00	36.97	45.78
SpaceThinker-Qwen2.5†	40.42	47.84	53.06	43.29	35.43	38.73	24.33	28.00	58.04	35.11	31.08
+ PseudoGD	42.20	55.41	49.42	49.41	40.95	39.09	26.80	26.45	67.65	35.64	44.58

Table 1: OmniSpatial benchmark results (%) on task-level evaluation. Results marked with † are cited from OmniSpatial (Mengdi Jia et al., 2025).

reasoning capabilities, we utilize the OmniSpatial evaluation set (Mengdi Jia et al., 2025) as our primary benchmark. We measure model performance across the four core dimensions defined by the benchmark: Dynamic Reasoning, Complex Spatial Logic, Spatial Interaction, and Perspective-Taking.

Implementation Details. To ensure a rigorous comparative analysis, we trained all models using the identical base architectures and datasets as their respective baselines, with the inclusion of PseudoGD being the sole experimental variable. All models were fine-tuned using LoRA (Hu et al., 2022). To balance reasoning capacity with the preservation of visual priors, we set the LoRA rank to 128 for the LLM and 4 for the vision encoder, while fully fine-tuning the multimodal projector to ensure robust cross-modal alignment. Additionally, the depth estimation and segmentation heads were implemented as lightweight modules integrated directly into the vision encoder, minimizing architectural complexity while facilitating geometric internalization. Further details are provided in the Appendix.

4.2 Experimental Results and Analysis

Superior Performance and Generalization. As shown in Table 1, PseudoGD consistently enhances spatial reasoning across diverse architectures, bridging the gap between 2D perception and 3D cognition. Notably, SpaceThinker-Qwen2.5 + PseudoGD achieved the best overall accuracy of 42.20% (+1.78%), with significant gains in Locate and Traffic Analysis. This confirms that internalizing metric depth empowers models to perform precise localization and dynamic reasoning.

Bridging Cognitive Bottlenecks. The most profound impact is observed in Hypothetical Perspective Taking, where SpaceThinker surged from 31.08% to 44.58%. This dramatic gain suggests the vision encoder has successfully internalized 3D structural information, overcoming the cognitive bottleneck of simulating unseen viewpoints without explicit 3D priors. Additionally, the universal improvement in the Locate metric proves that PseudoGD equips models with a physical sense of space, enabling reliable spatial grounding even in complex environments where semantic features alone are insufficient.

5 Conclusions

In this work, we addressed the limitations of current VLMs in 3D spatial reasoning, which stem from their reliance on 2D semantic priors. To mitigate this, we proposed PseudoGD, a framework integrating Pseudo Geometric Knowledge Distillation with a Joint Training strategy. This approach enables vision encoders to learn 3D geometric representations from monocular 2D inputs, effectively utilizing depth and segmentation cues without requiring explicit 3D training data. Our evaluation on the OmniSpatial benchmark demonstrates that this method consistently improves spatial reasoning capabilities across diverse architectures. The performance gains observed in Locate and Hypothetical Perspective Taking indicate that the model has effectively internalized physical scale and structural relationships. These findings suggest that explicitly distilling geometric features is a valid approach for enhancing the spatial understanding of VLMs.

Limitations

Although this study demonstrates that PseudoGD effectively enhances the spatial reasoning capabilities of VLMs by distilling geometric knowledge, several limitations remain.

First, the performance of our framework is fundamentally dependent on the quality of the teacher signals. Since we rely on Depth Pro (Bochkovskii et al., 2024) and SAM (Kirillov et al., 2023) to generate pseudo-labels for metric depth and structural segmentation, any errors or artifacts produced by these models inevitably propagate to the student encoder. Consequently, the model may exhibit performance degradation in scenarios where the teacher models struggle, such as scenes containing mirrors, transparent materials, or extreme lighting conditions that cause visual ambiguity.

Second, there is a limitation regarding the approximation of 3D geometry. While our method successfully empowers the vision encoder to internalize 3D cues from monocular 2D images, this remains an implicit approximation rather than an explicit measurement. Compared to approaches utilizing ground-truth 3D point clouds or multi-view geometry, our model may lack precision in fine-grained metric estimation for complex, heavily occluded structures. Future work is needed to bridge the gap between internalized 2D spatial cues and absolute 3D physical accuracy.

Third, the current framework is constrained to static single-image inference. Although the model showed improvements in the Dynamic Reasoning track of OmniSpatial, it infers motion and temporal relationships based solely on static visual evidence. Practical robotic applications often require continuous reasoning over temporal sequences to handle dynamic physical interactions. Extending the PseudoGD mechanism to video-based VLMs to capture temporal context remains a critical direction for future research.

Finally, the computational overhead during training is non-negligible. Unlike standard instruction tuning that relies on sparse token prediction, our joint training strategy involves aligning dense, pixel-level features from projection heads with teacher embeddings. This increases memory consumption and computational cost during the training phase, presenting a challenge for scalability when applying this method to ultra-large-scale multimodal models.

References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Aleksei Bochkovskii, Amaël Delaunoy, Hugo Germain, Marcel Santos, Yichao Zhou, Stephan R Richter, and Vladlen Koltun. 2024. Depth pro: Sharp monocular metric depth in less than a second. *arXiv preprint arXiv:2410.02073*.
- Nicolas Carion, Laura Gustafson, Yuan-Ting Hu, Shoubhik Debnath, Ronghang Hu, Didac Suris, Chaitanya Ryali, Kalyan Vasudev Alwala, Haitham Khedr, Andrew Huang, et al. 2025. Sam 3: Segment anything with concepts. *arXiv preprint arXiv:2511.16719*.
- Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 2024. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465.
- An-Chieh Cheng, Yang Fu, Yukang Chen, Zhijian Liu, Xiaolong Li, Subhashree Radhakrishnan, Song Han, Yao Lu, Jan Kautz, Pavlo Molchanov, et al. 2025. 3d aware region prompted vision language model. *arXiv preprint arXiv:2509.13317*.
- An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. 2024. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093.
- Alexey Dosovitskiy. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Ayzan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, et al. 2023. Palm-e: An embodied multimodal language model.
- Zhiwen Fan, Jian Zhang, Renjie Li, Junge Zhang, Runjin Chen, Hezhen Hu, Kevin Wang, Huaizhi Qu, Dilin Wang, Zhicheng Yan, et al. 2025. Vlm-3r: Vision-language models augmented with instruction-aligned 3d reconstruction. *arXiv preprint arXiv:2505.20279*.

414	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv	Junnan Li, Dongxu Li, Caiming Xiong, and Steven	471
415	Batra, and Devi Parikh. 2017. Making the v in vqa	Hoi. 2022. Blip: Bootstrapping language-image pre-	472
416	matter: Elevating the role of image understanding	training for unified vision-language understanding	473
417	in visual question answering. In <i>Proceedings of the</i>	and generation. In <i>International conference on ma-</i>	474
418	<i>IEEE conference on computer vision and pattern</i>	<i>chine learning</i> , pages 12888–12900. PMLR.	475
419	<i>recognition</i> , pages 6904–6913.		
420	Yining Hong, Chunru Lin, Yilun Du, Zhenfang Chen,	Ji Lin, Hongxu Yin, Wei Ping, Pavlo Molchanov, Mo-	476
421	Joshua B Tenenbaum, and Chuang Gan. 2023. 3d	hammad Shoeybi, and Song Han. 2024. Vila: On	477
422	concept learning and reasoning from multi-view im-	pre-training for visual language models. In <i>Proceed-</i>	478
423	ages. In <i>Proceedings of the IEEE/CVF Conference</i>	<i>ings of the IEEE/CVF conference on computer vision</i>	479
424	<i>on Computer Vision and Pattern Recognition</i> , pages	<i>and pattern recognition</i> , pages 26689–26699.	480
425	9202–9212.		
426	Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae	481
427	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,	Lee. 2024. Improved baselines with visual instruc-	482
428	Weizhu Chen, et al. 2022. Lora: Low-rank adap-	tuning. In <i>Proceedings of the IEEE/CVF con-</i>	483
429	tation of large language models. <i>ICLR</i> , 1(2):3.	<i>ference on computer vision and pattern recognition</i> ,	484
		pages 26296–26306.	485
430	Wenbo Hu, Jingli Lin, Yilin Long, Yunlong Ran, Lihan	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae	486
431	Jiang, Yifan Wang, Chenming Zhu, Runsen Xu, Tai	Lee. 2023. Visual instruction tuning. <i>Advances in</i>	487
432	Wang, and Jiangmiao Pang. 2025. G ₃ vlm: Geome-	<i>neural information processing systems</i> , 36:34892–	488
433	try grounded vision language model with unified 3d	34916.	489
434	reconstruction and spatial reasoning. <i>arXiv preprint</i>	Arjun Majumdar, Anurag Ajay, Xiaohan Zhang, Pranav	490
435	<i>arXiv:2511.21688</i> .	Putta, Sriram Yenamandra, Mikael Henaff, Sneha	491
436	Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun	Silwal, Paul Mccvay, Oleksandr Maksymets, Sergio	492
437	Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun	Arnaud, et al. 2024. Openeqa: Embodied question	493
438	Zhu, Baoxiong Jia, and Siyuan Huang. 2023. An em-	answering in the era of foundation models. In <i>Pro-</i>	494
439	bodyed generalist agent in 3d world. <i>arXiv preprint</i>	<i>ceedings of the IEEE/CVF conference on computer</i>	495
440	<i>arXiv:2311.12871</i> .	<i>vision and pattern recognition</i> , pages 16488–16498.	496
441	Drew A Hudson and Christopher D Manning. 2019.	Junhua Mao, Jonathan Huang, Alexander Toshev, Oana	497
442	Gqa: A new dataset for real-world visual reasoning	Camburu, Alan L Yuille, and Kevin Murphy. 2016.	498
443	and compositional question answering. In <i>Proceed-</i>	Generation and comprehension of unambiguous ob-	499
444	<i>ings of the IEEE/CVF conference on computer vision</i>	ject descriptions. In <i>Proceedings of the IEEE con-</i>	500
445	<i>and pattern recognition</i> , pages 6700–6709.	<i>ference on computer vision and pattern recognition</i> ,	501
446	Amita Kamath, Jack Hessel, and Kai-Wei Chang. 2023.	pages 11–20.	502
447	What’s “up” with vision-language models? investi-	Shaochen Zhang Wenyao Zhang Xinqiang Yu Jiawei	503
448	gating their struggle with spatial reasoning. <i>arXiv</i>	He He Wang Li Yi Mengdi Jia, Zekun Qi et al. 2025.	504
449	<i>preprint arXiv:2310.19785</i> .	Omnispacial: Towards comprehensive spatial reason-	505
450	Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi	ing benchmark for vision language models . In <i>Sub-</i>	506
451	Mao, Chloe Rolland, Laura Gustafson, Tete Xiao,	<i>mitted to The Fourteenth International Conference</i>	507
452	Spencer Whitehead, Alexander C Berg, Wan-Yen Lo,	<i>on Learning Representations</i> . Under review.	508
453	et al. 2023. Segment anything. In <i>Proceedings of</i>	Varun K Nagaraja, Vlad I Morariu, and Larry S Davis.	509
454	<i>the IEEE/CVF international conference on computer</i>	2016. Modeling context between objects for refer-	510
455	<i>vision</i> , pages 4015–4026.	ring expression understanding. In <i>European Confer-</i>	511
456	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin John-	<i>ence on Computer Vision</i> , pages 792–807. Springer.	512
457	son, Kenji Hata, Joshua Kravitz, Stephanie Chen,	Nikolay Nikolov, Giuliano Albanese, Sombit Dey, Alek-	513
458	Yannis Kalantidis, Li-Jia Li, David A Shamma, et al.	sandar Yanev, Luc Van Gool, Jan-Nico Zaeck, and	514
459	2017. Visual genome: Connecting language and vi-	Danda Pani Paudel. 2025. Spear-1: Scaling beyond	515
460	sion using crowdsourced dense image annotations.	robot demonstrations via 3d understanding. <i>arXiv</i>	516
461	<i>International journal of computer vision</i> , 123(1):32–	<i>preprint arXiv:2511.17411</i> .	517
462	73.		
463	Seonho Lee, Jiho Choi, Inha Kang, Jiwook Kim, Jun-	Maxime Oquab, Timothée Darcet, Théo Moutakanni,	518
464	sung Park, and Hyunjung Shim. 2025. 3d-aware	Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fer-	519
465	vision-language models fine-tuning with geometric	nandez, Daniel Haziza, Francisco Massa, Alaaeldin	520
466	distillation. <i>arXiv preprint arXiv:2506.09883</i> .	El-Nouby, et al. 2023. Dinov2: Learning robust vi-	521
467	Chengmeng Li, Junjie Wen, Yan Peng, Yaxin Peng,	sual features without supervision. <i>arXiv preprint</i>	522
468	Feifei Feng, and Yichen Zhu. 2025. Pointvla: Inject-	<i>arXiv:2304.07193</i> .	523
469	ing the 3d world into vision-language-action models.	Yiming Qin, Bomim Wei, Jiaxin Ge, Konstantinos	524
470	<i>arXiv preprint arXiv:2503.07511</i> .	Kallidromitis, Stephanie Fu, Trevor Darrell, and	525
		Xudong Wang. 2025. Chain-of-visual-thought:	526

527 Teaching vlms to see and think better with continuous
528 visual tokens. *arXiv preprint arXiv:2511.19418*.

529 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
530 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-
531 try, Amanda Askell, Pamela Mishkin, Jack Clark,
532 et al. 2021. Learning transferable visual models from
533 natural language supervision. In *International confer-
534 ence on machine learning*, pages 8748–8763. PmLR.

535 Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Rong-
536 hang Hu, Chaitanya Ryalı, Tengyu Ma, Haitham
537 Khedr, Roman Rädle, Chloe Rolland, Laura
538 Gustafson, et al. 2024. Sam 2: Segment any-
539 thing in images and videos. *arXiv preprint
540 arXiv:2408.00714*.

541 Chan Hee Song, Valts Blukis, Jonathan Tremblay,
542 Stephen Tyree, Yu Su, and Stan Birchfield. 2025. Ro-
543 bospatial: Teaching spatial understanding to 2d and
544 3d vision-language models for robotics. In *Proceed-
545 ings of the Computer Vision and Pattern Recognition
546 Conference*, pages 15768–15780.

547 Lin Sun, Bin Xie, Yingfei Liu, Hao Shi, Tiancai Wang,
548 and Jiale Cao. 2025. Geovla: Empowering 3d repre-
549 sentations in vision-language-action models. *arXiv
550 preprint arXiv:2508.09071*.

551 Tai Wang, Xiaohan Mao, Chenming Zhu, Runsen Xu,
552 Ruiyuan Lyu, Peisen Li, Xiao Chen, Wenwei Zhang,
553 Kai Chen, Tianfan Xue, et al. 2024. Embodiedscan:
554 A holistic multi-modal 3d perception suite towards
555 embodied ai. In *Proceedings of the IEEE/CVF Con-
556 ference on Computer Vision and Pattern Recognition*,
557 pages 19757–19767.

558 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,
559 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,
560 Chengen Huang, Chenxu Lv, et al. 2025. Qwen3
561 technical report. *arXiv preprint arXiv:2505.09388*.

562 Licheng Yu, Patrick Poirson, Shan Yang, Alexander C
563 Berg, and Tamara L Berg. 2016. Modeling context
564 in referring expressions. In *European conference on
565 computer vision*, pages 69–85. Springer.

566 Songsong Yu, Yuxin Chen, Hao Ju, Lianjie Jia, Fuxi
567 Zhang, Shaofei Huang, Yuhan Wu, Rundi Cui, Bing-
568 hao Ran, Zaibin Zhang, et al. 2025. How far are vlms
569 from visual spatial intelligence? a benchmark-driven
570 perspective. *arXiv preprint arXiv:2509.18905*.

571 Haoyi Zhu, Honghui Yang, Yating Wang, Jiange Yang,
572 Limin Wang, and Tong He. 2024. Spa: 3d spatial-
573 awareness enables effective embodied representation.
574 *arXiv preprint arXiv:2410.08208*.

575 Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu,
576 Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan
577 Welker, Ayzaan Wahid, et al. 2023. Rt-2: Vision-
578 language-action models transfer web knowledge to
579 robotic control. In *Conference on Robot Learning*,
580 pages 2165–2183. PMLR.

Appendix A Implementation Details 581

A.1 Common Configuration 582

LoRA Configuration. We apply Low-Rank Adaptation (LoRA) (Hu et al., 2022) to both the language model and the vision encoder, while fully fine-tuning the multimodal projector. All LoRA modules use a dropout rate of 0.05. The detailed configuration is summarized in Table 2. 583 584 585 586 587 588

Component	Rank	Alpha	LR
LLM	128	256	2×10^{-5}
Vision Encoder	4	8	1×10^{-5}
MM Projector	Full	–	2×10^{-5}

Table 2: LoRA configuration for all experiments. LLM LoRA targets all attention projections and feed-forward layers. Vision encoder LoRA targets the fused QKV projection and output projection in each of the 32 transformer blocks.

Geometric Distillation Heads. To facilitate 589 geometric knowledge distillation, we attach 590 lightweight prediction heads directly to the vision 591 encoder’s output features. The depth prediction 592 head is implemented as a Multi-Layer Perceptron 593 (MLP) that linearly projects the input features to a 594 1024-dimensional hidden layer with GELU activa- 595 tion, followed by a final linear projection to a scalar 596 output. To enforce physically meaningful metric 597 constraints, we apply a scaled activation function, 598 defined as $\text{Softplus}(x) \times 7.0 + 0.1$, ensuring posi- 599 tive depth values within the valid range of $[0.1, \infty)$ 600 meters. The resulting predictions are interpolated 601 to match the teacher’s 24×24 resolution. 602

The segmentation head is composed of a fea- 603 ture projector and a spatial upsampling module 604 designed to reconstruct high-resolution structural 605 details. The projector utilizes an MLP ($\rightarrow 2048 \rightarrow$ 606 256) with GELU activation to compress semantic 607 features into a structural embedding space. Subse- 608 quently, these features are spatially reshaped and 609 processed by Transposed Convolution layers fol- 610 lowed by a GELU activation and a standard Convo- 611 lution layer (kernel size 3, padding 1), ultimately 612 yielding 64×64 feature maps. Both heads are opti- 613 mized with a learning rate of 1×10^{-3} , signifi- 614 cantly higher than the base model parameters, to facilitate 615 the rapid adaptation of these randomly initialized 616 components. 617

618 A.2 SpaceQwen2.5 w/ PseudoGD 619 Configuration

620 **Base Model and Dataset.** Space-
621 Qwen2.5 w/ PseudoGD is built upon
622 Qwen/Qwen2.5-VL-3B-Instruct. The model
623 is fine-tuned using the remyxai/OpenSpaces
624 dataset, a spatial reasoning benchmark containing
625 approximately 10K question-answering samples.
626 Each sample consists of an RGB image paired
627 with a natural language question about spatial
628 relationships (e.g., relative positions, distances,
629 orientations) and a corresponding answer. The
630 dataset covers diverse indoor and outdoor scenes
631 with varying complexity levels. The dataset is
632 partitioned into 9.26K training, and 1.03K test
633 samples. All experiments are conducted using
634 bfloat16 mixed-precision training for memory
635 efficiency while maintaining numerical stability.

636 **Training Hyperparameters.** Training employs
637 a batch size of 4 with gradient accumulation steps
638 of 8, yielding an effective batch size of 32. The
639 maximum sequence length is set to 2048 tokens to
640 accommodate both visual tokens (variable length
641 due to dynamic resolution) and text tokens. We em-
642 ploy the AdamW optimizer with no weight decay.
643 The learning rate follows a cosine annealing sched-
644 ule with 3% linear warmup. Training proceeds for
645 approximately 3 epochs over the full dataset.

646 **Input Processing.** Images are processed at their
647 native resolution by Qwen2.5-VL’s dynamic reso-
648 lution mechanism, preserving fine-grained spatial
649 details. For teacher model inference, images are
650 center-cropped and resized to 224×224 pixels,
651 then normalized using ImageNet statistics ($\mu =$
652 $[0.485, 0.456, 0.406]$, $\sigma = [0.229, 0.224, 0.225]$)
653 to ensure compatibility with the pre-trained teacher
654 models.

655 A.3 SpaceThinker w/ PseudoGD 656 Configuration

657 **Base Model and Dataset.** Space-
658 Thinker w/ PseudoGD is based on
659 UCSC-VLAA/VLAA-Thinker-Qwen2.5VL-3B,
660 a variant of Qwen2.5-VL fine-tuned for explicit
661 chain-of-thought reasoning. Training is per-
662 formed on the remyxai/SpaceThinker dataset,
663 which contains approximately 12K samples with
664 structured reasoning annotations. The dataset is
665 partitioned into 11.4K training, and 1.25K test
666 samples. The model follows a two-stage output

format with explicit reasoning:

```
667 <think>[step-by-step reasoning]</think> 668  
669 <answer>[final answer]</answer> 670
```

671 This format encourages the model to externalize
672 its spatial reasoning process before providing an-
673 swers. The system prompt instructs: “*You should*
674 *first think about the reasoning process and then*
675 *provide the answer. Use <think>...</think> and*
676 *<answer>...</answer> tags.” All experiments use*
677 bfloat16 precision.

678 **Training Hyperparameters.** Due to the longer
679 output sequences required for explicit reasoning,
680 SpaceThinker is trained with a reduced per-GPU
681 batch size of 1 and gradient accumulation over 8
682 steps, yielding an effective batch size of 8. The
683 maximum sequence length remains 2048 tokens
684 with a warmup ratio of 3%. The model is trained
685 for 3 full epochs over the dataset.

686 A.4 SpaceLLaVA-13B w/ PseudoGD 687 Configuration

688 **Base Model and Dataset.** SpaceLLaVA-13B w/
689 PseudoGD is constructed upon the LLaVA-v1.5-
690 13B architecture, utilizing CLIP ViT-L/14 as the vi-
691 sion encoder which produces 1024-dimensional vi-
692 sual features. To validate the generalizability of our
693 geometric distillation approach, the model is fine-
694 tuned on the remyxai/vqasynth_spacellava
695 dataset, comprising approximately 28K synthetic
696 spatial QA samples. The dataset is partitioned into
697 25.2K training, and 2.8K test samples.

698 **Training Hyperparameters.** Training employs a
699 batch size of 4 with gradient accumulation steps of
700 8, yielding an effective batch size of 32. The model
701 is trained for 1 epoch using AdamW optimizer and
702 cosine annealing schedule. Mixed precision train-
703 ing (bfloat16) is employed for memory efficiency.

704 A.5 SpaceMantis w/ PseudoGD Configuration

705 **Base Model and Architecture.** SpaceMantis w/
706 PseudoGD extends the Mantis-8B-siglip-llama3
707 architecture, which integrates a SigLIP vision
708 encoder with the Llama-3-8B language model.
709 This configuration serves to evaluate the efficacy
710 of geometric distillation on a multi-image capa-
711 ble VLM framework underpinned by a distinct
712 vision backbone. The model is fine-tuned on
713 the remyxai/vqasynth_spacellava dataset. The
714 dataset is partitioned into 25.2K training, and 2.8K
715 test samples.

Training Hyperparameters. The model is trained with a batch size of 4 and gradient accumulation steps of 8, resulting in an effective batch size of 32. The optimizer is AdamW and cosine annealing schedule is applied. Mixed precision training (bfloat16) is employed for memory efficiency.

A.6 Loss Function

The total training objective combines three complementary loss terms as defined in Equation (2):

$$\mathcal{L}_{\text{total}} = \lambda_{\text{vqa}} \mathcal{L}_{\text{VQA}} + \lambda_{\text{depth}} \mathcal{L}_{\text{depth}} + \lambda_{\text{seg}} \mathcal{L}_{\text{seg}} \quad (2)$$

where \mathcal{L}_{VQA} is the standard cross-entropy loss for next-token prediction, and $\mathcal{L}_{\text{depth}}$, \mathcal{L}_{seg} are the geometric distillation losses described below. All balancing coefficients are set to $\lambda_{\text{vqa}} = \lambda_{\text{depth}} = \lambda_{\text{seg}} = 1.0$ by default, ensuring equal importance between semantic reasoning and geometric internalization.

Depth Distillation Loss. We adopt a scale-invariant logarithmic loss, which is robust to global scale ambiguities and focuses on relative depth relationships:

$$\mathcal{L}_{\text{depth}} = \frac{1}{n} \sum_i d_i^2 - \frac{1}{2n^2} \left(\sum_i d_i \right)^2, \quad (3)$$

$$d_i = \log(\hat{y}_i) - \log(y_i)$$

Here, \hat{y}_i and y_i denote the predicted and teacher depth values at spatial location i , respectively, and n is the total number of spatial locations. The second term penalizes systematic scale shifts, encouraging the model to learn relative depth orderings even when absolute scale information is noisy. Depth values are clamped to a minimum of 0.1 meters before taking the logarithm to ensure numerical stability.

Segmentation Distillation Loss. The segmentation distillation loss measures the alignment between predicted and teacher feature embeddings using cosine similarity:

$$\mathcal{L}_{\text{seg}} = 1 - \frac{1}{HW} \sum_{h,w} \cos(\hat{\mathbf{f}}_{h,w}, \mathbf{f}_{h,w}^{\text{teacher}}) \quad (4)$$

where $\hat{\mathbf{f}}_{h,w}, \mathbf{f}_{h,w}^{\text{teacher}} \in \mathbb{R}^{256}$ are the L_2 -normalized feature vectors at spatial position (h, w) . This loss encourages the student to learn semantically meaningful spatial representations that capture object boundaries and structural patterns, without requiring explicit segmentation labels.

A.7 Teacher Models

To empower the vision encoder with robust geometric inductive biases, we employ two complementary foundation models as teachers. First, for Metric Depth Distillation, we utilize Depth Pro (Bochkovskii et al., 2024), developed by Apple. Unlike conventional relative depth estimators, Depth Pro predicts absolute metric depth (in meters), enabling the student model to internalize a physical sense of scale. For compatibility with our projection architecture, the extracted depth maps are interpolated to a 24×24 spatial resolution. Second, for Structural Segmentation Distillation, we adopt the Segment Anything Model (SAM) (Kirillov et al., 2023) with a ViT-Base backbone. SAM provides semantic-agnostic structural embeddings that encode precise object boundaries. These 256-dimensional embeddings are similarly interpolated to a 64×64 resolution, ensuring that the encoder learns to capture fine-grained structural details within a standardized feature space.

A.8 Hardware and Software

All experiments are conducted on a cluster of $5 \times$ NVIDIA A100 80GB PCIe GPUs. Multi-GPU training is orchestrated using the Hugging Face Accelerate library with distributed data parallelism. The framework is implemented in PyTorch 2.0+, leveraging the Transformers library for model loading and the PEFT library for parameter-efficient fine-tuning.

A.9 Ablation Study on PseudoGD Weight.

Table 3 presents an ablation study on the loss weighting factor λ used for PseudoGD when applied to the SpaceMantis-8B model. The loss weight λ controls the relative contribution of geometric distillation during training, allowing us to analyze how strongly enforcing geometric supervision affects different categories of spatial reasoning.

Overall, increasing the PseudoGD weight (λ) consistently improves the aggregate performance, with the best overall accuracy achieved at $\lambda = 1.0$. This trend suggests that stronger geometric distillation provides more stable and globally beneficial supervision for spatial reasoning tasks. In particular, tasks under *Spatial Interaction* and *Complex Logic*, such as *Traffic*, *Geospatial Strategy*, and *Pattern Recognition*, show noticeable gains as the weight increases, indicating that these tasks benefit

Method	Overall	Dynamic Reasoning		Spatial Interaction			Complex Logic		Perspective Taking		
		Manipulate	Motion	Traffic	Locate	Geospatial Strategy	Pattern Recog.	Geometric Reasoning	Ego	Allo	Hypo.
SpaceMantis-8B	36.01	52.70	35.55	36.47	34.29	33.64	35.05	21.94	52.94	36.44	32.53
+ PseudoGD ($\lambda=0.1$)	36.73	55.41	32.95	36.47	35.24	33.64	28.87	23.87	56.86	40.43	33.73
+ PseudoGD ($\lambda=0.5$)	36.96	55.41	33.53	38.82	35.24	35.45	31.96	23.23	54.90	39.36	38.55
+ PseudoGD ($\lambda=1.0$)	37.18	54.05	33.24	42.35	35.24	36.36	34.02	21.94	51.96	38.83	43.37

Table 3: OmniSpatial benchmark results (%) for SpaceMantis variants with different PseudoGD lambdas. λ contains only λ_{depth} and λ_{seg} , and the same hyperparameter values were set for both lambdas. λ_{vqa} was always fixed at 1.

807 from explicit geometric constraints and relational
808 structure.

809 In contrast, lower PseudoGD weights (e.g., $\lambda =$
810 0.1) tend to favor *Perspective Taking* subtasks, espe-
811 cially *Egocentric* and *Allocentric* reasoning. This
812 behavior suggests that weaker geometric regular-
813 ization allows the model to rely more on implicit
814 visual cues and viewpoint-dependent reasoning,
815 rather than enforcing rigid geometric consistency.
816 However, this comes at the cost of degraded per-
817 formance on geometry-intensive subtasks, such as
818 *Pattern Recognition*.

819 The intermediate setting ($\lambda = 0.5$) exhibits a bal-
820 anced trade-off between these two regimes, achiev-
821 ing moderate improvements across most task cate-
822 gories without strongly biasing the model toward
823 either perspective-centric or geometry-centric rea-
824 soning. These results indicate that the PseudoGD
825 weight serves as an effective control knob for bal-
826 ancing geometric consistency and viewpoint flexi-
827 bility in multi-modal spatial reasoning models.