

EVO-0: VISION-LANGUAGE-ACTION MODEL WITH IMPLICIT SPATIAL UNDERSTANDING

Tao Lin*, Gen Li*, Yilei Zhong, Yanwen Zou, Bo Zhao†

¹School of AI, Shanghai Jiao Tong University, ²EvoMind Tech, ³IAAR-Shanghai
 taolin200108@gmail.com, bo.zhao@sjtu.edu.cn
<https://github.com/MINT-SJTU/Evo-VLA>

ABSTRACT

Vision-Language-Action (VLA) models have emerged as a promising framework for enabling generalist robots capable of perceiving, reasoning, and acting in the real world. These models usually build upon pretrained Vision-Language Models (VLMs), which excel at semantic understanding due to large-scale text pretraining. However, VLMs typically lack precise spatial understanding capabilities, as they are primarily tuned on 2D image-text pairs without 3D supervision. To address this limitation, recent approaches have incorporated explicit 3D inputs such as point clouds or depth maps, but this necessitates additional depth sensors or defective estimation. In contrast, our work introduces a plug-and-play module that implicitly injects 3D geometry features into VLA models by leveraging an off-the-shelf visual geometry foundation models. We design five spatially challenging tasks that require precise spatial understanding ability to validate effectiveness of our method. Extensive evaluations show that our method significantly improves the performance of state-of-the-art VLA models across diverse scenarios.

1 INTRODUCTION

Vision-Language-Action (VLA) models have recently attracted substantial attention and achieved notable progress. These models typically fine-tune pre-trained Vision-Language Models (VLMs) using robot manipulation data to leverage the vision-language generalization ability learned from large-scale image-text data. This paradigm has achieved impressive results success a wide range of real-world and simulated tasks.

However, existing VLA models exhibit a fatal limitation, which is the lack of precise 3D spatial understanding. This shortcoming can be largely attributed to two main factors: (1) the pre-training data and objectives of VLMs are primarily based on 2D image-text alignments, and (2) the robotic datasets used for fine-tuning are typically small-scale and contain only RGB observations, without explicit 3D spatial information. As a result, these models often struggle to capture the precise geometric and spatial relationships that are essential for effective interaction in the physical world. Recent studies Cai et al. (2024); Daxberger et al. (2025) have empirically validated this observation, showing that VLMs tend to generalize poorly when it comes to interpret 3D structures from visual inputs alone. This presents a critical bottleneck in scaling VLA models to more complex tasks and physically grounded scenarios.

To address this limitation, emerging approaches Li et al. (2025a); Qu et al. (2025); Zhen et al. (2024); Li et al. (2025b) have tied to incorporate 3D information into VLA models to enhance their spatial understanding capabilities. A common strategy is to explicitly inject depth information into the learning pipeline, such as point clouds or depth maps, either captured by depth sensors or predicted by depth estimation networks Cai et al. (2024); Bhat et al. (2023); Yang et al. (2024). While effective to some extent, these methods introduce new challenges. They often require additional depth sensors, which may not be available in practical settings. Moreover, the defective depth estimations probably introduce extra noise, affecting the reliability of the learned 3D representations.

*Equal contribution †Corresponding author

In this paper, we introduce *Evo-0*, a novel VLA architecture that explores an alternative strategy to enhance spatial understanding of VLA models in an implicit manner. Specifically, we leverage the powerful 3D perception ability of Visual Geometry Grounded Transformer (VGGT) Wang et al. (2025), which is trained on large-scale 2D-3D paired datasets. These 3D features can be obtained from the vanilla video inputs of robotic data with VGGT, which can be used to complement VLMs without relying on explicit depth inputs or estimation. To this end, we design a lightweight fusion module that integrates geometry-grounded features from VGGT with visual tokens in VLM, enabling the model to perceive object layouts and reason about spatial relations more effectively. We demonstrate the effectiveness of our method through comprehensive experiments on five spatially challenging real-world tasks, where our model consistently improves spatial understanding and outperforms the state-of-the-art VLA models.

In summary, we make two main contributions in this work: (1) We propose a plug-and-play module to enhance the spatial understanding of VLA models by implicitly injecting 3D geometric priors from the Visual Geometry Grounded Transformer (VGGT), and (2) We design and evaluate our method on five diverse and spatially challenging tasks, demonstrating consistent improvements over the strong baselines.

2 RELATED WORK

Vision-Language-Action Models. Recently, several studies Kim et al. (2024); Black et al. (2024); Bjorck et al. (2025); Li et al. (2025b); Liu et al. (2024); Brohan et al. (2023) have focused on building general-purpose robot policies by extending pre-trained vision-language models (VLMs) with action prediction capabilities. These models, known as vision-language-action (VLA) models, demonstrate strong performance and few-shot generalization across a wide range of embodied tasks.

Among them, OpenVLA Kim et al. (2024) is trained on 970k multi-robot demonstrations from the Open-X Embodiment O’Neill et al. (2024) dataset, demonstrates strong generalization across a wide range of tasks and embodiments, and supports efficient fine-tuning under limited computational resources. π_0 Black et al. (2024) adapts the PaliGemma Beyer et al. (2024) architecture for robotic control and introduces a flow-matching-based Lipman et al. (2022); Liu (2022) action expert module that enables accurate prediction of continuous actions. GR00T Bjorck et al. (2025) introduces an effective co-training strategy that jointly leverages web data, synthetic data, and real-world robot data within a unified framework, enabling broad generalization across tasks and embodiments.

Despite the promising progress, most existing VLA models primarily rely on 2D visual inputs and lack effective mechanisms for modeling the 3D spatial structure of the scene, which limits their spatial reasoning capabilities in complex manipulation tasks.

Robot Learning with 3D Information. In response to the spatial limitations of 2D-based VLA models, several recent approaches Cai et al. (2024); Zhen et al. (2024); Li et al. (2025a); Qu et al. (2025); Chen et al. (2024); Goyal et al. (2024); Jia et al. (2024) have explored integrating 3D information to enhance spatial understanding. For example, 3D-VLA Zhen et al. (2024) fuses 3D perception, reasoning, and action through a 3D-based large language model Hong et al. (2023), trained on a large-scale 3D dataset curated from existing embodied robotics benchmarks. To make 3D-aware policies applicable in real-world scenarios, methods such as SpatialVLA Qu et al. (2025) and PointVLA Li et al. (2025a) incorporate depth information captured from additional RGB-D cameras or depth estimation models, which enhances 3D scene understanding and enables more accurate perception of spatial relationships, object geometry, and depth-aware interactions.

Despite these advances, a fundamental limitation of current 3D-aware VLA methods lies in their reliance on explicit 3D inputs such as depth maps and point clouds, which require either specialized sensors or auxiliary estimation models. This dependency imposes constraints on scalability, deployment flexibility, and general applicability in diverse real-world environments.

To address this issue, we propose integrating VGGT Wang et al. (2025) into existing VLA models. While keeping the input as RGB images, VGGT implicitly models 3D structure by fusing spatial features from multi-view observations. Our approach serves as a bridge between pure 2D input models and explicit 3D perception methods, enhancing spatial understanding without requiring additional sensors or depth estimation modules.

3 METHOD

3.1 PRELIMINARIES

Vision-Language-Action Models. As a promising approach toward generalist robot policies, Vision-Language-Action (VLA) models have emerged as an increasingly popular research direction Black et al. (2024); Kim et al. (2024); Li et al. (2025a); Liu et al. (2025). VLAs aim to bridge the gap between high-level human instructions and low-level robotic actions by leveraging the rich multimodal priors encoded in large-scale pre-trained Vision-Language Models (VLMs), such as Paligemma Beyer et al. (2024), CLIP Radford et al. (2021), LLaMA Touvron et al. (2023a;b), and Flamingo Alayrac et al. (2022). These VLMs are trained on vast and diverse internet-scale image-text pairs, endowing them with strong world knowledge and the ability to ground natural language in visual concepts.

Unlike traditional imitation learning methods that typically train a task-specific policy from scratch, VLA models reuse this pretrained multimodal understanding to enable more flexible and scalable robotic behaviors. In particular, the VLM serves as a general-purpose semantic encoder, while a downstream module—commonly referred to as the *action expert*—learns to map the fused representations into robot control commands. This modular design separates general world understanding from task-specific actuation, allowing the model to generalize better across instructions and visual environments.

Formally, at each timestep t , the VLA model receives multi-view visual observations $\{I_t^i\}_{i=1}^N$ and a language instruction L , which are jointly encoded by the VLM to produce a contextual embedding z_t . This embedding is then concatenated with robot-specific states S_t (e.g., joint angles, gripper status, or end-effector pose), and passed to the action expert to generate the low-level control command A_t . The entire pipeline thus defines a conditional distribution $p(A_t | I_t^i, L, S_t)$.

Compared to standard imitation learning policies, which are typically trained on a specific task, the VLA framework improves *semantic grounding*, *modality fusion*, and *generalization capability*. This enables robots not only to follow diverse and abstract language instructions but also to adapt to new tasks and visual scenes with minimal fine-tuning.

Visual Geometry Foundation Models. Unlike traditional SLAM or depth estimation pipelines that rely on finely-tuned modules and sensors, Visual Geometry Foundation Models (VGFM) Leroy et al. (2024); Wang et al. (2024; 2025); Li et al. (2025c) are a class of vision models trained to reconstruct 3D structural information from 2D visual inputs. Since VGFM are trained with geometric supervision, they have the ability to recover fine-grained spatial structure from multi-view monocular inputs. These models provide strong structural priors for downstream tasks such as spatial understanding, especially when explicit 3D sensors are unavailable.

Given a set of multi-view images $\{I^i\}_{i=1}^N$, a typical VGFM predicts a 3D point cloud P representing the scene as

$$f_{\text{VGFM}}(\{I^i\}_{i=1}^N) = P. \tag{1}$$

These geometry-aware models complement vision-language systems by injecting 3D structural cues, enhancing spatial grounding from purely 2D observations such as video frames.

Recently, Visual Geometry Grounded Transformer (VGGT) Wang et al. (2025) has introduced a novel feed-forward architecture and demonstrated impressive performance in 3D attributes prediction. It takes an arbitrary number of image views as input and alternates between frame-wise and global self-attention to model spatial consistency. Given a sequence of N RGB images $\{I^i\}_{i=1}^N$, where each $I_i \in \mathbb{R}^{3 \times H \times W}$, the model outputs a set of 3D annotations for each frame, including predicted camera poses g_i , depth maps D_i , point maps P_i , and 3D point tracks T_i , i.e.,

$$f(\{I_i\}_{i=1}^N) = (g_i, D_i, P_i, T_i)_{i=1}^N. \tag{2}$$

3.2 PROPOSED VLA ARCHITECTURE

Recent 3D-based VLA models, such as PointVLA Li et al. (2025a) and SpatialVLA Qu et al. (2025), often employ explicit 3D inputs like point clouds or depth maps to enhance spatial understanding. While effective, these approaches typically require additional sensors and preprocessing, and are

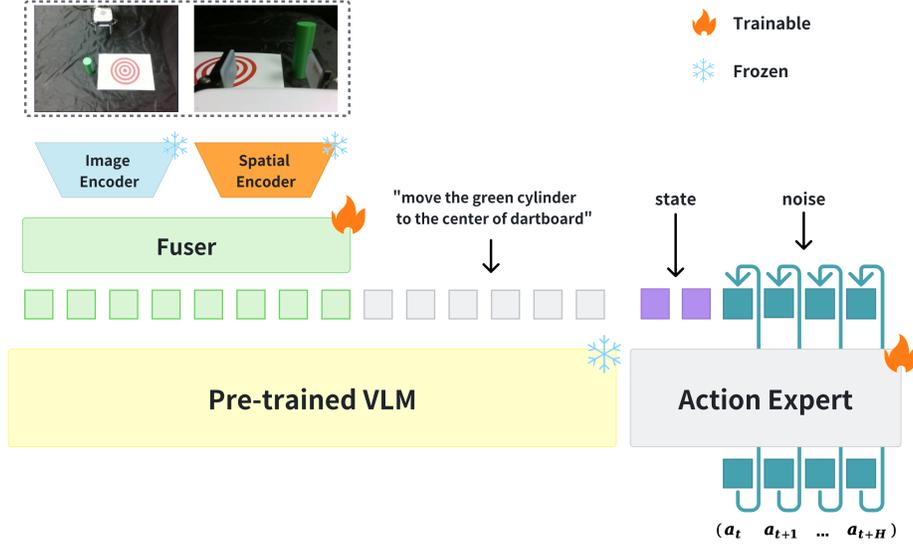


Figure 1: Architecture of Evo-0.

often sensitive to variations in camera viewpoints. In contrast, VGGT presents a promising alternative for implicitly introducing spatial awareness, benefiting from its diverse training data and elegant feed-forward architecture. Recent studies have successfully applied VGGT to VLM architectures Wu et al. (2025) and SLAM systems Maggio et al. (2025), demonstrating that geometry-grounded visual tokens can improve spatial understanding in both multimodal learning and classical robotic perception.

Motivated by these findings, we hypothesize that introducing geometry-aware visual representations from VGGT into the action prediction pipeline can enrich spatial context, leading to more precise and generalizable policy learning without requiring explicit point cloud or depth inputs. To evaluate this hypothesis, we build our model upon $\pi 0$ Black et al. (2024), a state-of-the-art open-source VLA model, and incorporate geometry-aware features from VGGT into its visual embedding stream. The architecture is described in Figure 1. Specifically, we utilize VGGT as a spatial encoder and extract tokens from its final layer:

$$\mathcal{E}_l(\{I_i\}_{i=1}^N) = t_c, t_r, t_{3D}, \quad (3)$$

where N is the number of views, l denotes the layer index, and t_c , t_r , and t_{3D} denote the camera, register, and 3D tokens, respectively. We extract the 3D tokens t_{3D} to inject spatial information, as they are originally trained to conduct 3D tasks in VGGT. These tokens capture rich geometric representations, including depth-aware context, temporally consistent object trajectories, and spatial correspondences across views.

To integrate the VGGT-derived token features into the vision-language pipeline, we introduce a lightweight fuser module that combines embeddings from the Vision Transformer Dosovitskiy et al. (2020) and the VGGT encoder. Specifically, the fuser consists of a single cross-attention layer, where the 2D visual tokens $t_{2D} \in \mathbb{R}^{N \times M_{2D} \times d_{2D}}$ serve as queries, and the VGGT-derived tokens $t_{3D} \in \mathbb{R}^{N \times M_{3D} \times d_{3D}}$ act as keys and values. Here, M_{2D} and M_{3D} denote the number of tokens from ViT and VGGT encoder, respectively. The 2D visual tokens are then updated as follows:

$$Q = t_{2D}W_Q, \quad K = t_{3D}W_K, \quad V = t_{3D}W_V, \quad (4)$$

$$t^i = \text{softmax}\left(\frac{Q^i(K^i)^\top}{\sqrt{d}}\right)V^i, \quad (5)$$

$$t = \text{Concat}_{i=1}^N(t^i), \quad (6)$$

where $W_Q \in \mathbb{R}^{d_{2D} \times d}$, and $W_K, W_V \in \mathbb{R}^{d_{3D} \times d}$ are trainable projection matrices shared across views. Each view $i \in 1, \dots, N$ is processed independently via the cross-attention module, and the resulting tokens are concatenated to form the fused output t .



Figure 2: Illustration of the task setup.

The fused tokens are then forwarded to the PaliGemma Beyer et al. (2024) vision-language model, which jointly attends over both the geometry-enhanced visual input and the language tokens to predict actions. To maintain computational efficiency and minimize disruption to the pretrained VLM backbone, we freeze the core VLM parameters and insert lightweight Low-Rank Adaptation (LoRA) Hu et al. (2022) layers. During training, only the fuser module, LoRA layers, and the flow-matching action expert are fine-tuned, enabling effective adaptation with minimal overhead.

4 EXPERIMENTS

4.1 IMPLEMENTATION DETAILS

Our framework builds upon the open-source VLA model Black et al. (2024). For each task, we collect 100 expert demonstrations using tele-operation. To promote diversity and robustness, the positions of objects and targets are randomly perturbed during data collection. The model is trained using AdamW with weight decay of 10^{-10} . The cosine learning rate schedule is used with a learning rate of 2.5×10^{-5} , warmup over 1000 steps, and decay to 2.5×10^{-6} . Training is performed using bfloat16 mixed precision on a single NVIDIA A800 GPU (80GB) with a batch size of 32.

4.2 TASK SETUP

We design five tasks for real-world robot evaluation, which span a range of spatial understanding challenges, from fine-grained geometric alignment to pick-and-place and transparent object interaction. In particular, each task has a low tolerance for spatial error, as minor inaccuracies in the spatial predictions can lead to task failure. This makes them well-suited for assessing if representations from VGGT can enhance VLA’s spatial understanding.

A detailed description of the five tasks is provided below, with an illustration shown in Figure 2.

1. **Centering a cylinder on a target.** The robot is required to align a cylindrical object precisely at the center of a marked target area on the table. This task resembles target shooting: the target has concentric rings, and scoring is based on which ring the center of the cylinder falls into. The closer to the center, the higher the score.
2. **Peg-in-hole insertion.** This task requires the robot to insert a cylindrical peg into one of three tightly fitting holes on a board. This necessitates accurate alignment in 3D space, as small tilting or offset could cause task failure.
3. **Middle bottle grasping.** Three bottles are closely placed in a row, and the robot is instructed to pick the middle one. This setup mimics a grocery store scenario, where items are densely arranged on shelves. Success is defined as picking up the middle bottle without touching or knocking over the adjacent ones.
4. **Can pick-and-place.** In this task, the robot must pick up a standard can and place it in a designated spot on a shelf. The location of the placement is varied across trials in both position and height, requiring the model to generalize spatial understanding to different configurations.
5. **Transparent object pick-and-place.** The task setup is similar to the previous one, but involves transparent objects such as glass bottles. This presents additional challenge, since transparent materials are often poorly captured by RGB sensors and are prone to glare, making them difficult to perceive and localize.

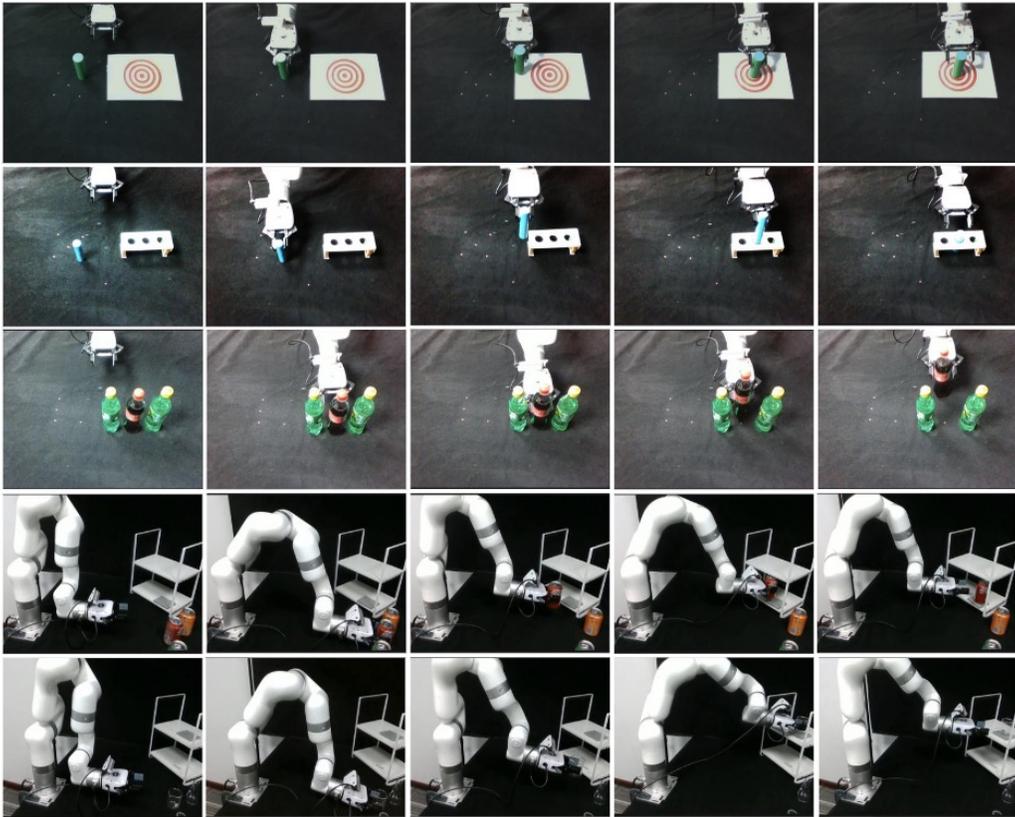


Figure 3: Qualitative results of our model in real-world tasks.

	Task1 (15)	Task2 (15)	Task3 (15)	Task4 (10)	Task5 (20)	Average (75)
π_0	59.33%	20.00%	13.30%	20.00%	30.00%	28.53%
Ours	68.67%	66.67%	26.70%	60.00%	65.00%	57.41%

Table 1: Success rates for five real-world tasks. The number in parentheses represents the number of trails for each task. In particular, Task 1 is evaluated using scores; for example, 68.67% denotes that our model achieves an average score of 3.43 out of a maximum of 5.

We evaluate all tasks using binary task completion except for Task 1. For Task 1, we adopt a more fine-grained evaluation inspired by target shooting: the innermost ring yields the highest score (5), while outer rings correspond to decreasing accuracy (4 to 1). A score of 0 is assigned if the robot fails to grasp the object. This scoring formulation captures subtle differences in spatial precision that would be lost in a binary metric. We report the overall success rate (or average score in the case of Task 1) for each task. The object positions are marked by stickers to ensure fair comparison and reproducibility.

4.3 REAL-WORLD EXPERIMENT

Quantitative Results. We evaluate the effectiveness of our proposed method by comparing it against the baseline model π_0 . The quantitative results are presented in the Table 1. Across all tasks, our method achieves consistent improvements over the baseline, indicating that the implicit 3D geometry features contribute positively to task performance. Notably, our model demonstrates the largest performance gain on Task 2 (peg-in-hole insertion), a particularly challenging task that demands accurate spatial reasoning. Furthermore, Task 3 (middle bottle grasping) poses a substantial challenge due to the narrow margin between adjacent bottles, requiring the gripper to perform

careful, collision-free insertion and grasping. Our method exhibits reasonable improvement on this task compared to the baseline, demonstrating enhanced spatial understanding and control in cluttered environments. Overall, we achieve a 28.88% performance gain in the average success rate.

Qualitative Results. In Figure 3, we present visualizations of task executions across different tasks. These visual results further complement the quantitative findings, showcasing our model’s enhanced spatial awareness and manipulation precision. For instance, in the cylinder-centering and peg-in-hole insertion tasks, our model reliably achieves stable grasping and precise alignment with the target area. In contrast, the baseline π_0 often fails to establish a proper grasp on the cylinder from the initial step, leading to unsuccessful or unstable placement attempts.

5 CONCLUSION

In this paper, we explore using implicit 3D representations to enhance spatial understanding in Vision-Language-Action (VLA) models. By leveraging features from the Visual Geometry Grounded Transformer (VGGT), trained on large-scale 2D–3D paired data, we inject strong geometric priors into VLA models without relying on explicit 3D inputs. Through extensive experiments across five spatially challenging tasks, we demonstrate that our approach significantly outperforms baseline models, validating the effectiveness of the proposed implicit geometric prior integration. Our method offers a simple and efficient solution for enhancing spatial understanding in VLA systems.

REFERENCES

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. URL <https://arxiv.org/abs/2204.14198>.
- Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. *arXiv preprint arXiv:2407.07726*, 2024.
- Shariq Farooq Bhat, Reiner Birkel, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023.
- Johan Bjorck, Fernando Castañeda, Nikita Cherniadev, Xingye Da, Runyu Ding, Linxi Fan, Yu Fang, Dieter Fox, Fengyuan Hu, Spencer Huang, et al. Gr00t n1: An open foundation model for generalist humanoid robots. *arXiv preprint arXiv:2503.14734*, 2025.
- Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. π_0 : A vision-language-action flow model for general robot control, 2024. URL <https://arxiv.org/abs/2410.24164>, 2024.
- Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- Wenxiao Cai, Yaroslav Ponomarenko, Jianhao Yuan, Xiaoqi Li, Wankou Yang, Hao Dong, and Bo Zhao. Spatialbot: Precise spatial understanding with vision language models. *arXiv preprint arXiv:2406.13642*, 2024.
- Shizhe Chen, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Sugar: Pre-training 3d visual representations for robotics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18049–18060, 2024.
- Erik Daxberger, Nina Wenzel, David Griffiths, Haiming Gang, Justin Lazarow, Gefen Kohavi, Kai Kang, Marcin Eichner, Yinfei Yang, Afshin Dehghan, et al. Mm-spatial: Exploring 3d spatial understanding in multimodal llms. *arXiv preprint arXiv:2503.13111*, 2025.

- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Ankit Goyal, Valts Blukis, Jie Xu, Yijie Guo, Yu-Wei Chao, and Dieter Fox. Rvt-2: Learning precise manipulation from few demonstrations. *arXiv preprint arXiv:2406.08545*, 2024.
- Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. *Advances in Neural Information Processing Systems*, 36:20482–20494, 2023.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Yueru Jia, Jiaming Liu, Sixiang Chen, Chenyang Gu, Zhilue Wang, Longzan Luo, Lily Lee, Pengwei Wang, Zhongyuan Wang, Renrui Zhang, et al. Lift3d foundation policy: Lifting 2d large-scale pretrained models for robust 3d robotic manipulation. *arXiv preprint arXiv:2411.18623*, 2024.
- Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. In *European Conference on Computer Vision*, pp. 71–91. Springer, 2024.
- Chengmeng Li, Junjie Wen, Yan Peng, Yaxin Peng, Feifei Feng, and Yichen Zhu. Pointvla: Injecting the 3d world into vision-language-action models. *arXiv preprint arXiv:2503.07511*, 2025a.
- Peiyan Li, Yixiang Chen, Hongtao Wu, Xiao Ma, Xiangnan Wu, Yan Huang, Liang Wang, Tao Kong, and Tieniu Tan. Bridgevla: Input-output alignment for efficient 3d manipulation learning with vision-language models. *arXiv preprint arXiv:2506.07961*, 2025b.
- Zhengqi Li, Richard Tucker, Forrester Cole, Qianqian Wang, Linyi Jin, Vickie Ye, Angjoo Kanazawa, Aleksander Holynski, and Noah Snavely. Megasam: Accurate, fast and robust structure and motion from casual dynamic videos. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 10486–10496, 2025c.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Jiaming Liu, Hao Chen, Pengju An, Zhuoyang Liu, Renrui Zhang, Chenyang Gu, Xiaoqi Li, Ziyu Guo, Sixiang Chen, Mengzhen Liu, et al. Hybridvla: Collaborative diffusion and autoregression in a unified vision-language-action model. *arXiv preprint arXiv:2503.10631*, 2025.
- Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint arXiv:2209.14577*, 2022.
- Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024.
- Dominic Maggio, Hyungtae Lim, and Luca Carlone. Vggt-slam: Dense rgb slam optimized on the sl(4) manifold, 2025. URL <https://arxiv.org/abs/2505.12549>.
- Abby O’Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 6892–6903. IEEE, 2024.
- Delin Qu, Haoming Song, Qizhi Chen, Yuanqi Yao, Xinyi Ye, Yan Ding, Zhigang Wang, JiaYuan Gu, Bin Zhao, Dong Wang, et al. Spatialvla: Exploring spatial representations for visual-language-action model. *arXiv preprint arXiv:2501.15830*, 2025.

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 5294–5306, 2025.
- Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 20697–20709, 2024.
- Diankun Wu, Fangfu Liu, Yi-Hsin Hung, and Yueqi Duan. Spatial-mlm: Boosting mllm capabilities in visual-based spatial intelligence, 2025. URL <https://arxiv.org/abs/2505.23747>.
- Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10371–10381, 2024.
- Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang, Xin Yan, Yilun Du, Yining Hong, and Chuang Gan. 3d-vla: A 3d vision-language-action generative world model. *arXiv preprint arXiv:2403.09631*, 2024.