

# SPEAR-1: Scaling Beyond Robot Demonstrations via 3D Understanding

Nikolay Nikolov Giuliano Albanese Sombit Dey Aleksandar Yanev

Luc Van Gool Jan-Nico Zaech Danda Pani Paudel

{nikolay.nikolov, giuliano.albanese, sombit.dey, aleksandar.yanев  
luc.vangool, jan-nico.zaech, danda.paudel}@insait.ai

INSAIT, Sofia University "St. Kliment Ohridski"

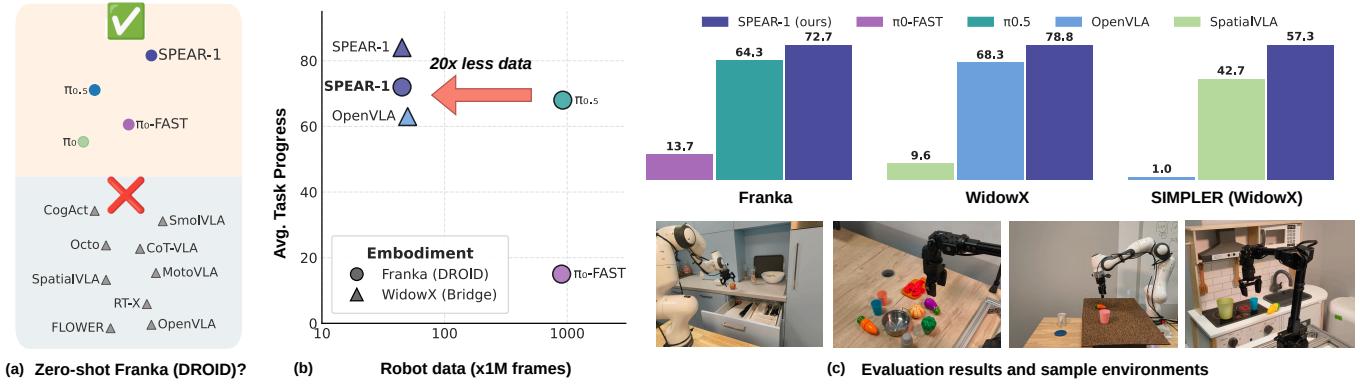


Figure 1. (a) Most VLAs fail to show zero-shot performance on the challenging Franka (DROID) setup in unseen environments, without task or environment-specific fine-tuning. (b) SPEAR-1 operates in this challenging setup, outperforms  $\pi_0$ -FAST [33] and matches  $\pi_0.5$  [6] on Franka (DROID) embodiment zero-shot in unseen environments while using  $20\times$  less robot demonstrations data. It also shows strong performance on WidowX (Bridge). (c) SPEAR-1 evaluation results on different embodiments and in different environments.

## Abstract

Robotic Foundation Models (RFMs) hold great promise as generalist, end-to-end systems for robot control. Yet their ability to generalize across new environments, tasks, and embodiments remains limited. We argue that a major bottleneck lies in their foundations: most RFMs are built by fine-tuning internet-pretrained Vision-Language Models (VLMs). However, these VLMs are trained on 2D image-language tasks and lack the 3D spatial reasoning inherently required for embodied control in the 3D world. Bridging this gap directly with large-scale robotic data is costly and difficult to scale. Instead, we propose to enrich easy-to-collect non-robotic image data with 3D annotations and enhance a pretrained VLM with 3D understanding capabilities. Following this strategy, we train SPEAR-VLM, a 3D-aware VLM that infers object coordinates in 3D space from a single 2D image. Building on SPEAR-VLM, we introduce our main contribution, **SPEAR-1**: a robotic foun-

dation model that integrates grounded 3D perception with language-instructed embodied control. Trained on  $\sim 45M$  frames from 24 Open X-Embodiment datasets, SPEAR-1 outperforms or matches state-of-the-art models such as  $\pi_0$ -FAST and  $\pi_0.5$ , while it uses  $20\times$  fewer robot demonstrations. This carefully-engineered training strategy unlocks new VLM capabilities and as a consequence boosts the reliability of embodied control beyond what is achievable with only robotic data. We make our model weights and 3D-annotated datasets publicly available.

## 1. Introduction

Vision-Language-Action (VLA) models have emerged as a promising paradigm for building generalist, end-to-end systems for robot control. Their success relies on two factors: (1) the strong visual-linguistic understanding inherited from internet-scale pretraining of the underlying Vision Language Model (VLM), which provides broad “com-

mon sense” knowledge, and (2) training on large, diverse datasets of robot demonstrations.

Despite this progress, the landscape of generalist VLA policies remains fragmented in terms of generalization – across embodiments, environments, and tasks. This becomes especially prominent in zero-shot performance in *unseen* real-world environments with variations in camera positions and out-of-distribution backgrounds, such as the typical deployment scenarios of the Franka (DROID) setup [2]. In contrast, as shown in Fig. 1 (a), most existing VLAs (*e.g.* OpenVLA [19], CogAct [21], SpatialVLA [35], MotoVLA [40]) achieve high zero-shot performance<sup>1</sup> in “toy” environments with seen camera positions, but struggle with zero-shot performance in *unseen* challenging Franka scenarios and depend on task- or environment-specific fine-tuning. Recent efforts such as  $\pi_0$  and  $\pi_{0.5}$  [6] push toward broader generalization, yet at the cost of closed large-scale robotic data. We introduce **SPEAR-1**, which advances these desired generalization capabilities while being substantially more data-efficient. Quantitatively (see Fig. 1 (b)), **SPEAR-1** outperforms  $\pi_0$ -FAST [33] and matches  $\pi_{0.5}$  on multiple robot embodiments using  $20\times$  fewer demonstrations, which is especially important given the high cost and logistical difficulty of collecting real-world robotic data.

We achieve this efficiency by introducing explicit 3D awareness into the vision-language backbone before any robot training. The model incorporates a pretrained depth encoder and is optimized on 3D-aware vision-language tasks such as distance estimation and 3D bounding box prediction, embedding control-relevant spatial reasoning directly into its representations. Achieving such integration is non-trivial: aligning 3D geometric cues with high-level linguistic and visual features requires detailed multimodal dataset annotations and precise cross-modal calibration, as naive fusion often degrades both semantic understanding and spatial accuracy. In contrast, existing VLAs rely on 2D VLMs that excel at semantic perception, but lack geometric understanding, forcing them to learn 3D structure implicitly from large-scale robot demonstrations. This dependence on costly and embodiment-specific data limits scalability and generalization across environments, underscoring the difficulty and significance of **SPEAR-1**’s design.

In our progression from spatial understanding to embodied control, we introduce a staged training pipeline, as shown in Fig. 2. In *Stage 1*, we develop a 3D-aware vision-language model, **SPEAR-VLM**, which extends a pretrained VLM by learning spatial reasoning from non-robotic 2D images annotated with 3D cues. This stage establishes a perceptual backbone that encodes geometric relations while preserving the rich semantic priors of large-

scale pretraining. In *Stage 2*, we introduce an *action expert* that maps the grounded visual-language representations to motor actions. This stage demands well-tuned vision-language-action modeling choices and a carefully-crafted multi-embodiment data processing strategy to learn precise low-level robot controls. Together, these stages bridge the gap between internet-scale 2D perception and embodied 3D interaction, progressively transforming passive spatial understanding into actionable behavior.

Unlike previous works that address the challenge of 3D knowledge for robot control [35, 40], **SPEAR-1** demonstrates improvement on a foundation model level, with an end-to-end policy across multiple different environments and robots. It is capable of achieving state-of-the-art robot control on multiple robot embodiments without requiring target evaluation environment fine-tuning. Furthermore, **SPEAR-1** demonstrates how significant amounts of robot demonstration data can be “replaced” by non-robotic 3D-annotated image data. In summary, our work makes the following contributions:

- **SPEAR-VLM**: a VLM with *control-inspired 3D capabilities* (*e.g.* localizing objects in 3D), trained on carefully-crafted VQA tasks and enriched 2D-image non-robotic data. Importantly, **SPEAR-VLM** directly boosts downstream VLA performance.
- **SPEAR-1**: an open-weight *robotics foundation model with 3D understanding*, which significantly outperforms or matches the strongest state-of-the-art baselines, trained with  $20\times$  more robot demonstration data
- **Extensive experimental validation**: we demonstrate strong generalization across diverse settings with a substantial reduction in reliance on hard-to-collect robotic data. Notably, using only **200k non-robotic 2D images**, **SPEAR-1** surpasses models trained with more than **900M additional frames of robotic demonstrations**.

## 2. Related Work

**Spatial Understanding for VLMs.** The majority of existing VLMs trained on large-scale datasets have been limited to flat 2D image understanding [4, 17, 26, 41, 43, 46]. Our work extends the PaliGemma VLM [4] by integrating the MoGe monocular depth estimator [47] as a supplementary vision backbone and by training on manipulation-relevant 3D tasks to enhance the VLM’s 3D understanding. Previously, Chen et al. [8] used a similar data annotation approach for training a 3D-aware VLM, but they do not integrate a pretrained depth estimator and neither their model nor their dataset is publicly accessible. Additionally, unlike SpatialVLM [8] or RoboSpatial [39], trained on high-level spatial relationships, our **SPEAR-VLM** focuses on explicit 3D coordinate prediction, a pretraining task much closer to embodied control. SpatialBot [7] also previously proposed a spatially-aware VLM targeting robot control, but

<sup>1</sup>For more details on the definition of “zero-shot performance in unseen environments” see Appendix A.4

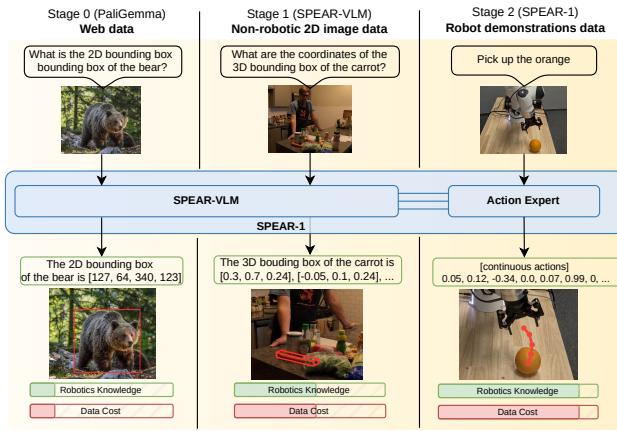


Figure 2. **SPEAR-1 stages of training.** **Stage 0:** General VLM pretraining on web scale data, e.g. PaliGemma. **Stage 1:** Integrate a mono depth vision encoder to build **SPEAR-VLM** and train it on embodied-inspired VQA tasks, e.g. 3D bounding box or object-to-object distance estimation. We use 2D images from non-robotic data, enriched with 3D annotations. **Stage 2:** Add an *action expert* to train **SPEAR-1** on robot demonstration data, e.g. OpenX [32]. Each stage boosts the model’s robotics-relevant knowledge and capabilities, but the abundance and diversity of data decreases.

their method involves a multi-step VLM inference process and was never shown to integrate into a VLA for generalist robotic control.

**Vision-Language-Action Models.** Recently, multiple works have developed generalist robot policies [5, 6, 10, 19, 32, 33, 53] trained on multiple robot embodiments. SPEAR-1 builds on top of the  $\pi_0$  architecture, which combines a pretrained PaliGemma VLM and an action expert module, but we initialize the underlying VLM from our SPEAR-VLM to integrate pretrained 3D understanding. Previously, SpatialVLA [35] proposed integrating a monocular depth encoder [48] in the VLA, but without any VLM alignment or pretraining and therefore learning 3D capabilities entirely from hard-to-collect robotic data. MolmoAct [20] recently proposed a spatially-aware VLA, but the approach involves ‘reasoning’ at inference time, rendering the method too slow for real-time control. Most closely related, Gemini Robotics 1.0 [44] follows a similar 3D pre-training method to fine-tune the significantly larger Gemini 2.0 [34] and distill into a smaller VLA with reasoning capabilities. With most of the method’s details undisclosed, our work still differs in several important aspects: (1) we investigate the benefits of 3D pretraining in isolation, (2) train much smaller open-access model on limited, less diverse open data from OpenX [32], and, most importantly, (3) we demonstrate the ability to reduce the need for robotic data with non-robotic 2D images.

### 3. Method

In this section, we describe SPEAR-1 and its training recipe in detail. In section 3.1 we describe the architecture, data generation pipeline, and training procedure of our 3D-aware SPEAR-VLM. This stage aims to enhance the 3D spatial understanding capabilities of an off-the-shelf VLM through fine-tuning on 3D spatial perception tasks. We then proceed, in section 3.2 to detail the architecture and training procedure of SPEAR-1, which comprises a pre-training and post-training stage.

#### 3.1. SPEAR-VLM

Our approach considers the architecture of recent robotics foundational models that are based on VLMs, pretrained on large corpora of internet-scale text-image data. The architecture of those models usually consists of a vision encoder, a vision-to-text-embedding feature projector, and a LLM. The majority of the tasks on which VLMs are usually trained are limited to 2D [4, 17, 25]. To extend the capabilities of a pre-trained VLM to 3D understanding, we propose (1) extending the model architecture by adding a monocular depth encoder and (2) training the VLM on VQA tasks that require explicit 3D reasoning.

**SPEAR-VLM Architecture.** Our model builds on PaliGemma [4] as backbone, but the same method can be used with any late-fusion VLM [1, 11, 27]. PaliGemma consists of three main components: (1) a SigLIP **visual encoder** [50], (2) a linear **projector** that maps the visual tokens predicted by the visual encoder to the language model input space and (3) a Gemma **language model** [42]. To enable the model to perceive accurate depth, we integrate the MoGe [47] depth encoder as an additional vision encoder. We choose MoGe due to its affine-invariant modeling approach, capable of fitting cameras with different intrinsics. Our intuition is that affine-invariant depth should generalize better across environments thus being better suited for learning generalist robot control. Similar to the MoGe decoder inputs, we concatenate the intermediate features from the last 4 layers of the MoGe ViT encoder along the feature dimension and project them to the LLM embedding space via a randomly-initialized linear projector. The visual input to the LLM consists of the averaged outputs of the SigLIP and MoGe projectors. To encode 3D information into text we extend the PaliGemma tokenizer with  $N = 1024$  3D tokens (see Fig. 3 and Appendix A.1.3).

**3D pretraining tasks.** Given the above architecture, we propose a pre-training scheme to enable the model to leverage the depth information in MoGe’s encoder features and acquire 3D spatial understanding capabilities. To embed as much control-relevant 3D knowledge in the SPEAR-VLM as possible, we design VQA tasks inspired by the embodied tasks a VLA needs to learn, e.g. ‘Output the vertices of the 3D bounding box of object X’ or ‘Output the xyz compo-

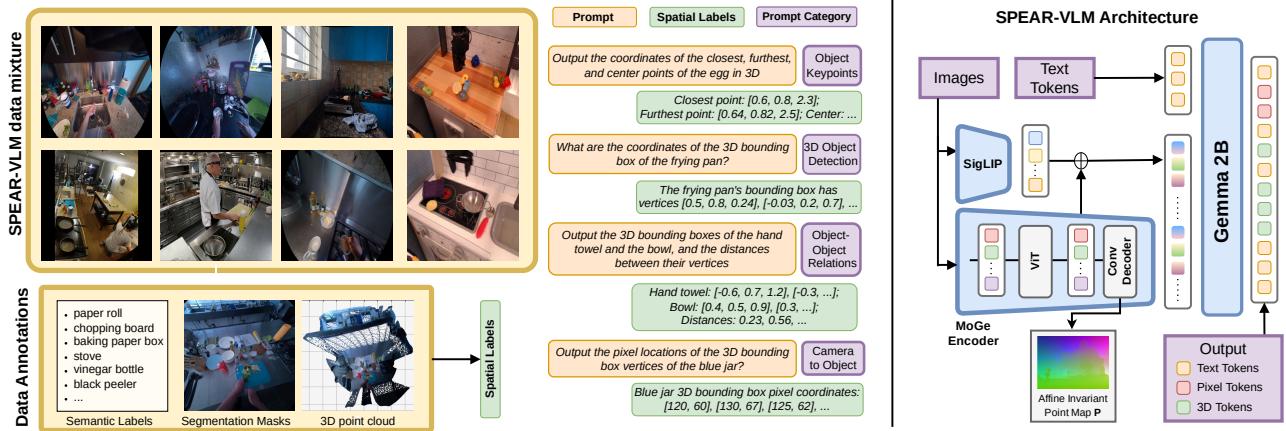


Figure 3. **SPEAR-VLM overview.** Left: Training data mixture, auto computed spatial annotations and example question-answer pairs from each category. Right: High-level architecture with fusion between SigLIP and MoGe encoders and PaliGemma embeddings expansion with 3D tokens. This design equips SPEAR-VLM with explicit 3D understanding that serves as a strong foundation for SPEAR-VLA.

nents of the distance between object  $X$  and object  $Y$ '. These VLM pre-training tasks ensure learning semantic 3D localization, object-to-object spatial relations, and 3D coordinate system geometry (Fig. 3). For a full list of question-answer pairs, see Appendix A.1.1.

**3D Vision-Question-Answering Data.** As few open datasets contain the annotations needed for the proposed training scheme, we devise the following semi-automatic annotation pipeline to enrich existing datasets with the necessary annotations: *object-level segmentation masks*, *semantic labels* and *projected 3D point cloud*. Importantly, our pipeline requires only 2D images as input and off-the-shelf vision foundation models:

1. Use Gemini [9] to detect 2D bounding boxes and semantic labels for the objects in the image.
2. Prompt SAM2 [36, 37] with the detected bounding boxes to produce instance-level segmentation masks.
3. Obtain 3D point cloud annotations for the entire image via MoGe direct point cloud predictions [47].

To construct a training example, we randomly sample a templated text prompt and a set of objects from the image. We then filter the annotated MoGe 3D point cloud with the object mask to obtain the object 3D point cloud. Based on the segmented point cloud, we compute the oriented 3D bounding box and construct the question-answer pair.

We focus on indoor environments and annotate the "cooking" and "bike repair" parts of EgoExo4D [13] that already have segmentation masks, resulting in 200k images. For visual diversity, we further annotate 30k frames of the Bridge-V2 [45] robot demonstration dataset, downsampled to 10% in the VLM training data mixture.

**Training process.** Similar to LLaVa [25], we train SPEAR-VLM in two stages. In the first stage, we initialize from PaliGemma and MoGe weights, with the MoGe pro-

jector and the LLM 3D token embeddings initialized randomly. We train only the randomly initialized weights and SigLIP projector, keeping everything else frozen. In the second and longer stage, we keep only SigLIP and MoGe encoders frozen and we scale the next-token-prediction loss for 3D tokens by a factor  $\lambda = 2$ .

### 3.2. SPEAR-1

SPEAR-1 follows a similar overall architecture as  $\pi_0$  [5], however, we build on SPEAR-VLM, use a rotation formulation in flow matching on the  $\mathbb{S}^3$  manifold of unit quaternions, and several data & engineering improvements. Design decisions were ablated on small-scale experiments on BridgeData V2 [45] due to the cost of training on the entire OpenX mixture. We summarize these key decisions and learning in the following.

**Preliminaries.** Formally, we aim to learn a function  $\pi(\cdot)$  mapping an observation  $\mathbf{o}_t$  to a sequence of robot actions  $\mathbf{A}_t = [\mathbf{a}_t, \mathbf{a}_{t+1}, \dots, \mathbf{a}_{t+H-1}]$  over a horizon  $H$ . The observation is defined as  $\mathbf{o}_t = [\mathbf{I}_t^1, \dots, \mathbf{I}_t^n, \mathbf{p}_t, \mathbf{l}_t]$ , where  $\mathbf{I}_t^i$  is the  $i$ -th image observation from an uncalibrated camera,  $\mathbf{p}_t$  is a vector containing the robot state comprising of the end-effector pose and gripper state,  $\mathbf{l}_t$  is a vector of language tokens representing the language instruction.

**Architecture.** We follow the broadly accepted architecture introduced in  $\pi_0$ : a Flow Matching action expert that processes proprioception observations and predicts the robot actions by attending to the VLM's intermediate key-value pairs. For full details, see Appendix A.2.2 and  $\pi_0$  [5].

**Flow Matching Formulation.** The action sequence prediction is supervised via conditional flow matching [23, 24, 28]. The model takes as input the observation  $\mathbf{o}_t$ , the flow-matching step  $\tau \in [0, 1]$  and a sequence of noisy actions  $\mathbf{A}_t^\tau = [\mathbf{a}_t^\tau, \dots, \mathbf{a}_{t+H-1}^\tau]$  and predicts a denoising vector

$\mathbf{v}_\theta(\mathbf{A}_t^\tau, \mathbf{o}_t)$ . We denote the decomposed action of translation, rotation and gripper components as  $\mathbf{a}_t = [\mathbf{x}_t, \mathbf{q}_t, \mathbf{g}_t]$ . We use the square brackets operator  $[\cdot]$  on the predicted denoising vector  $\mathbf{v}_\theta$  and the denoising vector field  $\mathbf{u}$  to denote a specific component, e.g.  $\mathbf{u}[\mathbf{x}_t]$  corresponds to the translation component of the denoising vector field.

We follow a flow matching formulation in linear space for translation and on the  $\mathbb{S}^3$  manifold of unit quaternions for rotation. For simplicity, we omit the gripper component as it follows the same linear formulation as translation.

During training, we sample a random timestep  $\tau \sim \mathcal{B}(\alpha, \beta)$  and random noise  $\mathbf{x}_\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\mathbf{q}_\epsilon \sim \mathcal{U}(\mathbb{S}^3)$ . “Noisy actions” are computed by linear interpolation for translation  $\mathbf{x}_t^\tau = \tau \mathbf{x}_t + (1 - \tau) \mathbf{x}_\epsilon$  and spherical linear interpolation on the  $\mathbb{S}^3$  manifold for quaternion rotation

$$\mathbf{q}_t^\tau = \frac{\sin((1 - \tau)\theta)}{\sin \theta} \mathbf{q}_\epsilon + \frac{\sin(\tau\theta)}{\sin \theta} \mathbf{q}_t, \quad (1)$$

with  $\theta = \cos^{-1}(\mathbf{q}_\epsilon \cdot \mathbf{q}_t)$ . The “noisy action sequence”  $\mathbf{A}_t^\tau$  is passed as input to the model and trained to output the denoising vector field  $\mathbf{u}(\mathbf{A}_t^\tau | \mathbf{A}_t) = \frac{d\mathbf{A}_t^\tau}{d\tau}$ . Training is supervised with the conditional flow-matching loss, equivalent to the MSE loss for translation

$$\mathcal{L}_{\mathbb{R}^3}(\theta) = \left\| \mathbf{v}_\theta(\mathbf{A}_t^\tau, \mathbf{o}_t)[\mathbf{X}_t] - \mathbf{u}(\mathbf{A}_t^\tau | \mathbf{A}_t)[\mathbf{X}_t] \right\|^2. \quad (2)$$

For rotations, we combine a cosine loss between the velocity predictions  $\mathbf{v}_\theta(\mathbf{A}_t^\tau, \mathbf{o}_t)[\mathbf{q}] \in \mathbb{R}^4$  and the denoising vector field  $\mathbf{u}(\mathbf{A}_t^\tau | \mathbf{A}_t)[\mathbf{q}]$ , and a geodesic loss [12, 14] between a target quaternion  $\mathbf{q}_t^{\tau+\delta} \in \mathbb{S}^3$  computed from Eq. (1) at time  $t + \delta$ , and a quaternion prediction  $\mathbf{q}_{\theta,t}^{\tau+\delta} = \mathbf{q}_t^\tau \otimes \mathbf{q}_{\theta,t}^\delta \in \mathbb{S}^3$ , with  $\mathbf{q}_{\theta,t}^\delta \in \mathbb{S}^4$  computed by integrating  $\mathbf{v}_\theta[\mathbf{q}_t] \in \mathbb{R}^4$  over a small integration step  $\delta \sim \mathcal{U}(0.01, 1 - \tau)$ . The total loss is the sum of the translation and rotation loss

$$\mathcal{L}(\theta) = \mathbb{E}_{p(\mathbf{A}_t | \mathbf{o}_t), q(\mathbf{A}_t^\tau | \mathbf{A}_t)} [\mathcal{L}_{\mathbb{R}^3}(\theta) + \mathcal{L}_{\mathbb{S}^3}(\theta)]. \quad (3)$$

During inference, we generate actions by integrating the learned vector field from  $\tau = 0$  to  $\tau = 1$ , starting with random noise  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\mathbf{q}_0 \sim \mathcal{U}(\mathbb{S}^3)$  and using Euler integration in linear space for translations

$$\mathbf{x}_t^{\tau+\delta} = \mathbf{x}_t^\tau + \delta \mathbf{v}_\theta^\mathbf{x}(\mathbf{A}_t^\tau, \mathbf{o}_t), \quad (4)$$

and on the  $\mathbb{S}^3$  manifold for rotations

$$\mathbf{q}_t^{\tau+\delta} = \mathbf{q}_t^\tau \otimes \mathbf{q}_t^\delta (\mathbf{v}_\theta^\mathbf{q}(\mathbf{A}_t^\tau, \mathbf{o}_t)). \quad (5)$$

See Appendix A.2.3 for more details on flow matching.

**Image Resolution.** We select a resolution of  $280 \times 210$  for the main external camera and  $112 \times 112$  for the wrist camera, and resize images either by a central crop or padding. Importantly, unlike prior work [19], we do not distort the aspect ratio of the images by naive resizing as this also

distorts camera intrinsics and negatively affects depth and point cloud estimates. As wrist cameras contain less information than the external camera, we use a lower resolution, without losing important details, and reduce training and inference compute.

**Fine-tuning vision encoders.** As previously observed by ReVLA [10], robotics training can degrade the representations of the pre-trained vision encoders. We experiment with various configurations of vision encoder training and find the optimal setting to keep both SigLIP and MoGe vision encoders trainable during VLM training, but freeze MoGe in the VLA training stage.

**Control frequency & Data normalization.** We use an action chunk of size  $H = 5$  and frequency of 5Hz. For datasets not providing observations at 5Hz we resample the action targets via linear interpolation. We design data normalization to encourage learning motion across datasets, instead of “memorizing” each dataset separately. For target control normalization, we use global quantile normalization with statistics computed across the entire training mixture.

**Rotations.** We investigate various rotation representations including Euler angles, rotation matrices and unit quaternions. This is run in combination of using different rotation losses, including MSE or cosine for velocity predictions and geodesic and/or chordal loss [14] for integrated rotation predictions, as well as end-effector or robot base reference frames. We use Gram-Schmidt orthonormalization [12] to ensure valid rotation matrix predictions, but we find half-space unit quaternions to produce better results overall. We also find our proposed formulation on the manifold of unit quaternions  $\mathbb{S}^3 \rightarrow \mathbb{S}^3$  to be more stable and effective than linear flow matching  $\mathbb{R}^4 \rightarrow \mathbb{S}^3$ .

**Evaluation and Checkpointing.** We ablate all design choices by evaluating on the SIMPLER WidowX environments [22]. We set the same seed and enable deterministic CUDA operations for all VLA ablations to reduce training variance. We further resort to exponential moving average (EMA) checkpointing, which significantly stabilizes final checkpoint performance. For further details and ablations, see Appendix A.2.4.

## 4. Experimental evaluation

We evaluate the performance of SPEAR-1 as a generalist policy for robot manipulation and compare it to open-weights state-of-the-art VLA models. Our experiments aim to answer the following research questions:

1. Does 3D VLM pretraining improve the downstream VLA performance on robot control tasks?
2. How well does SPEAR-1 compare against state-of-the-art VLA models?

To answer these questions, we evaluate SPEAR-1 on a variety of manipulation tasks in simulation and multiple real-world environments on several robot embodiments.

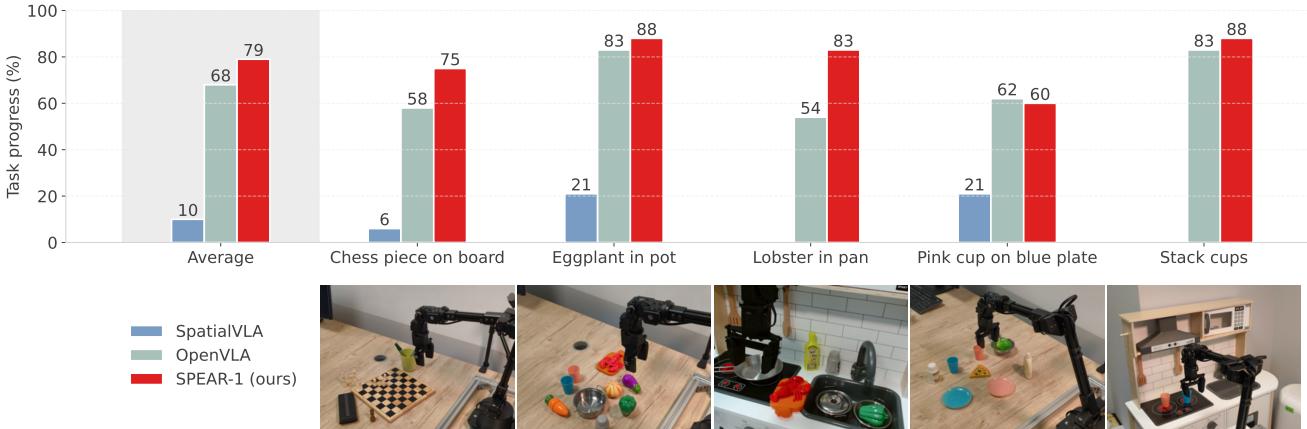


Figure 4. **Real world evaluation on WidowX.** SPEAR-1 is able to achieve 10% higher average task progress across all tasks than OpenVLA, a strong baseline in this setting. Bottom images correspond to the real-world tasks, whose performances are reported above.

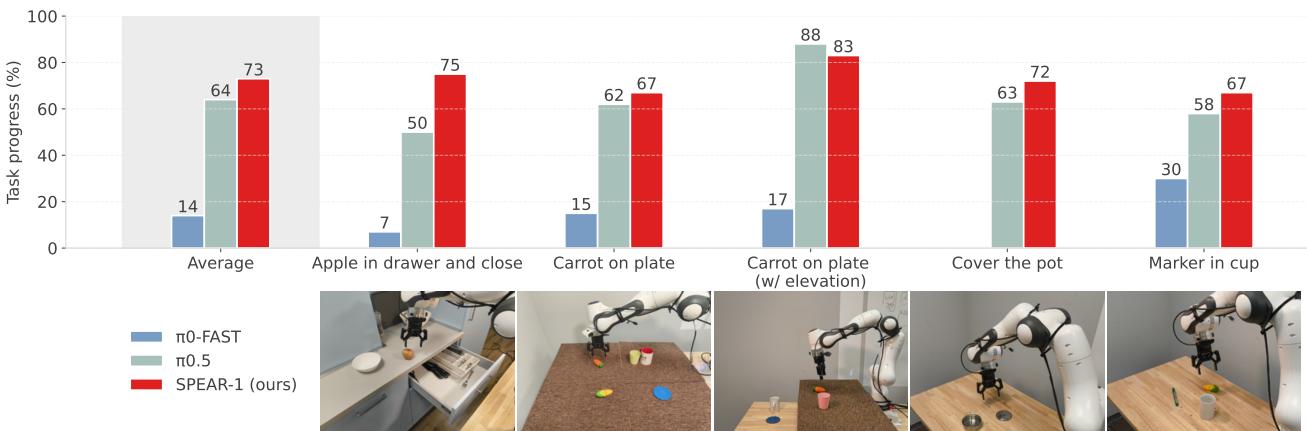


Figure 5. **Real world evaluation on Franka.** We find that without any fine-tuning on the target environment, SPEAR-1 noticeably outperforms  $\pi_0$ -FAST, and matches  $\pi_{0.5}$ , even though both baselines are trained on  $20\times$  more robotic data from significantly more diverse environments. The bottom row shows challenging, varied Franka environments where SPEAR-1 maintains strong zero-shot performance.

#### 4.1. Implementation details

**VLM training.** We train SPEAR-VLM with a batch size of 512 for 2k steps during the first stage and 10k steps for the second, for a total of 18hrs on 16 Nvidia H200 GPUs.

**VLA pre-training.** For VLA training, we start from SPEAR-VLM and randomly initialized action expert. We provide two camera views as inputs to the model: external, with resolution 280x210, and wrist, with resolution 112x112. When the wrist camera is not available, we feed a black image. We train on 32 H200 GPUs with batch size 2048 for 300k steps ( $\sim 6$  days) on a data mixture comprising 24 datasets (see Appendix A.2.1) from the Open X-Embodiment (OXE) collection [32].

**VLA post-training.** For WidowX real-world and SIMPLER simulation and Franka real-world experiments, we

additionally fine-tune our OXE pre-trained SPEAR-1 for 50k steps on the Bridge V2 [45] and DROID [18] datasets respectively. We refer to these versions as **SPEAR-1 (Bridge)** and **SPEAR-1 (DROID)** respectively.

#### 4.2. 3D ablations: SPEAR-VLM vs PaliGemma

We first evaluate whether 3D VLM pretraining improves VLA performance on downstream tasks and what design choices matter. Due to the cost of training on the entire OXE mixture (Tab. 4), we train only on specific subsets.

**SIMPLER WidowX experiments.** We perform an ablation study by training on a subset of the Bridge V2 [45] dataset, containing demonstrations only from a single kitchen sink environment, and evaluate the models in the SIMPLER [22] WidowX environments (see Appendix A.2.5 for details). This induces a distribution shift, which allows to demon-

Experiment	VLM Architecture	3D tasks	VLM training		VLA training		Put Carrot on Plate	Put Eggplant in Yellow Basket	Put Spoon on Towel	Stack Green Block on Yellow Block	Avg. Success Rate
			SigLIP	MoGe	SigLIP	MoGe					
1. no 3D	<b>PaliGemma</b>	None	train	–	train	–	25%	0%	54%	4.1%	20.8%
2. no OBJ	<b>SPEAR-VLM</b>	<b>points</b>	frozen	frozen	train	frozen	37.5%	0%	45.8%	4.1%	20.8%
3. no VLM-T	SPEAR-VLM	<b>objects</b>	frozen	frozen	train	frozen	41.7%	0%	70.8%	0%	29.1%
4. no VLA-MF	SPEAR-VLM	objects	frozen	frozen	train	<b>train</b>	29.1%	0%	41.7%	4.1%	18.8%
5. SPEAR-VLM	SPEAR-VLM	objects	<b>train</b>	<b>train</b>	train	frozen	50%	0%	79.1%	12.5%	<b>35.4%</b>

Table 1. **SPEAR-VLM 3D ablations.** Ablations on a single environment subset of Bridge V2 [45] to show the impact of 3D VLM pretraining. Without object-level 3D tasks (OBJ), 3D VLM pretraining does not show improvement over PaliGemma (1 vs. 2). Training MoGe during VLA training (VLA-MF) significantly degrades performance (3 vs. 4). Our study shows that the best training configuration is obtained when both SigLIP and MoGe are trained during VLM pretraining (5), followed by the frozen MoGe during VLA training.

Method	Carrot on Plate (Dist)	Carrot on Plate (Elev.)	Marker in Cup (Dist)	Avg. Task Progress
$\pi_0$ -PaliGemma (DROID)	0%	32%	<b>67%</b>	34%
$\pi_0$ -SPEAR-VLM (DROID)	<b>42%</b>	<b>52%</b>	43%	<b>46%</b>

Table 2. **SPEAR-VLM vs. PaliGemma for the downstream VLA tasks.** The experiments were conducted on the Franka (DROID) platform and the models were from trained from scratch on DROID. SPEAR-VLM achieves noticeable improvements. “Carrot on Plate” is not a part of DROID. This indicates SPEAR-VLM leads to better generalization.

strate the benefits of 3D pretraining when evaluating in unseen environments. In contrast, training on entire Bridge V2 leads to nearly the same performance for all models due to the close match between training and evaluations.

Results are reported in Tab. 1. First, we note that simply using SPEAR-VLM architecture and training without object-level prompts, but only 3D coordinates of random pixels (row 2), does not lead to any meaningful change in VLA performance over the baseline  $\pi_0$  model based on PaliGemma (row 1). However, training SPEAR-VLM on all 3D object-level tasks (Fig. 3), we observe a significant improvement in performance (row 3 and 5 vs. row 1). We also observe the importance of training SigLIP and MoGe encoders both during VLM and VLA training (row 3-5), with the best performance achieved when both are fine-tuned during VLM training and frozen MoGe during VLA training (row 5). We hypothesize this is because SigLIP has been trained only for image level semantics, while MoGe has been trained for dense and detailed depth prediction, which is much closer to the nature of robotic manipulation.

**Real-world Franka experiments.** To further validate the benefits of 3D VLM pretraining, we run comparisons by training on the DROID dataset [18]. Due to the higher cost, we train only 2 models: one initialized from the base PaliGemma VLM and the other from our 3D-aware SPEAR-VLM. We refer to the resulting models as  $\pi_0$ -PaliGemma (DROID) and  $\pi_0$ -SPEAR-VLM (DROID) respectively. We compare the performance of both VLAs on three of the four tasks from our Franka experiments (Section 4.4). The results are reported in Tab. 2. We can observe

that  $\pi_0$ -SPEAR-VLM (DROID) is able to outperform the baseline by more than 10% on average. We note that the task “Carrot on plate” is not a part of the DROID training dataset, thus shows the improved generalization capabilities of SPEAR-VLM. The lower scores of both models on the variations tabletop/elevations are likely due to workspace 3D position being out-of-distribution. Even in this case,  $\pi_0$ -SPEAR-VLM (DROID) is able to successfully complete the task in some cases while  $\pi_0$ -PaliGemma (DROID) fails every time.

### 4.3. Simulation experiments

We evaluate SPEAR-1 on the WidowX environments of the SIMPLER simulation benchmark [22], and compare it with OpenVLA [19] and SpatialVLA [35]. We report the results in Tab. 3. Our model is able to outperform the baselines by more than 10%. In our experience, we found SIMPLER simulation results only to be indicative of relative performance of the models on the real WidowX robot, but not necessarily of absolute performance. Therefore, we focus on real-world evaluations.

Model	Put Carrot on Plate	Put Eggplant in Yellow Basket	Put Spoon on Towel	Stack Green Block on Yellow Block	Avg. Success Rate
OpenVLA	0%	4.1%	0%	0%	1.0%
SpatialVLA	25.0%	<b>100.0%</b>	16.7%	29.2%	42.7%
SPEAR-1 (ours)	<b>58.3%</b>	62.5%	62.5%	<b>45.82%</b>	<b>57.3%</b>

Table 3. **SIMPLER [22] simulation evaluations.** SpatialVLA numbers are from [35]. SPEAR-1 outperforms the considered baselines by more than 10%.

### 4.4. Real-world experiments

We conduct evaluations on a total of 10 manipulation tasks across two robot platforms: WidowX and Franka Research 3. The tasks are designed to assess the ability of the evaluated models to generalize to unseen environments and objects. We design the tasks to be challenging for all models. For more details about the selected tasks see Appendix A.3.

**Evaluation protocol.** For each task we define M initial conditions by varying the starting position of the objects in the scene. We execute N trials for each initial condition, for a total of  $N \times M$  trials per task. For each model, we evaluate

and report the average task progress across all tasks, configurations, and trials. To that end, following previous works [3, 32], we define a scoring rubric with partial scoring for each task (see Appendix A.3 for scoring rubrics details).

**WidowX experiments.** Our hardware setup for this set of experiments closely matches the Bridge V2 setup [45], with a single external camera positioned on the side of the robot arm, pointing toward the workspace. For this set of experiments, 5 tasks are evaluated, with  $M = 4$ ,  $N = 3$ , for a total of 60 trials per model. We compare the performance of SPEAR-1 with OpenVLA [19], using the publicly released implementation and model weights. In this setting, we do not compare against  $\pi_0$  [5],  $\pi_0$ -FAST [33] and  $\pi_{0.5}$  [6] due to the lack of publicly accessible weights for the WidowX platform. The results are reported in Fig. 4. SPEAR-1 is able to achieve 10% higher average task progress across all tasks than OpenVLA, a very strong baseline in this setting.

**Franka experiments.** Our hardware setup for this set of experiments is similar to that of DROID [18]. We design 5 tasks, with  $M = 5$  and  $N = 3$ , for a total of 75 trials per model. We found that the wrist camera view is crucial for training and deployment on DROID. To ensure a fair comparison, we compare against open-weights models that use both the external and wrist camera. Specifically, we compare SPEAR-1 with the DROID-finetuned variants of  $\pi_0$ -FAST [5, 33], a strong autoregressive baseline, and  $\pi_{0.5}$  [6], one of the latest state-of-the-art robotic foundation model optimized for open-world generalization.

The results of our real-world experiments are reported in Fig. 5. Without any fine-tuning on the target environment, SPEAR-1 is able to significantly outperform  $\pi_0$ -FAST, and match  $\pi_{0.5}$ . We note that both baselines do not integrate any sort of specialized 3D-aware training and are trained on at least 900M more robot demonstration frames collected in diverse environments. In contrast, SPEAR-1 is trained on  $\sim$ 45M frames, approximately  $20\times$  less robotics data. These results indicate the importance of 3D-based knowledge and pretraining for VLA’s generalization. As an architecturally close comparison,  $\pi_0$ -FAST integrates a specialized action tokenization compared to  $\pi_0$  and was the first generalist policy trained on the DROID [18] to be successfully evaluated zero-shot in unseen environments, without fine-tuning. In comparison, SPEAR-1, which also follows the  $\pi_0$  architecture, can reach  $\sim 5\times$  higher performance than  $\pi_0$ -FAST without fine-tuning and without the large-scale robotic data used by  $\pi_0$ -FAST.

Apart from architectural enhancements and co-training on top of  $\pi_0$ -FAST,  $\pi_{0.5}$  integrates high-level semantic sub-task prediction and robotic data mixture explicitly focused on environment diversity and generalization. Qualitatively and quantitatively, we find  $\pi_{0.5}$  to perform better at environment generalization than  $\pi_0$ -FAST and match SPEAR-1’s performance on our set of evaluation tasks. This suggests

that 3D VLM pretraining on non-robotic data from diverse environments is a more scalable way to boost robotic models’ generalization capabilities without the need for large-scale robotic data collection in diverse environments.

## 5. Discussion and Limitations

As highlighted by our experimental evaluation, SPEAR-1 achieves state-of-the-art performance in multiple zero-shot robot control scenario, both in simulation and in the real-world. Nevertheless, our approach still presents several limitations. The proposed 3D VLM pre-training strategy is not well suited for deformable objects or objects with complex shapes.

Future work could explore the use of different 3D priors to better capture the geometry of such objects. Additionally, the coordinates of the 3D bounding boxes labels computed using MoGe’s affine-invariant depth predictions are not in metric space. Further investigation is required to analyze the implications of this design choice on downstream performance, as well as to explore how ground truth point cloud labels or state-of-the-art metric-depth estimators could be integrated to help resolve this limitation.

While we have showed the benefits of 3D VLM pre-training on downstream robot control tasks, the scaling laws relating the latter to the quantity and quality of 3D pre-training data are still not well understood. Due to resource and time constraints, we leave this investigation for future work. Another limitation of SPEAR-1 is the need to fine-tune on the target embodiment to achieve satisfactory results. We plan to explore how to alleviate this requirement in future work. It also remains to be seen how well SPEAR-1 generalizes to orders of magnitude more tasks and environments against models such as  $\pi_{0.5}$  trained on significantly more diverse robot data.

## 6. Conclusion

In this work we introduced SPEAR-1 and SPEAR-VLM that demonstrate a path towards building generalist robot policies from data beyond robot teleoperation only.

Our method targets the VLM backbone with **SPEAR-VLM**, a 3D-aware VLM trained on 2D images from non-robotic datasets enriched with 3D annotations. To embed control-relevant 3D knowledge in SPEAR-VLM, we train it on VQA questions, inspired by embodied tasks. Stepping on this foundation, we built **SPEAR-1**, a robotic foundation model that can be deployed robustly across multiple robot platforms and environments, and matches or outperforms state-of-the-art foundation models which have been trained on  $20\times$  more robot demonstrations.

Our work supports the hypothesis that enhancing VLM capabilities with non-robotic embodied knowledge is a scalable way to *reduce dependence on hard-to-collect robot*

demonstrations and build future robotic foundation models.

## Acknowledgments

Project Lead: Nikolay Nikolov, Project Manager: Jan-Nico Zaech, PI: Danda Pani Paudel, Luc Van Gool

We thank Alexander-Marc Spiridonov, Anna-Maria Hacheva and Yutong Hu for feedback and helpful technical discussions. We also thank Hristo Venev for engineering support and Kamen Pavlov for help with figures and visualizations.

## References

- [1] 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 Samangoeei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy 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. *arXiv preprint arXiv:2204.14198*, 2022. 3
- [2] Pranav Atreya, Karl Pertsch, Tony Lee, Moo Jin Kim, Arhan Jain, Artur Kuramshin, Clemens Eppner, Cyrus Neary, Edward Hu, Fabio Ramos, et al. Roboarena: Distributed real-world evaluation of generalist robot policies. In *Proceedings of the Conference on Robot Learning (CoRL 2025)*, 2025. 2, 16
- [3] Jose Barreiros, Andrew Beaulieu, Aditya Bhat, Rick Cory, Eric Cousineau, Hongkai Dai, Ching-Hsin Fang, Kunimatsu Hashimoto, Muhammad Zubair Irshad, Masha Itkina, et al. A careful examination of large behavior models for multitask dexterous manipulation. *arXiv preprint arXiv:2507.05331*, 2025. 8
- [4] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Koppula, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and Xiaohua Zhai. PaliGemma: A versatile 3B VLM for transfer. *arXiv preprint arXiv:2407.07726*, 2024. 2, 3
- [5] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolò Fusai, Lachy Groom, Karol Hausman, Brian Ichter, et al. pi\_0: A vision-language-action flow model for general robot control. *arXiv preprint arXiv:2410.24164*, 2024. 3, 4, 8, 14, 15
- [6] Kevin Black, Noah Brown, James Darpinian, Karan Dhabalnia, Danny Driess, Adnan Esmail, Michael Robert Equi, Chelsea Finn, Niccolò Fusai, Manuel Y. Galliker, Dibya Ghosh, Lachy Groom, Karol Hausman, brian ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky.  $\pi_{0.5}$ : a vision-language-action model with open-world generalization. In *9th Annual Conference on Robot Learning*, 2025. 1, 2, 3, 8, 15, 16
- [7] Wenxiao Cai, Iaroslav 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. 2
- [8] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 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, 2024. 2
- [9] Gheorghe Comanici, Eric Bieber, Mike Schaeckermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blissestein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025. 4
- [10] Sombit Dey, Jan-Nico Zaech, Nikolay Nikolov, Luc Van Gool, and Danda Pani Paudel. Revla: Reverting visual domain limitation of robotic foundation models. *arXiv preprint arXiv:2409.15250*, 2024. 3, 5
- [11] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 3
- [12] Andreas René Geist, Jonas Frey, Mikel Zobro, Anna Levina, and Georg Martius. Learning with 3d rotations, a hitchhiker’s guide to SO(3). In *Forty-first International Conference on Machine Learning*, 2024. 5
- [13] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first-and third-person perspectives. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19383–19400, 2024. 4, 13
- [14] Richard Hartley, Jochen Trumpf, Yuchao Dai, and Hongdong Li. Rotation averaging. *International journal of computer vision*, 103(3):267–305, 2013. 5
- [15] John Hewitt. Initializing new word embeddings for pretrained language models. <https://www.cs.columbia.edu/~johnhew/vocab-expansion.html>. 13
- [16] HuggingFace. Huggingface transformers documentation. [https://huggingface.co/docs/transformers/en/main\\_classes/model](https://huggingface.co/docs/transformers/en/main_classes/model). 13
- [17] Siddharth Karamcheti, Suraj Nair, Ashwin Balakrishna, Percy Liang, Thomas Kollar, and Dorsa Sadigh. Prismatic

- vlms: Investigating the design space of visually-conditioned language models. *arXiv preprint arXiv:2402.07865*, 2024. 2, 3
- [18] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. In *Robotics: Science and Systems*, 2024. 6, 7, 8, 15, 16
- [19] 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. 2, 3, 5, 7, 8, 15, 16
- [20] Jason Lee, Jiafei Duan, Haoquan Fang, Yuquan Deng, Shuo Liu, Boyang Li, Bohan Fang, Jieyu Zhang, Yi Ru Wang, Sangho Lee, et al. Molmoact: Action reasoning models that can reason in space. *arXiv preprint arXiv:2508.07917*, 2025. 3
- [21] Qixiu Li, Yaobo Liang, Zeyu Wang, Lin Luo, Xi Chen, Mozheng Liao, Fangyun Wei, Yu Deng, Sicheng Xu, Yizhong Zhang, et al. Cogact: A foundational vision-language-action model for synergizing cognition and action in robotic manipulation. *arXiv preprint arXiv:2411.19650*, 2024. 2, 15, 16
- [22] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, Sergey Levine, Jiajun Wu, Chelsea Finn, Hao Su, Quan Vuong, and Ted Xiao. Evaluating real-world robot manipulation policies in simulation. *arXiv preprint arXiv:2405.05941*, 2024. 5, 6, 7, 15
- [23] Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. 4
- [24] Yaron Lipman, Marton Havasi, Peter Holderith, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv preprint arXiv:2412.06264*, 2024. 4
- [25] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems*, pages 34892–34916. Curran Associates, Inc., 2023. 3, 4
- [26] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 2
- [27] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 3
- [28] Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint arXiv:2209.14577*, 2022. 4
- [29] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 13
- [30] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 12
- [31] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Charles Xu, Jianlan Luo, Tobias Kreiman, You Liang Tan, Lawrence Yunliang Chen, Pannag Sanketi, Quan Vuong, Ted Xiao, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In *Proceedings of Robotics: Science and Systems*, Delft, Netherlands, 2024. 15
- [32] Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023. 3, 6, 8, 15
- [33] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025. 1, 2, 3, 8, 15, 16
- [34] Sundar Pichai, Demis Hassabis, and Koray Kavukcuoglu. Introducing gemini 2.0: our new ai model for the agentic era. <https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/>. 3
- [35] 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. 2, 3, 7, 12, 15, 16
- [36] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024. 4, 13
- [37] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kun-chang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024. 4
- [38] Moritz Reuss, Hongyi Zhou, Marcel Rühle, Ömer Erdinç Yağmurlu, Fabian Otto, and Rudolf Lioutikov. FLOWER: Democratizing generalist robot policies with efficient vision-language-action flow policies. In *7th Robot Learning Workshop: Towards Robots with Human-Level Abilities*, 2025. 15, 16
- [39] Chan Hee Song, Valts Blukis, Jonathan Tremblay, Stephen Tyree, Yu Su, and Stan Birchfield. Robospatial: Teaching spatial understanding to 2d and 3d vision-language models for robotics. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 15768–15780, 2025. 2
- [40] Alexander Spiridonov, Jan-Nico Zaech, Nikolay Nikolov, Luc Van Gool, and Danda Pani Paudel. Generalist robot ma-

- nipulation beyond action labeled data. In *9th Annual Conference on Robot Learning*, 2025. 2, 15
- [41] Andreas Steiner, André Susano Pinto, Michael Tschannen, Daniel Keysers, Xiao Wang, Yonatan Bitton, Alexey Gritsenko, Matthias Minderer, Anthony Sherbondy, Shangbang Long, et al. Paligemma 2: A family of versatile vlms for transfer. *arXiv preprint arXiv:2412.03555*, 2024. 2
- [42] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024. 3, 14
- [43] Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025. 2
- [44] Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montserrat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna, Robert Baruch, Maria Bauza, Michiel Blokzijl, et al. Gemini robotics: Bringing ai into the physical world. *arXiv preprint arXiv:2503.20020*, 2025. 3
- [45] Homer Walke, Kevin Black, Abraham Lee, Moo Jin Kim, Max Du, Chongyi Zheng, Tony Zhao, Philippe Hansen-Estruch, Quan Vuong, Andre He, Vivek Myers, Kuan Fang, Chelsea Finn, and Sergey Levine. Bridgedata v2: A dataset for robot learning at scale. In *Proceedings of the Conference on Robot Learning (CoRL)*, 2023. 4, 6, 7, 8, 13, 14, 15
- [46] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 2
- [47] Ruicheng Wang, Sicheng Xu, Cassie Dai, Jianfeng Xiang, Yu Deng, Xin Tong, and Jiaolong Yang. Moge: Unlocking accurate monocular geometry estimation for open-domain images with optimal training supervision. *arXiv preprint arXiv:2410.19115*, 2024. 2, 3, 4, 13
- [48] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 3
- [49] Michał Zawalski, William Chen, Karl Pertsch, Oier Mees, Chelsea Finn, and Sergey Levine. Robotic control via embodied chain-of-thought reasoning. In *8th Annual Conference on Robot Learning*, 2024. 15
- [50] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023. 3
- [51] Qingqing Zhao, Yao Lu, Moo Jin Kim, Zipeng Fu, Zhuoyang Zhang, Yecheng Wu, Zhaoshuo Li, Qianli Ma, Song Han, Chelsea Finn, et al. Cot-vla: Visual chain-of-thought reasoning for vision-language-action models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1702–1713, 2025. 15, 16
- [52] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing. *arXiv:1801.09847*, 2018. 13
- [53] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, Quan Vuong, Vincent Vanhoucke, Huong Tran, Radu Soricu, Anikait Singh, Jaspiar Singh, Pierre Sermanet, Pannag R. Sanketi, Grecia Salazar, Michael S. Ryoo, Krista Reymann, Kanishka Rao, Karl Pertsch, Igor Mordatch, Henryk Michalewski, Yao Lu, Sergey Levine, Lisa Lee, Tsang-Wei Edward Lee, Isabel Leal, Yuheng Kuang, Dmitry Kalashnikov, Ryan Julian, Nikhil J. Joshi, Alex Irpan, Brian Ichter, Jasmine Hsu, Alexander Herzog, Karol Hausman, Keerthana Gopalakrishnan, Chuyuan Fu, Pete Florence, Chelsea Finn, Kumar Avinava Dubey, Danny Driess, Tianli Ding, Krzysztof Marcin Choromanski, Xi Chen, Yevgen Chebotar, Justice Carbalal, Noah Brown, Anthony Brohan, Montserrat Gonzalez Arenas, and Kehang Han. Rt-2: Vision-language-action models transfer web knowledge to robotic control. In *Proceedings of The 7th Conference on Robot Learning*, pages 2165–2183. PMLR, 2023. 3

## A. Appendix

The appendix is organized as follows:

- In Sec. A.1 we provide more details on the VLM pre-training including VQA tasks, encoder fusion strategies, 3D tokenization and data annotation pipeline.
- In Sec. A.2 we provide more details on the VLA training including data mixture, architecture, flow matching and design decision ablation results.
- In Sec. A.3 we provide the scoring rubrics for real-world evaluation tasks
- In Sec. A.4 we discuss the differences between Bridge V2 and DROID datasets for zero-shot control evaluations in unseen environments.

Dataset	Weight
austin_buds_dataset	0.5
austin_sailor_dataset	2.0
austin_sirius_dataset	0.5
berkeley_autolab_ur5	1.0
berkeley_cable_routing	0.1
berkeley_fanuc_manipulation	1.0
bridge	18.0
dlr_edan_shared_control	0.1
droid	35.0
fmb	1.5
fractal20220817.data	12.0
furniture_bench_dataset	1.5
iamlab_cmu_pickup_insert	0.3
kuka	4.0
language_table	1.5
nyu_franka_play_dataset	0.3
roboset (kinesthetic)	2.0
roboset (teleoperation)	5.0
roboturk	3.0
stanford_hydra_dataset	3.0
taco_play	2.0
toto	1.5
ucsd_kitchen_dataset	0.2
utaustin_mutex	3.0
viola	1.0

Table 4. Open X-Embodiment data mixture for SPEAR-1 pre-training

## A.1. VLM training

### A.1.1. VQA tasks for VLM pre-training

The Visual Question Answering (VQA) tasks used during VLM pre-training are inspired by VLA embodied tasks and aim to embed as much control-relevant 3D information into the VLM as possible. We use templated question-answer pairs grouped in the following categories:

- **3D keypoints prediction:** Output the 3D coordinates of the closest, furthest and center points of an object with

respect to the camera frame.

- **3D bounding prediction:** Output the vertices of the 3D bounding box of an object.
- **Object-to-object distance prediction:** Output the direct distance between object X and object Y in 3D space as well as its  $xyz$  components.
- **Object-to-object bounding box prediction:** Output the distance between the bounding box vertices and the centers of object X and object Y.
- **Backprojection:** Locate the vertices of the 3D bounding box of an object on the 2D image.
- **Chain-of-thought comparisons:** What is the distance from the camera to object X? What is the distance from the camera to object Y? Which object is closer to the camera?

To further encourage the model to ‘reason’ over the information provided and attend to the right objects, in a single training example we use a random number (between 1 and 4) of question-answer pairs corresponding to different prompts and objects in the scene. To resolve ambiguities, if two instances of the same type of object appear in the image, we filter them out and never ask questions about them.

### A.1.2. VLM encoder fusion strategies

We experimented with 2 different strategies to combine the outputs of the SigLIP and MoGe encoders:

1. Concatenating the visual features predicted by both encoders and projecting them via a linear layer to the LLM embedding space. In particular, for SigLIP we take only the tokens at the last layer of the vision encoder, while for MoGe we take the tokens at the last 4 layers of the encoder, following the approach used by MoGe architecture to decode the features to a 3D point cloud.
2. Using MoGe’s predicted 3D point cloud  $\mathbf{P}$  in the camera ego pose (in an affine-invariant space) and adding them to the SigLIP encoder features, similar to SpatialVLA [35]. In particular, MoGe’s 3D point cloud output  $\mathbf{P} \in \mathbb{R}^{H \times W \times 3}$  is embedded to  $\mathbf{P}' \in \mathbb{R}^{h \times w \times d}$  through a projector  $\psi(\cdot)$ , composed of normalization, convolution, sinusoidal embedding  $\gamma(x) = (x, \sin(2^0\pi x), \cos(2^0\pi x), \dots, \sin(2^{L-1}\pi x), \cos(2^{L-1}\pi x))$  [30] and an MLP. Finally, the features  $\mathbf{F}' = \mathbf{F} + \mathbf{P}'$  are fed to PaliGemma’s SigLIP linear projector, where  $\mathbf{F} \in \mathbb{R}^{h \times w \times d}$  denotes the features at the SigLIP encoder output.

During our preliminary VLM evaluations we found the first strategy to demonstrate qualitatively better performance on bounding box prediction tasks. In particular, models trained with the second approach struggled to consistently output “grammatically” correct bounding boxes, e.g. they would output 22 or 23 3D tokens instead of the required 24. We therefore used the first approach for all VLM pre-training experiments in the main paper.

<b>Dataset</b>	<b>Domain / Subset</b>	<b># Annotated Images</b>	<b>Segmentation Masks</b>
EgoExo4D [13]	Cooking & Bike Repair	~200k	GT
Bridge [45]	Robot Demonstrations	~30k	SAM2 Generated
<b>Total</b>		~230k	

Table 5. Annotated image counts for training dataset construction, with segmentation mask availability.

<b>Task</b>	<b>0.25</b>	<b>0.50</b>	<b>0.75</b>	<b>1.00</b>
Carrot on Plate (w/ distractors & elevations)	Reach carrot	Pick up carrot	Drop on/near plate	Correctly place on plate
Marker in Cup (w/ distractors)	Reach marker	Pick up marker	Drop on/near cup	Place inside cup
Cover the Pot	-	Pick up lid	Drop lid on pot	Correctly cover pot
Apple in drawer and close	Pick up the apple	Put the apple in the drawer	Half-close the drawer	Fully close the drawer

Table 6. Scoring rubric for Franka evaluation tasks.

<b>Task</b>	<b>0.25</b>	<b>0.50</b>	<b>0.75</b>	<b>1.00</b>
Eggplant in pot	Reach the eggplant	Pick up the eggplant	Drop the eggplant near the pot	Drop the eggplant on the pot
Pink cup on blue plate	Reach the pink cup	Pick up the pink cup	Drop the pink cup near the plate	Place the pink cup correctly
Chess piece on board	-	Go to the brown chess piece	Pick up the brown chess piece	Drop it on the board
Lobster in the pan	-	Pick up the lobster	Drop the lobster near the pan	Place the lobster inside the pan
Corn between cups	-	Pick up the corn	-	Place the corn between the cups

Table 7. Scoring rubric for WidowX evaluation tasks.

### A.1.3. 3D tokenization

To encode 3D information into text we extend the PaliGemma tokenizer with  $N = 1024$  3D tokens, as 3D coordinates are conceptually different from the existing visual and language tokens. This is in line with PaliGemma’s approach of extending Gemma’s tokenizer to pixel locations. Each 3D token corresponds to a quantized *distance value* in the range  $[z_{\min}, z_{\max}]$ , where  $z_{\min}$  and  $z_{\max}$  are computed as the 1st and 99th quantiles of the 3D point cloud distribution along any of the  $xyz$  coordinates.

We found the *distance values* in the data to approximately follow a Normal distribution. Therefore, to allow for more accurate tokenization, we compute non-uniform bins with fine-grained discretization around the mean and spread out widths near the tails such that the distribution of 3D tokens is approximately uniform.

We initialize the new token embedding weights from a multivariate normal distribution that has the mean and covariance of the pretrained embeddings [15, 16].

### A.1.4. VQA data annotation pipeline

We follow the method described in Section 3.1 in order to enrich 2D images with semantics, segmentation masks and 3D point clouds. We also experimented with GroundingDINO [29] instead of Gemini, but we found the semantic

labels produced by GroundingDINO to be a lot less accurate and consistent. We found that prompting SAM2 [36] with 2D bounding boxes near the target objects, leads to segmentation masks of high quality.

We also found that MoGe [47] outputs depths at different scales depending on the input image size. Therefore, we resize all our images to 840x630 for MoGe point cloud annotations.

For 3D bounding box estimation, after filtering the 3D point cloud with a segmentation mask, we run statistical outlier removal and estimate an oriented 3D bounding box around the remaining points using Open3D [52]. To facilitate learning, we order all 8 bounding box vertices in a consistent way, starting based on their spatial coordinates with respect to the camera frame.

## A.2. VLA training

### A.2.1. Data mixture

We report the VLA training data mixture in Tab. 4. The sampling weights are chosen manually based on dataset size, visual and task diversity, and quality of language annotations.



Figure 6. **3D ablation environments on WidowX.** (a) Training data subset from Bridge V2 [45]. (b) - (e) SIMPLER evaluation environments.

### A.2.2. VLA training details

**Reference Frames.** In this work we focus on learning position control of fixed-base single-arm manipulators. Each control in the sequence  $\mathbf{A}_t = [\mathbf{a}_t, \mathbf{a}_{t+1}, \dots, \mathbf{a}_{t+H-1}]$  is defined as a delta with respect to the current end-effector cartesian pose  $\Delta_{EE} = [\Delta_T, \Delta_R]$ . The translation component,  $\Delta_T$  is in robot base frame and the rotation component,  $\Delta_R$ , is in end-effector frame. The gripper action is binary.

**Action Chunking.** During VLA training we use an action chunk of size  $H = 5$  and frequency of 5Hz. As not all datasets in Open X-Embodiment provide action labels at 5Hz, we downsample or upsample the actions accordingly via linear interpolation. This is done with the goal to encourage the model to share knowledge across datasets with different control frequencies and embodiments instead of ‘memorizing’ each dataset separately.

**Architecture.** Similar to  $\pi_0$  [5], SPEAR-1 combines a VLM, which processes the image-language inputs, with an *action expert* module, which processes robot proprioception observations and predicts the robot action sequence conditioned on the VLM transformer’s intermediate key-value pairs. The action expert has the same architecture and number of layers as the Gemma [42] transformer and configuration downsized to *token\_size* = 4096, *hidden\_size* = 4096, for a total of  $\sim 300M$  parameters, which is exactly the same as  $\pi_0$  [5]. Corresponding layers in the VLM and the action expert have a shared attention operation with block-wise causal attention over the blocks  $[\mathbf{l}_t, \mathbf{l}_t], [\mathbf{p}_t], [\mathbf{a}_{t+1}, \dots, \mathbf{a}_{t+H-1}]$ . Within each block, there is full bidirectional attention and the tokens in each block can attend to tokens in previous blocks, but cannot attend to the tokens in future blocks. During training, only the action sequence prediction is supervised and gradient updates are propagated back to the VLM parameters through the shared attention layers.

### A.2.3. Flow matching details

To address the inherent double coverage of 3D rotations by the unit quaternion group  $\mathbb{S}^3$ , we ensure that all quaternions used during training and inference lie in the same half-space defined by  $\Re(\mathbf{q}) = \mathbf{q}_w > 0$ .

**Quaternion integration.** Given a unit quaternion  $\mathbf{q}_t \in$

$\mathbb{S}^3$  and its time derivative  $\dot{\mathbf{q}}_t \in \mathbb{R}^4$ , we can compute the angular velocity of rotation via  $\omega_t = 2.0 \cdot \Im(\mathbf{q}_t^* \otimes \dot{\mathbf{q}}_t) \in \mathbb{R}^3$ . For a small time step  $\Delta t$ , the corresponding delta rotation is given by a rotation vector around the unit axis  $\omega = \omega / \|\omega\|$  over an angle  $\Delta\phi = \|\omega\|\Delta t$ . The corresponding delta quaternion is given by

$$\Delta \mathbf{q} = \left[ \cos\left(\frac{\Delta\phi}{2}\right), \omega \sin\left(\frac{\Delta\phi}{2}\right) \right]. \quad (6)$$

The integrated unit quaternion is then given by  $\mathbf{q}_{t+\Delta t} = \mathbf{q}_t \otimes \Delta \mathbf{q} \in \mathbb{S}^3$ ,

**Rotation losses.** The denoising vector field for quaternions  $\mathbf{u}_t(\mathbf{q}_t^\tau | \mathbf{q}_t) \in \mathbb{R}^4$  is computed as:

$$\begin{aligned} \mathbf{u}_t(\mathbf{q}_t^\tau | \mathbf{q}_t) &= \frac{d\mathbf{q}_t^\tau}{d\tau} = \\ &= \frac{\theta}{\sin \theta} \left[ -\cos((1-\tau)\theta) \mathbf{q}_\epsilon + \cos(\tau\theta) \mathbf{q}_t \right]. \end{aligned} \quad (7)$$

The cosine loss is applied directly on the velocity predictions and has the form:

$$\mathcal{L}_t^{\cos}(\theta) = 1 - \mathbf{v}_\theta(\mathbf{A}_t^\tau, \mathbf{o}_t)[\mathbf{q}] \cdot \mathbf{u}(\mathbf{A}_t^\tau | \mathbf{A}_t)[\mathbf{q}]. \quad (8)$$

The geodesic loss is applied on an integrated rotation prediction  $\mathbf{q}_{\theta,t}^{\tau+\delta} \in \mathbb{S}^3$ , derived by integrating the noised input quaternion  $\mathbf{q}_t^\tau$  with the predicted rotation velocity  $\mathbf{v}_\theta(\mathbf{A}_t^\tau, \mathbf{o}_t)[\mathbf{q}]$  over a small time step  $\delta$ . We follow the integration method described above. The target is given by the ground truth interpolated quaternion at time  $t + \delta$ , denoted as  $\mathbf{q}_t^{\tau+\delta} \in \mathbb{S}^3$ . The closed form geodesic loss is given by:

$$\mathcal{L}_t^{\text{geo}}(\theta) = \min |\mathbf{q}_t^{\tau+\delta} \pm \mathbf{q}_{\theta,t}^{\tau+\delta}|. \quad (9)$$

The complete rotation loss is given by:

$$\mathcal{L}_{\mathbb{S}^3}(\theta) = \sum_{k=t}^{t+H} [\mathcal{L}_k^{\cos}(\theta) + \mathcal{L}_k^{\text{geo}}(\theta)]. \quad (10)$$

### A.2.4. VLA design decisions details

We present more details and results on the design choices we explored for VLA training.

Experiment	Put Carrot on Plate	Put Eggplant in Yellow Basket	Put Spoon on Towel	Stack Green Block on Yellow Block	Avg. Success Rate
224×224	<b>70.8%</b>	70.8%	79.1%	8.3%	57.25%
280×210	62.5%	<b>75.0%</b>	<b>83.3%</b>	<b>12.5%</b>	<b>58.3%</b>

Table 8. **Image resolution ablations.** Different resolutions lead to comparable results on SIMPLER WidowX tasks.

Experiment	Put Carrot on Plate	Put Eggplant in Yellow Basket	Put Spoon on Towel	Stack Green Block on Yellow Block	Avg. Success Rate
trainable SigLIP	<b>75.0%</b>	<b>100.0%</b>	79.1%	<b>37.5%</b>	<b>72.9%</b>
frozen SigLIP	62.5%	54.1%	83.3%	25%	56.3%
frozen-trainable SigLIP	66.6%	83.3%	<b>100.0%</b>	33.3%	70.8%
lower lr SigLIP	58.3%	58.3%	79.1%	29.1%	56.3%

Table 9. **Vision encoder training.** Trainable SigLIP outperforms other strategies on SIMPLER WidowX tasks. Frozen SigLIP followed by switching on gradients is comparable.

Experiment	Put Carrot on Plate	Put Eggplant in Yellow Basket	Put Spoon on Towel	Stack Green Block on Yellow Block	Avg. Success Rate
99-th quantile	54.1%	79.1%	79.1%	<b>33.3%</b>	61.5%
min-max const	<b>66.6%</b>	<b>87.5%</b>	<b>87.5%</b>	20.8%	<b>65.6%</b>
mean-std	45.8%	79.1%	45.8%	25.0%	49.0%

Table 10. **Translation controls normalization.** Normalizing translation controls with min-max constants outperforms other strategies on SIMPLER WidowX tasks.

Flow matching	Velocity Loss	Rotation loss	Put Carrot on Plate	Put Eggplant in Yellow Basket	Put Spoon on Towel	Stack Green Block on Yellow Block	Avg. Success Rate
linear	MSE	geodesic	41.6%	<b>100.0%</b>	41.6%	16.6%	50.0%
linear	cos	geodesic	41.6%	87.5%	50.0%	29.1%	52.1%
$\mathbb{S}^3$	MSE	geodesic	<b>45.8%</b>	62.5%	<b>75.0%</b>	<b>45.8%</b>	57.3%
$\mathbb{S}^3$	cos	geodesic	<b>45.8%</b>	79.1%	66.6%	<b>45.8%</b>	<b>59.4%</b>

Table 11. **Linear vs  $\mathbb{S}^3$  Flow Matching for rotations.**  $\mathbb{S}^3$  flow matching consistently outperforms linear flow matching on SIMPLER WidowX tasks.

**Image Resolution.** Image resolution ablations are presented in Tab. 8. We observe that square vs 4:3 aspect ratio does not significantly affect performance.

**Fine-tuning vision encoders.** Ablations on fine-tuning vision encoders are presented in Tab. 9. Trainable SigLIP strongly outperforms a frozen SigLIP as well as SigLIP with a lower learning rate compared to the rest of the weights. Freezing SigLIP and fine-tuning for additional 2k steps (frozen-trainable SigLIP) leads to comparable performance to trainable SigLIP, but requires an additional hyperparameter tuning.

**Controls normalization.** Ablations on translation controls normalization are presented in Tab. 10. Mean-std normalization is significantly worse than other forms of normalization. Min-max normalization with const values is slightly better than per-dataset min-max normalization with 1st and 99th quantiles.

**Rotations.** Partial ablations on rotation representations are presented in Tab. 11.  $\mathbb{S}^3$  flow matching consistently outperforms linear flow matching. Cosine distance leads to slightly better performance than MSE for rotation velocity prediction.

### A.2.5. 3D ablation details

**SIMPLER WidowX experiments.** For SIMPLER [22] 3D ablations on WidowX, we train on a subset of the Bridge V2 [45] dataset, containing demonstrations only from a single kitchen sink environment. The resulting subset comprises  $\sim 41\%$  of the original Bridge V2 dataset. We train each VLA for 30k steps with batch size 512. Example images from the training and evaluation environments are shown in Fig. 6.

**Franka experiments.** For the 3D ablations on a real-world Franka robot, we train on the entire DROID [18] dataset for 100k steps with batch size 2048. Both models take as input both side and wrist cameras.

### A.3. Real-world robot task description and scoring

We provide the detailed task progression scoring for all real-world evaluations on the WidowX and Franka in Tab. 6 and Tab. 7 respectively.

### A.4. Zero-shot control: Bridge V2 vs. DROID

Model	Zero-shot control embodiments in real-world <b>unseen</b> environment
<b>RT-1-X</b> [32]	WidowX
<b>RT-2-X</b> [32]	WidowX, Google Robot
<b>Octo</b> [31]	WidowX, Google Robot?
<b>OpenVLA</b> [19]	WidowX
<b>SpatialVLA</b> [35]	WidowX
<b>CogACT</b> [21]	WidowX
<b>FLOWER</b> [38]	WidowX
<b>MotoVLA</b> [40]	WidowX
<b>CoT-VLA</b> [51]	WidowX
$\pi_0$ [5]	<b>Franka</b> , Others?
$\pi_0$ -FAST [33]	<b>Franka</b> , Others?
$\pi_{0.5}$ [6]	<b>Franka</b> , Mobile Fibocom, Mobile Galaxea, Others?
<b>Gemini Robotics 1.5</b>	Bimanual Franka, Aloha, Apollo humanoid
<b>RDT1</b>	Bimanual UR5, Aloha
<b>RDT2</b>	Bimanual UR5, Bimanual Franka
<b>SmolVLA</b>	S0101
<b>GROOT-N1.5</b>	None
<b>SPEAR-1 (ours)</b>	WidowX, <b>Franka</b>

Table 12. Most works on generalist models for robot manipulation evaluate zero-shot control on Bridge V2 + WidowX using in-distribution environments. Only few do so on DROID + Franka in **unseen** environments.

Most works on generalist models for robot manipulation [19, 21, 32, 49, 51] evaluate the zero-shot control capabilities of their policies by pretraining on the Bridge V2 dataset [45] and deploying on the WidowX robot in environments close to the training distribution. Bridge V2, however, is not very diverse in the number of environments, objects,

and camera viewpoints. As a result, we observe that models pre-trained on Bridge V2 only perform well on WidowX environments when the deployment scenario is similar to what is seen in the dataset (e.g. in the blue toy sink environment), but are usually very sensitive to variations in camera position and OOD backgrounds and objects. In addition, the WidowX arm has a very low payload and short reach, which makes it unable to manipulate objects beyond the items in a toy kitchen set. The DROID dataset [18], on the other hand, is significantly more diverse, and the Franka arm used for data collection is more capable. Furthermore, DROID demonstrations are collected primarily in real-world environments instead of toy environments, the number of unique scenes is  $20\times$  higher, and the camera viewpoints vary significantly. Therefore, we posit that pretraining on DROID and deploying on Franka is a superior experimental setup to showcase generalization to more realistic real-world scenarios, as shown by [2]. The diversity and richness of DROID, however, is at the same time a challenge. Training generalist control policies on DROID that perform well zero-shot on a Franka robot in unseen environment, is a task that, to the best of our knowledge, has been tackled successfully only by a handful of works so far [6, 33]. In contrast, multiple other works that pre-train on DROID, resort to mixing or fine-tuning on demonstrations collected from the specific target environment in order to achieve good performance [19, 21, 35, 38, 51]. Therefore, as suggested also by [33], we argue that achieving state-of-the-art performance on zero-shot control on the DROID setup by pre-training on DROID is a significantly stronger result than pre-training on Bridge V2 and deploying on WidowX.