

OG-VLA: 3D-Aware Vision Language Action Model via Orthographic Image Generation

Ishika Singh¹, Ankit Goyal², Stan Birchfield², Dieter Fox², Animesh Garg³, Valts Blukis²

¹University of Southern California, ²NVIDIA, ³Georgia Institute of Technology

Abstract: We introduce OG-VLA, a novel architecture and learning framework that combines the generalization strengths of Vision Language Action models (VLAs) with the robustness of 3D-aware policies. We address the challenge of mapping natural language instructions and multi-view RGBD observations to quasi-static robot actions. 3D-aware robot policies achieve state-of-the-art performance on precise robot manipulation tasks, but struggle with generalization to unseen instructions, scenes, and objects. On the other hand, VLAs excel at generalizing across instructions and scenes, but can be sensitive to camera and robot pose variations. We leverage prior knowledge embedded in language and vision foundation models to improve generalization of 3D-aware keyframe policies. OG-VLA projects input observations from diverse views into a point cloud which is then rendered from canonical orthographic views, ensuring input view invariance and consistency between input and output spaces. These canonical views are processed with a vision backbone, a Large Language Model (LLM), and an image diffusion model to generate images that encode the next position and orientation of the end-effector on the input scene. Evaluations on the ARNOLD and COLOSSEUM benchmarks demonstrate state-of-the-art generalization to unseen environments, with over 40% relative improvements while maintaining robust performance in seen settings. We also show real-world adaption in 3 to 5 demonstrations along with strong generalization. Videos and resources at <https://og-vla.github.io>

1 Introduction

We study the problem of mapping natural language instructions and multiple posed RGBD observations to robot actions, with specific focus on quasi-static manipulation tasks that can be decomposed into a sequence of end-effector keyframes. This category encompasses a wide variety of tasks such as pick-and-place, opening/closing doors and containers, manipulating buttons, valves, switches, and more. Building robust policies that solve such tasks in *unseen* environments remains an open challenge that could enable numerous industrial and household applications, from cleaning and sorting robots to machine tending.

VLAs have recently demonstrated successful generalization to concepts unseen in the robotics training data [1, 2, 3, 4], such as manipulating novel objects based on language instructions. While achieving generalization breakthroughs, they require massive training datasets [3, 5, 4] and typically accept a single RGB view input. As a result, despite the large amount of training data, the resulting systems remain sensitive to variations in camera and robot poses, which hurt their adaptability to new applications [6]. They also lack explicit visual reasoning when predicting actions (often as LLM tokens), constraining their ability to perform precise, generalizable 3D spatial reasoning.

In contrast, 3D-aware keyframe based policies learn effectively from few demonstrations and generalize well to novel camera poses and object placements [7, 8, 9, 10] as confirmed in the camera perturbation evaluation study by Pumacay et al. [11]. This success stems from a 3D scene representation within the model, such as a voxel map [7], canonical orthographic views [8], or point

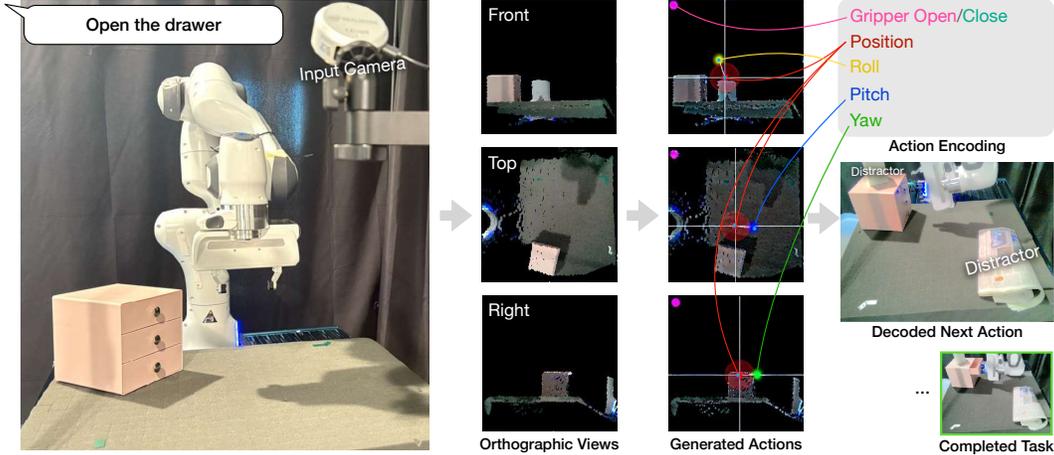


Figure 1: OG-VLA illustration. OG-VLA represents robot end-effector keyframes with easy-to-decode annotations on orthographic images in a set of canonical views. This output encoding enables action prediction via image generation, and using canonical views achieves invariance to input camera poses. The red hotspot is the predicted end-effector position in each image. In this example, the hotspot is indicating the 3D point for approaching the drawer handle to open it. The yellow, blue and green hotspots work in tandem with the red hotspot to encode the three axes of end-effector orientation. The color of the hotspot on the top-left encodes the gripper open/close state. The system is robust to distractors and changing lighting conditions.

clouds [10]. However, unlike VLAs, these systems overfit to the training scenes and objects, failing to accept instructions that refer to new, previously unseen objects.

We propose OG-VLA (Orthographic-image Generation Vision-Language-Action model), a novel robot policy architecture that combines the generalization strengths of VLAs with robustness of 3D-aware policies. OG-VLA uses LLMs and image generation to map posed RGBD images and language to 6-DOF end-effector pose keyframes in a sequence. The system comprises four components: a point cloud renderer that renders a scene reconstruction to canonical orthographic views, a vision backbone that encodes these views into visual embeddings, an LLM that predicts action tokens, and an image diffusion model that decodes these action tokens to predict actions on each of the orthographic views through image generation, which we decode to final 3D poses. The LLM and image diffusion models are trained end-to-end, so that they work together to produce consistent and precise predictions required for robotic manipulation. Figure 1 illustrates our method.

We conduct simulation and real-robot experiments. In simulation, the ARNOLD [12] benchmark tests generalization to unseen objects and environments, while COLOSSEUM [11] evaluates robustness to variations in camera poses, object poses, colors, and distractors. We show significantly improved performance on generalization tests of both benchmarks, achieving state-of-the-art on ARNOLD [12]. In real-world experiments, we show our method’s ability to learn manipulation tasks from as few as 3 to 5 demonstrations, highlighting its suitability for kinesthetic teaching and rapid adaptation to new domains. We also present a detailed study of our model architecture design choices, providing insights into its functioning as well as potential future improvements.

2 Related Work

Large Language Models have seen an explosion in research on their use-cases in robotics, such as task planning [13, 14, 15], reward generation for reinforcement learning [16, 17], and interfacing with vision models [18]. However, although these methods generalize at a high level, they assumed access to low-level skills such pick, place, open, and close, along with perception and scene representation systems. These skills are hard to build, and their development is in line with the goals of this work.

3D-aware keyframe-based multi-task policies have shown the ability to learn complex manipulation behaviors from as little as ten demonstrations per task. Examples of these works include PerAct [7] based on voxel grids, RVT [8, 9] based on orthonormal views, and Act3D [10] based visual feature point clouds as the scene representation. These systems significantly improve upon image-based policies [19, 20] in the amount of data they require and their robustness to new object placements at test-time. However, these systems have been trained from scratch on specific tasks, robots, and environments, and struggle to generalize to new objects and scenes, or instructions.

Vision Language Action models (VLAs) leverage large-scale prior knowledge from LLMs and vision foundation models to create robot policies that generalize to new concepts and objects at test time [1, 21, 2, 3, 22, 4]. However, these systems require huge amount of demonstration data to train. For example, Kim et al. [3] use over 900k demos from the Open-X Embodiment dataset [23]. Despite the large quantity of data, the resulting systems are sensitive to changes in the 3D environment, such as different camera poses relative to the robot system. Our system trains with significantly less data to achieve state-of-the-art performance and generalization.

Generative image and video models have been explored for robot policies as well. Du et al. [24] explore mapping generated videos of robots back to control via inverse dynamics. Genima [25] makes this easier by drawing robot joint annotations as textured spheres on the video, while RT-Trajectory [26] generates trajectory annotations on images. In contrast, we predict annotations on 3D canonical views, which allows us to more easily solve for 3D end-effector poses even in free space, and improves generalization from few demonstrations by enabling SE(3) data augmentation. We show that our model can do free space reasoning for tasks such as *lift the bottle 30cm off the ground* (Figure 5). Methods that directly predict heatmaps on images have also been previously used for language-guided navigation tasks [27, 28].

3 Method: Orthographic-image Generation Vision-Language-Action model

At deployment time, the input to our system is a language instruction l , and a set of observations $O_k = \{I_k, D_k, P_k, K_k\}$, where I_k is an RGB image, D_k is a corresponding depth image, P_k is the camera pose, and K_k are the camera intrinsics, with a camera index k . The output of our system is an end-effector state $s = \langle p, \omega \rangle$, which consists of a position target p , rotation target ω . To complete a task, we sequentially execute our system, at each time-step using a motion planner to reach the predicted s , and obtain the next set of observations. Figure 2 shows our model architecture.

3.1 Multi-Modal Vision and Language Model

At the core of our system is a Large Language Model (LLM). The LLM takes as input a sequence of input tokens (vectors) $\langle t_1, \dots, t_i \rangle$, and generates a sequence of output tokens (vectors) $\langle t_{i+1}, \dots, t_{i+j} \rangle$. We use three types of input and output tokens: (1) text tokens, computed from a text tokenizer and embedding table, (2) input image tokens, which are either patch tokens or the image CLS token [29], computed by a visual encoder, and projected to LLM space through a learned MLP input projection, and (3) output image (action) tokens, which we add to the LLM vocabulary and decode as a special token using an additional MLP decoder. The output image tokens represent the next robot action. We use an image diffusion model to decode the image tokens into actions by producing images that contain annotations that illustrate the gripper position and rotation over a set of input views of the scene. The end-effector state s is decoded from these image annotations.

3.2 3D-Aware Reasoning with Orthogonal Orthographic Projections

To imbue the LLM with 3D-awareness, we unproject all input camera images into a point cloud in a canonical workspace. We then render the scene from a fixed set of views (independent of the input camera poses) before feeding them to the LLM. This brings the input and output in the same space, and our selection of the views (orthogonal views such as "front", "top", "left", "right", rendered in orthographic mode) ensures no ambiguity in output.

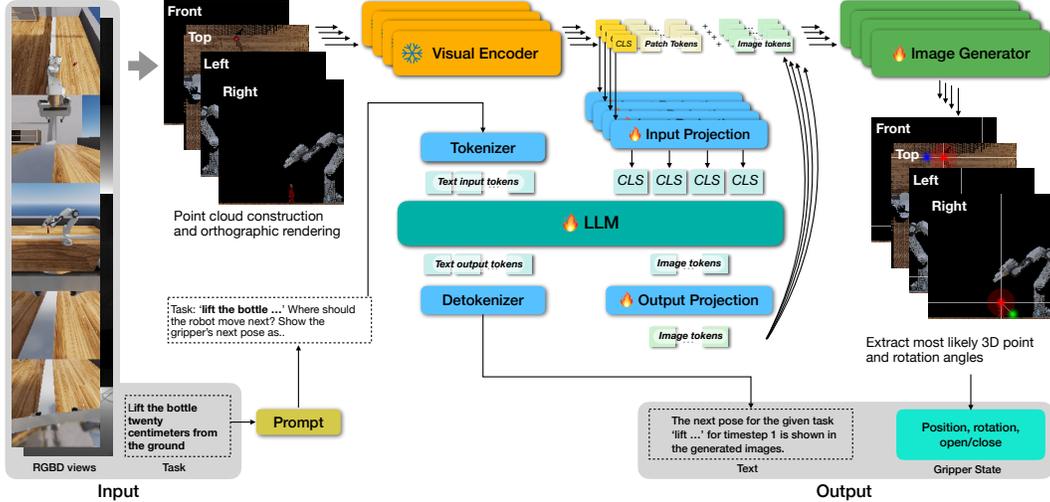


Figure 2: Model Overview. The input to our system is a task instruction and multiple RGB-D views of the scene. We build a point cloud from the input views and re-project it to orthographic projections from orthonormal views. The orthonormal views are fed into a Visual Encoder to derive a set of CLS and patch embeddings. CLS embeddings are projected into the LLM latent space and concatenated with a tokenized prompt that queries the next end-effector state and specifies the output format. The LLM outputs image token embeddings to condition the IMAGEGENERATOR, which are projected to the IMAGEGENERATOR’s input latent space, and then concatenated with skip-connected visual features. The IMAGEGENERATOR generates heatmaps-one per orthographic view-indicating the next end-effector pose. We decode the heatmaps by interpreting them as probabilities, and inferring the most likely 3D position across all views and one rotation angle per view.

Input Reprojection to Canonical Views. For each camera observation $\{I_k, D_k, P_k, K_k\}$ with N_k valid depth pixels, we compute a point cloud $C_k \in \mathbb{R}^{N_k \times 6}$, where each contains the RGB color and 3D coordinate in a fixed reference frame. We compute an aggregated point cloud $C = \bigcup_{k=1}^K C_k$ overall all input cameras. Next, we define a set of m canonical cameras $\{P_c^C, K_c^C\}_{c=1, \dots, m}$, where each P_c^C is a camera pose and K_c^C are the intrinsics for a given orthographic camera c . We then render the point cloud to an RGB image $I_c^C \in \mathbb{R}^{h \times w \times 3}$ seen by each orthographic camera.

Our system allows adapting the set of canonical camera parameters based on application and workspace geometry. In this work, we use four orthographic cameras that view the workspace from front, left, right and top directions, such that the workspace fills the camera image. Figure 2 shows the input RGBD views and the resulting orthographic images. While we use point clouds, our method may accommodate any 3D representation (e.g., neural radiance fields [30, 31] or novel-view synthesis methods [32]) that supports canonical view rendering.

LLM Input and Output. The canonical views $\{I_c^C\}_{c=1, \dots, m}$ constitute the visual input to our VLM system. We process each view with a VISUALENCODER and obtain an image embedding (CLS token) e_c^{CLS} , as well as a sequence of image patch embeddings $\langle e_c^1, \dots, e_c^n \rangle$. Next, we apply an input projection neural network to map each CLS embedding e_c^{CLS} , $c \in \{1, \dots, m\}$ into a token t_c^{CLS} compatible with the LLM input space. Finally, the input sequence to the LLM is the following sequence of tokens: $\langle \text{PROMPT}(l), t_{CLS}^1, \dots, t_{CLS}^m \rangle$, where $\text{PROMPT}(l)$ is a function that constructs a prompt for the instruction l and tokenizes it in a way compatible with the LLM.

We feed the input sequence to the LLM to produce an output sequence of format: $\langle t_a^1, t_a^2, t_a^3, t_a^4, t_c^1, t_c^2, \dots, t_c^j \rangle$, where $t_a^{(i)}$ are four image tokens that together represent the next action, and $t_c^{(j)}$ are accompanying text tokens as a response to the input prompt. We show the abbreviated prompt in Figure 2. Full prompt and text response is provided in the Appendix B.

Image Token Decoding for Action Prediction. Successful robot manipulation, such as picking objects or grabbing drawer handles, requires very precise end-effector position predictions. Although the LLM can output gripper poses in text with generally close predictions, we find them to lack the precision needed for practical applications. We thus propose decoding the image tokens to action

annotations over each of the canonical images. From these annotations, we can infer the next 3D gripper position and orientation by decoding and aggregating the image annotations in 3D.

First, we project each of the output image tokens t_a^i back to the visual embedding space with an output projection layer to obtain a set of embeddings e_a^i ; $i \in \{1, \dots, 4\}$. Next, we use these embeddings as well as the output of the VISUALENCODER (all patch embeddings $\langle e_c^1, \dots, e_c^n \rangle$, as well as the CLS token e_c^{CLS}) to condition an IMAGEGENERATOR network for each orthographic camera $c \in \{1, \dots, m\}$. The network outputs an RGB image with action predictions overlaid on the input canonical views:

$$H^c = \text{IMAGEGENERATOR}(e_a^i, e_c^x),$$

$$i \in \{1, \dots, 4\} \quad c \in \{1, \dots, m\} \quad x \in \{CLS, 1, \dots, 256\} \quad (1)$$

where $H_c \in \{\mathbb{R}^{h,w,3}\}$ is a reconstruction of the input canonical image I_c^C with overlaid annotations that encode the gripper position and orientation. Each image is aligned with one of the m canonical input views c , as shown in Figure 2. While any image generator can be used, we use StableDiffusion 1.5 [33]. e_a^i acts as the textual conditioning and e_c^x as the visual conditioning for the IMAGEGENERATOR. In practice, we use $e_a^i + \text{CLIP}(\text{PROMPT}(l))$ for textual conditioning, essentially adding a residual skip connection from a direct text encoding, a design choice informed by empirical experiments. CLIP is the textual encoder used for pretraining the StableDiffusion 1.5 model.

Extracting 3D Position and Rotation from Generated Images. We generate gripper position and Euler rotation as Gaussian distributions overlaid in different color channels on the input orthographic views. We infer a 3D position p^{hm} that best explains the predictions in each of the canonical views by solving the optimization problem:

$$p^{hm} = \arg \max_p \prod_{c=1, \dots, m} (H_c[\text{CAMERAPROJECTION}(p, P_c^C, K_c^C)] + \epsilon), \quad (2)$$

where CAMERAPROJECTION projects the 3D point p to 2D image coordinates w.r.t. the orthographic camera c . The square brackets represent a 2D pixel-wise indexing operation with interpolation to support sub-pixel coordinates. ϵ is a small value we add to allow decoding in situations where one of the heatmaps is zero for all 3D points.

We predict Euler rotations on the images such that the rotation along an axis is overlaid on the canonical view along that axis using 3 different colors. We present x-axis rotation on the front view, y-axis rotation on left and right views, and z-axis rotation on the top view as shown in Figure 2. Rotations are overlaid as Gaussian distributions at 30 pixel radius with reference to a horizontal line drawn from the translation point towards the right of the image. To decode rotation, we first extract the pixel location of the most likely rotation (r_x^c, r_y^c) and translation point (p_x^c, p_y^c) from the Gaussian distributions using a filtering operation for each view, and then compute the rotation angles using the arctangent function. The gripper’s open/close state is encoded with binary colored hotspots on the top-left of the image, as shown in Figure 1.

3.3 Training and Implementation Details

Our architecture builds on X-VILA [34], a multi-modal chat model supporting language, visual, video, and audio modalities. OG-VLA is initialized from X-VILA’s pretrained weights to leverage its any-to-any modality alignment, focusing on text-image input to text-image output alignment. We leave study of enriching human-robot interaction using the other modalities to future work. We train OG-VLA using DeepSpeed [35]. We freeze the visual encoder and tune the LLM, input and output projection networks, and IMAGEGENERATOR with end-to-end gradient flow. Following X-VILA, we use ImageBind [36] as the visual encoder, linear layers for input and output projections, and Stable Diffusion 1.5 [37] for the IMAGEGENERATOR. The LLM is based on Vicuna-7b v1.5 architecture. ImageBind encodes images into 256 patches (16×16) with a CLS token. Each training sample includes a natural language instruction l , visual observations I_k, D_k, P_k, K_k , and ground-truth gripper state \hat{s} . Following prior work [7, 8], we augment each keyframe with N SE(3)-transformed perturbations: translation in $[\pm 0.1\text{m}, \pm 0.1\text{m}, \pm 0.1\text{m}]$, rotation in $[\pm 0^\circ, \pm 0^\circ, \pm 90^\circ]$. All

Model	Pickup Object	Reorient Object	Open Drawer	Close Drawer	Open Cabinet	Close Cabinet	Pour Water	Transfer Water	Overall
6D-CLIPort [38]	6.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.8
-Novel Object	8.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0
-Novel Scene	10.4	0.0	0.0	0.0	0.0	1.3	0.0	0.0	1.5
-Novel State	0.0	0.0	0.0	0.8	0.0	0.0	0.0	0.0	0.1
PerAct [7]	83.3 ± 2.4	16.7 ± 6.2	30.0 ± 10.8	31.7 ± 8.5	25.0 ± 0.0	30.0 ± 0.0	36.7 ± 6.2	18.3 ± 2.4	34.0 ± 3.1
-Novel Object	75.0 ± 0.0	3.3 ± 2.4	0.0 ± 0.0	23.3 ± 13.1	0.0 ± 0.0	0.0 ± 0.0	30.0 ± 4.1	1.7 ± 2.4	16.7 ± 2.6
-Novel Scene	75.0 ± 4.1	13.3 ± 2.4	13.3 ± 9.4	30.0 ± 14.1	0.0 ± 0.0	6.7 ± 2.4	26.7 ± 6.2	3.3 ± 2.4	21.0 ± 3.1
-Novel State	16.7 ± 2.4	1.7 ± 2.4	5.0 ± 0.0	11.7 ± 6.2	0.0 ± 0.0	0.0 ± 0.0	5.0 ± 0.0	11.7 ± 2.4	6.5 ± 1.2
OG-VLA@30k	86.7 ± 2.9	15.0 ± 8.7	38.3 ± 2.9	51.7 ± 2.9	0.0 ± 0.0	16.7 ± 2.9	25.0 ± 5.0	16.7 ± 7.6	31.2 ± 2.9
-Novel Object	85.0 ± 5.0	0.0 ± 0.0	1.7 ± 2.9	55.0 ± 13.2	1.7 ± 2.9	5.0 ± 5.0	18.3 ± 2.9	6.7 ± 7.6	21.7 ± 0.7
-Novel Scene	70.0 ± 2.9	1.7 ± 2.8	26.7 ± 11.5	36.7 ± 5.8	1.7 ± 2.9	1.7 ± 2.9	16.7 ± 11.5	8.3 ± 2.9	20.8 ± 1.3
-Novel State	0.0 ± 0.0	13.3 ± 7.6	13.3 ± 2.9	20.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	8.3 ± 7.6	13.3 ± 2.9	8.5 ± 1.9
OG-VLA@100k	88.3 ± 2.4	16.7 ± 9.4	48.3 ± 2.4	56.7 ± 2.4	6.7 ± 4.7	23.3 ± 16.5	33.3 ± 6.2	28.3 ± 2.4	37.7 ± 0.6
-Novel Object	65.0 ± 8.2	15.0 ± 4.1	1.7 ± 2.4	58.3 ± 12.5	0.0 ± 0.0	5.0 ± 4.1	45.0 ± 8.2	8.3 ± 4.7	24.8 ± 1.2
-Novel Scene	75.0 ± 7.1	13.3 ± 8.5	31.7 ± 4.7	51.7 ± 2.4	1.7 ± 2.4	5.0 ± 4.1	26.7 ± 2.4	25.0 ± 7.1	28.8 ± 0.5
-Novel State	0.0 ± 0.0	13.3 ± 2.4	25.0 ± 7.1	15.0 ± 4.1	0.0 ± 0.0	0.0 ± 0.0	6.7 ± 4.7	20.0 ± 7.1	10.0 ± 0.9

Table 1: Success rate on ARNOLD [12]. Success rates for 8 tasks and 4 test splits are shown for 2 baseline models and our model at specified training iterations (30k and 100k). The first row for each model is the Novel Pose split (the Test set in ARNOLD). The numbers in bold represent the best system for each task and test split.

models are trained on 8×A100 GPUs with batch size 64. All models are trained for one run, and the LLM is LoRA-finetuned. Inference is done on a single A100 GPU. For inference-time image sampling from Stable Diffusion 1.5, we use 100 steps and a guidance scale of 7.0. These parameters result in reasonable sampled image quality and avoid jitters in the sampled image from the model’s latent space.

4 Simulation Experiments

4.1 Benchmarks and Datasets

We evaluate our method on the ARNOLD [12] and COLOSSEUM [11] benchmarks, which test language-grounded robot task learning in realistic 3D scenes with emphasis on testing generalization. ARNOLD uses five input cameras (front, base, left, wrist top, wrist bottom) and includes eight language-conditioned tasks (see Table 1), each with four generalization test splits: (1) Novel Pose (held-out object/robot placements), (2) Novel Object (unseen objects), (3) Novel Scene (unseen scenes with seen objects), and (4) Novel State (unseen goal states). ARNOLD tasks follow a two-keyframe format—grasping and manipulating (e.g., pull drawer to 50% open)—and do not require gripper state prediction. These tasks demand both object pose and free space reasoning. COLOSSEUM has four cameras (front, left shoulder, right shoulder, wrist) and features 20 language-conditioned tabletop tasks (e.g., close box, empty dishwasher) with 2–13 gripper keyframes (average 6), requiring gripper open/close prediction. It evaluates generalization via the *all perturbation* test set, which simultaneously alters object/table/background appearance, lighting, camera pose, and adds distractors.

We train on ARNOLD’s training split with 8 tasks, ~500 demos/task, and 2 keyframes/demo (7100 keyframes). We separately train on COLOSSEUM training split with 100 demos/task. For SE(3) augmentations, we set $N = 10$ for ARNOLD (~70k training samples) and $N = 5$ for COLOSSEUM (~1M training samples). All ARNOLD models are trained for 30k iterations and the final result is reported at 30k and 100k iterations (30k takes 1.5 days, and 100k takes 5 days to train). We train on COLOSSEUM for 250k iterations due to the larger training data size. We noticed that further training beyond 250k iterations degrades image generation quality. We conduct ablation studies at 30k on ARNOLD due to compute constraints.

Each ARNOLD test split includes 20 episodes; COLOSSEUM uses 25. We evaluate each ARNOLD model over 3 runs to account for simulator’s motion planner noise during keypoint execution, and we report means and standard deviations. Baseline models were trained by the respective benchmark au-

thors. We re-run PerAct [7] for 3 runs on ARNOLD for fair comparison, and report 6D-CLIPort [38] results from ARNOLD [12] due to lack of checkpoint release.

4.2 Simulation Results

ARNOLD We present our model results in Table 1. OG-VLA outperforms the baseline (PerAct) for most of the tasks’ Novel Pose split (seen objects and scenes), with a task-averaged relative improvement of 10.8%. We present our result at 30k and 100k iterations. OG-VLA improves significantly across generalization splits (Novel Object, Scene, and State), with relative increases of 20.0% and 46.5% respectively at 30k and 100k iterations in overall success rates. This demonstrates the effectiveness of our method in adapting pretrained visual and textual priors to robotics applications.

COLOSSEUM Figure 3 shows task-averaged success rate on *all perturbation* test set on the COLOSSEUM benchmark. We compare with baselines reported in COLOSSEUM, including R3M [19], MVP [20], 3DDA [39], RVT [8], and PerAct [7]. OG-VLA improves upon the baselines by a relative increase of 45.8%. The absolute performance remains low however (10.5%). We observe that COLOSSEUM tasks are much harder to learn due to longer sequence of keyframe predictions, leading to error accumulation for an imitation learning based method.

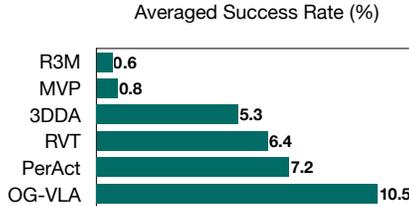


Figure 3: Evaluation on the COLOSSEUM benchmark [11]. Task-averaged success rate shows that OG-VLA outperforms all baselines on the hardest generalization test set (*all perturbation*).

4.3 Design Choice Ablations

We conducted extensive ablations on the ARNOLD benchmark to validate our design choices. We highlight the main findings in this section, with detailed results and analysis in Appendix C.

Action Prediction Approaches We experimented with two common alternatives to our image generation approach: (1) direct text-based action prediction and (2) adding additional action tokens to the LLM vocabulary that are decoded with MLP decoders directly to gripper states. Both alternatives failed to learn effective policies under our training conditions, likely due to insufficient visual reasoning capabilities for precise manipulation tasks with limited training data.

Image Generation Modes We compared three image generation approaches: (1) generation without scene reconstruction (black background), (2) generation with scene reconstruction (action annotations overlaid on the input image), and (3) generation with faded reconstruction (rescaled to 0-127 range, keeping colors in range 128-255 reserved for action annotations). Generation with faded reconstruction performed best overall on generalization splits (e.g. 23.8% vs. 12.8% without reconstruction and 20.8% with full reconstruction on the ARNOLD novel scene split). Generating full reconstructions complicates action decoding due to occasional color collisions. Generating black backgrounds was unstable during training, likely due to the challenge of adapting a model of natural images towards generating plain backgrounds.

Architecture Components Our ablations revealed several key insights: (1) Unlike prior work, tiling orthographic views decreased performance (24.8% vs. 31.2%); (2) Removing the LLM significantly reduced performance (20.0% vs. 31.2%), likely due to its strong priors and its role in conditioning consistent generation across views; (3) Directly bypassing the instruction to the Image Generator decreased performance by 9.5%, suggesting the LLM’s image token outputs preserve crucial task information. These findings validate the importance of each component in OG-VLA.

5 Real Robot Experiments

Experimental Setup We perform real-world experiments on a Franka Emika Panda arm mounted on a tabletop, with a single front-facing camera, as shown in Figure 1. We collect 3–5 demos for 4 real-world tasks—22 demos in total—with human-annotated keyframes and a motion planner [40] to achieve annotated keyframes. We train both the baseline and our model on this dataset. For our model, we augment each keyframe with 10 SE(3) perturbations and finetune the Arnold-pretrained OG-VLA@30k checkpoint for another 10k iterations with a batch size of 64. For π_0 -FAST [4] baseline, we use the provided pretrained checkpoint trained on 10k+ hours of robot data across various robot setups with actions in both joint and end-effector control spaces. We finetune it on our dataset with joint angle actions for 30k iterations with a batch size of 32, as used in most of their pre-training and finetuning experiments [4, 22]. During inference, the model predicts an action chunk of 10 actions; we execute the last action in the sequence for faster execution.

Quantitative Results We report our results in Table 2. Each test set success rate is averaged over 10 episodes. For novel object variation, we use unseen colored and shaped objects. For novel scene variation, we introduce distractors and change lighting, background, and table appearance. Please refer to Appendix D for training and test scene details. Our results show that OG-VLA can adapt to new tasks with only 3–5 demonstrations and generalize well to novel poses, objects, and scenes. Although we attempted to compare with π_0 -FAST, it failed at all tasks, learning only to reach the block with 30% success. We hypothesize that this is due to the small number of demonstrations and lack of support for SE(3) data augmentation without our 3D representation.

Model	Pickup Object	Put Object in Drawer	Open Drawer	Close Drawer
OG-VLA@10k	100.0	90.0	60.0	90.0
- Novel Object	80.0	70.0	30.0	50.0
- Novel Scene	90.0	80.0	50.0	90.0

Table 2: Real world success rate (%). Success rates are reported for 4 real world tasks and novel pose, novel object, and novel scene test splits, averaged over 10 episodes each.

Qualitative Results We show qualitative examples of real world evaluations in Figure 4 for the task ‘put object in the drawer’. For training demonstrations, we use a blue cube. During evaluation, we replace it with bottle or perturb the scene by placing a newspaper under the cube also make the scene brighter. OG-VLA can generalize to manipulating different objects as well as to unseen scenes for a given task.

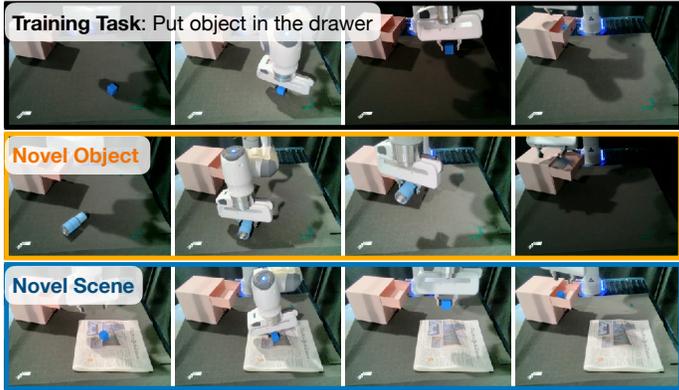


Figure 4: Qualitative example. Showing generalization of OG-VLA to unseen scenarios.

6 Conclusion

We introduced OG-VLA, a novel architecture and learning framework that combines the generalization strengths of Vision-Language-Action (VLA) models with the robustness of 3D-aware keyframe policies for robotic manipulation. By leveraging foundation models in language and vision, OG-VLA improves generalization to unseen instructions, objects, and scenes while maintaining precise control. Our approach ensures input-view invariance by projecting multi-view RGBD observations into a canonical point cloud representation, which is processed through a vision backbone, an LLM, and an image diffusion model to generate actions as images. OG-VLA achieves

state-of-the-art performance on robotic manipulation tasks, particularly on scene and object generalization tests. Moreover, OG-VLA can adapt to real-world tasks in 3-5 demonstrations with strong generalization to unseen objects and scenes. These results highlight the effectiveness of integrating pretrained visual and text priors into a structured 3D-aware framework for robot learning. Future work will explore extending OG-VLA to complex long-horizon tasks, external data augmentation and co-training techniques.

7 Limitations

OG-VLA has several strengths, but is not free of limitation. One challenge arises from the reliance on orthographic canonical views, which, while effective in many scenarios, can struggle in situations with severe occlusions, such as when multiple objects are stacked on a shelf against a wall. Occlusions may lead to partial or incorrect scene representations, potentially affecting the downstream task performance. Moreover, keyframe-based policies cover a broad range of tasks, however, are not able to solve several dynamic tasks, such as object tossing, or pressing with a desired force. Another consideration is the computational cost associated with OG-VLA. The current training procedure is computationally intensive, requiring substantial time and resources. We aim to address this by optimizing the model architecture, and exploring strategies such as distillation or parameter-efficient fine-tuning to accelerate learning while maintaining performance. Markovian policy learning also makes long horizon tasks quite challenging where several keyframes need to be predicted.

Real World Limitations Single camera input makes some of the orthographic views redundant, uninformative or even noisy, which affects downstream performance. In this case, even if the predictions are correct in a few views, the noisy predictions affect the eventually decoded keyframe.

References

- [1] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson, Q. Vuong, T. Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- [2] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, X. Chen, K. Choromanski, T. Ding, D. Driess, A. Dubey, C. Finn, et al. RT-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.
- [3] M. Kim, K. Pertsch, S. Karamcheti, T. Xiao, A. Balakrishna, S. Nair, R. Rafailov, E. Foster, G. Lam, P. Sanketi, Q. Vuong, T. Kollar, B. Burchfiel, R. Tedrake, D. Sadigh, S. Levine, P. Liang, and C. Finn. OpenVLA: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024.
- [4] K. Pertsch, K. Stachowicz, B. Ichter, D. Driess, S. Nair, Q. Vuong, O. Mees, C. Finn, and S. Levine. FAST: Efficient action tokenization for vision-language-action models, 2025.
- [5] O. X.-E. Collaboration, A. O’Neill, A. Rehman, A. Maddukuri, A. Gupta, A. Padalkar, A. Lee, A. Pooley, A. Gupta, A. Mandlekar, A. Jain, A. Tung, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Gupta, A. Wang, A. Kolobov, A. Singh, A. Garg, A. Kembhavi, A. Xie, A. Brohan, A. Raffin, A. Sharma, A. Yavary, A. Jain, A. Balakrishna, A. Wahid, B. Burgess-Limerick, B. Kim, B. Schölkopf, B. Wulfe, B. Ichter, C. Lu, C. Xu, C. Le, C. Finn, C. Wang, C. Xu, C. Chi, C. Huang, C. Chan, C. Agia, C. Pan, C. Fu, C. Devin, D. Xu, D. Morton, D. Driess, D. Chen, D. Pathak, D. Shah, D. Büchler, D. Jayaraman, D. Kalashnikov, D. Sadigh, E. Johns, E. Foster, F. Liu, F. Ceola, F. Xia, F. Zhao, F. V. Frujeri, F. Stulp, G. Zhou, G. S. Sukhatme, G. Salhotra, G. Yan, G. Feng, G. Schiavi, G. Berseth, G. Kahn, G. Yang, G. Wang, H. Su, H.-S. Fang, H. Shi, H. Bao, H. B. Amor, H. I. Christensen, H. Furuta, H. Walke, H. Fang, H. Ha, I. Mordatch, I. Radosavovic, I. Leal, J. Liang, J. Abou-Chakra, J. Kim, J. Drake, J. Peters, J. Schneider, J. Hsu, J. Bohg, J. Bingham, J. Wu, J. Gao, J. Hu,

- J. Wu, J. Wu, J. Sun, J. Luo, J. Gu, J. Tan, J. Oh, J. Wu, J. Lu, J. Yang, J. Malik, J. Silvério, J. Hejna, J. Booher, J. Tompson, J. Yang, J. Salvador, J. J. Lim, J. Han, K. Wang, K. Rao, K. Pertsch, K. Hausman, K. Go, K. Gopalakrishnan, K. Goldberg, K. Byrne, K. Oslund, K. Kawaharazuka, K. Black, K. Lin, K. Zhang, K. Ehsani, K. Lekkala, K. Ellis, K. Rana, K. Srinivasan, K. Fang, K. P. Singh, K.-H. Zeng, K. Hatch, K. Hsu, L. Itti, L. Y. Chen, L. Pinto, L. Fei-Fei, L. Tan, L. J. Fan, L. Ott, L. Lee, L. Weihs, M. Chen, M. Lepert, M. Memmel, M. Tomizuka, M. Itkina, M. G. Castro, M. Spero, M. Du, M. Ahn, M. C. Yip, M. Zhang, M. Ding, M. Heo, M. K. Srirama, M. Sharma, M. J. Kim, N. Kanazawa, N. Hansen, N. Heess, N. J. Joshi, N. Suenderhauf, N. Liu, N. D. Palo, N. M. M. Shafiullah, O. Mees, O. Kroemer, O. Bastani, P. R. Sanketi, P. T. Miller, P. Yin, P. Wohlhart, P. Xu, P. D. Fagan, P. Mitranon, P. Sermanet, P. Abbeel, P. Sundaresan, Q. Chen, Q. Vuong, R. Rafailov, R. Tian, R. Doshi, R. Mart'in-Mart'in, R. Bajjal, R. Scalise, R. Hendrix, R. Lin, R. Qian, R. Zhang, R. Mendonca, R. Shah, R. Hoque, R. Julian, S. Bustamante, S. Kirmani, S. Levine, S. Lin, S. Moore, S. Bahl, S. Dass, S. Sonawani, S. Song, S. Xu, S. Haldar, S. Karamcheti, S. Adebola, S. Guist, S. Nasiriany, S. Schaal, S. Welker, S. Tian, S. Ramamoorthy, S. Dasari, S. Belkhale, S. Park, S. Nair, S. Mirchandani, T. Osa, T. Gupta, T. Harada, T. Matsushima, T. Xiao, T. Kollar, T. Yu, T. Ding, T. Davchev, T. Z. Zhao, T. Armstrong, T. Darrell, T. Chung, V. Jain, V. Vanhoucke, W. Zhan, W. Zhou, W. Burgard, X. Chen, X. Chen, X. Wang, X. Zhu, X. Geng, X. Liu, X. Liangwei, X. Li, Y. Pang, Y. Lu, Y. J. Ma, Y. Kim, Y. Chebotar, Y. Zhou, Y. Zhu, Y. Wu, Y. Xu, Y. Wang, Y. Bisk, Y. Dou, Y. Cho, Y. Lee, Y. Cui, Y. Cao, Y.-H. Wu, Y. Tang, Y. Zhu, Y. Zhang, Y. Jiang, Y. Li, Y. Li, Y. Iwasawa, Y. Matsuo, Z. Ma, Z. Xu, Z. J. Cui, Z. Zhang, Z. Fu, and Z. Lin. Open X-Embodiment: Robotic learning datasets and RT-X models, 2023.
- [6] Y. Li, Y. Deng, J. Zhang, J. Jang, M. Memmel, C. R. Garrett, F. Ramos, D. Fox, A. Li, A. Gupta, and A. Goyal. HAMSTER: Hierarchical action models for open-world robot manipulation. In *ICLR*, 2025.
- [7] M. Shridhar, L. Manuelli, and D. Fox. Perceiver-Actor: A multi-task transformer for robotic manipulation. In *Proceedings of the 6th Conference on Robot Learning (CoRL)*, 2022.
- [8] A. Goyal, J. Xu, Y. Guo, V. Blukis, Y.-W. Chao, and D. Fox. RVT: Robotic view transformer for 3d object manipulation. *Proceedings of the 7th Conference on Robot Learning (CoRL)*, 2023.
- [9] A. Goyal, V. Blukis, J. Xu, Y. Guo, Y.-W. Chao, and D. Fox. RVT-2: Learning precise manipulation from few demonstrations. *arXiv preprint arXiv:2406.08545*, 2024.
- [10] T. Gervet, Z. Xian, N. Gkanatsios, and K. Fragkiadaki. Act3D: 3d feature field transformers for multi-task robotic manipulation. In *7th Annual Conference on Robot Learning*, 2023.
- [11] W. Pumacay, I. Singh, J. Duan, R. Krishna, J. Thomason, and D. Fox. The COLOSSEUM: A benchmark for evaluating generalization for robotic manipulation. In *RSS 2024 Workshop: Data Generation for Robotics*, 2024.
- [12] R. Gong, J. Huang, Y. Zhao, H. Geng, X. Gao, Q. Wu, W. Ai, Z. Zhou, D. Terzopoulos, S.-C. Zhu, et al. ARNOLD: A benchmark for language-grounded task learning with continuous states in realistic 3d scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023.
- [13] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman, A. Herzog, D. Ho, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, E. Jang, R. J. Ruano, K. Jeffrey, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, K.-H. Lee, S. Levine, Y. Lu, L. Luu, C. Parada, P. Pastor, J. Quiambao, K. Rao, J. Rettinghouse, D. Reyes, P. Sermanet, N. Sievers, C. Tan, A. Toshev, V. Vanhoucke, F. Xia, T. Xiao, P. Xu, S. Xu, M. Yan, and A. Zeng. Do as i can and not as i say: Grounding language in robotic affordances. In *arXiv preprint arXiv:2204.01691*, 2022.

- [14] I. Singh, V. Blukis, A. Mousavian, A. Goyal, D. Xu, J. Tremblay, D. Fox, J. Thomason, and A. Garg. ProgPrompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11523–11530, 2023.
- [15] J. Liang, W. Huang, F. Xia, P. Xu, K. Hausman, B. Ichter, P. Florence, and A. Zeng. Code as Policies: Language model programs for embodied control. In *arXiv preprint arXiv:2209.07753*, 2022.
- [16] Y. J. Ma, W. Liang, G. Wang, D.-A. Huang, O. Bastani, D. Jayaraman, Y. Zhu, L. Fan, and A. Anandkumar. Eureka: Human-level reward design via coding large language models. In *ICLR*, 2024.
- [17] W. Yu, N. Gileadi, C. Fu, S. Kirmani, K.-H. Lee, M. Gonzalez Arenas, H.-T. Lewis Chiang, T. Erez, L. Hasenclever, J. Humplik, B. Ichter, T. Xiao, P. Xu, A. Zeng, T. Zhang, N. Heess, D. Sadigh, J. Tan, Y. Tassa, and F. Xia. Language to rewards for robotic skill synthesis. *Arxiv preprint arXiv:2306.08647*, 2023.
- [18] W. Huang, C. Wang, R. Zhang, Y. Li, J. Wu, and L. Fei-Fei. VoxPoser: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023.
- [19] S. Nair, A. Rajeswaran, V. Kumar, C. Finn, and A. Gupta. R3M: A universal visual representation for robot manipulation. In *6th Annual Conference on Robot Learning*, 2022.
- [20] I. Radosavovic, T. Xiao, S. James, P. Abbeel, J. Malik, and T. Darrell. Real-world robot learning with masked visual pre-training. In *6th Annual Conference on Robot Learning*, 2022.
- [21] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. Joshi, R. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath, I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao, M. Ryoo, G. Salazar, P. Sanketi, K. Sayed, J. Singh, S. Sontakke, A. Stone, C. Tan, H. Tran, V. Vanhoucke, S. Vega, Q. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. RT-1: Robotics transformer for real-world control at scale. In *arXiv preprint arXiv:2212.06817*, 2022.
- [22] K. Black, N. Brown, D. Driess, A. Esmail, M. Equi, C. Finn, N. Fusai, L. Groom, K. Hausman, B. Ichter, S. Jakubczak, T. Jones, L. Ke, S. Levine, A. Li-Bell, M. Mothukuri, S. Nair, K. Pertsch, L. X. Shi, J. Tanner, Q. Vuong, A. Walling, H. Wang, and U. Zhilinsky. π_0 : A vision-language-action flow model for general robot control, 2024.
- [23] A. Padalkar, A. Pooley, A. Jain, A. Bewley, A. Herzog, A. Irpan, A. Khazatsky, A. Rai, A. Singh, A. Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023.
- [24] Y. Du, S. Yang, B. Dai, H. Dai, O. Nachum, J. Tenenbaum, D. Schuurmans, and P. Abbeel. Learning universal policies via text-guided video generation. In *Advances in Neural Information Processing Systems*, volume 36, pages 9156–9172, 2023.
- [25] M. Shridhar, Y. L. Lo, and S. James. Generative image as action models. In *Proceedings of the 8th Conference on Robot Learning (CoRL)*, 2024.
- [26] J. Gu, S. Kirmani, P. Wohlhart, Y. Lu, M. G. Arenas, K. Rao, W. Yu, C. Fu, K. Gopalakrishnan, Z. Xu, et al. Rt-trajectory: Robotic task generalization via hindsight trajectory sketches. *arXiv preprint arXiv:2311.01977*, 2023.
- [27] P. Anderson, A. Shrivastava, D. Parikh, D. Batra, and S. Lee. Chasing ghosts: Instruction following as bayesian state tracking. *Advances in neural information processing systems*, 32, 2019.

- [28] V. Blukis, Y. Terme, E. Niklasson, R. A. Knepper, and Y. Artzi. Learning to map natural language instructions to physical quadcopter control using simulated flight. In *Proceedings of the Conference on Robot Learning (CoRL)*, 2019.
- [29] J. Devlin. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [30] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [31] A. Yu, V. Ye, M. Tancik, and A. Kanazawa. Pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4578–4587, 2021.
- [32] J. Kulhánek, E. Derner, T. Sattler, and R. Babuška. Viewformer: Nerf-free neural rendering from few images using transformers. In *European Conference on Computer Vision*, pages 198–216. Springer, 2022.
- [33] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models, 2021.
- [34] H. Ye, D.-A. Huang, Y. Lu, Z. Yu, W. Ping, A. Tao, J. Kautz, S. Han, D. Xu, P. Molchanov, et al. X-VILA: Cross-modality alignment for large language model. *arXiv preprint arXiv:2405.19335*, 2024.
- [35] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He. ZeRO: Memory optimizations toward training trillion parameter models, 2020. URL <https://github.com/microsoft/DeepSpeed>.
- [36] R. Girdhar, A. El-Nouby, Z. Liu, M. Singh, K. V. Alwala, A. Joulin, and I. Misra. ImageBind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023.
- [37] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, June 2022.
- [38] K. Zheng, X. Chen, O. Jenkins, and X. E. Wang. VLMbench: A compositional benchmark for vision-and-language manipulation. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- [39] T.-W. Ke, N. Gkanatsios, and K. Fragkiadaki. 3d diffuser actor: Policy diffusion with 3d scene representations. *Arxiv*, 2024.
- [40] Y. Zhu, A. Joshi, P. Stone, and Y. Zhu. VIOLA: Imitation learning for vision-based manipulation with object proposal priors. *arXiv preprint arXiv:2210.11339*, 2022.
- [41] X. Li, C. Mata, J. Park, K. Kahatapitiya, Y. S. Jang, J. Shang, K. Ranasinghe, R. Burgert, M. Cai, Y. J. Lee, and M. S. Ryoo. LLaRA: Supercharging robot learning data for vision-language policy. In *International Conference on Learning Representations*, 2025.
- [42] W. Yuan, J. Duan, V. Blukis, W. Pumacay, R. Krishna, A. Murali, A. Mousavian, and D. Fox. RoboPoint: A vision-language model for spatial affordance prediction for robotics. *arXiv preprint arXiv:2406.10721*, 2024.

Task Types	Goal States	Success Ranges	Training Data
Pickup Object	10, 20, 30, 40 (cm)	± 5 cm	623
ReorientObject	0, 45, 135, 180 ($^\circ$)	$\pm 20^\circ$	355
Open Drawer	25, 50, 75, 100 (%)	$\pm 10\%$	554
Close Drawer	0, 25, 50, 75 (%)	$\pm 10\%$	671
Open Cabinet	25, 50, 75, 100 (%)	$\pm 10\%$	319
Close Cabinet	0, 25, 50, 75 (%)	$\pm 10\%$	478
Pour Water	25, 50, 75, 100 (%)	$\pm 10\%$	312
Transfer Water	20, 40, 60, 80 (%)	$\pm 10\%$	259

Table 3: Overview of the 8 tasks in ARNOLD. Each task features 4 goal states specified by human language, one of which is reserved for novel state evaluation and the other three are seen in the training dataset. The task is considered successful when the object state remains in the success range for two seconds. For Transfer Water, an additional condition of only less than 10% spillage of the original amount of water in the cup is imposed.

Tasks	Examples of Templates
Pickup Object	<i>Raise</i> [value_object] [value_height] <i>above the ground</i>
Reorient Object	<i>Reorient</i> [value_object] [value_degree] <i>away from the up axis</i>
Open Drawer	<i>Open the</i> [value_position] [value_object] [value_percent]
Close Drawer	<i>Close the</i> [value_position] [value_object] [value_percent]
Open Cabinet	<i>Open the</i> [value_position] [value_object] [value_percent]
Close Cabinet	<i>Close the</i> [value_position] [value_object] [value_percent]
Pour Water	<i>Pour</i> [value_percent] <i>water out of</i> [value_object]
Transfer Water	<i>Transfer</i> [value_percent] <i>water to</i> [value_object]

Table 4: Examples of instruction templates used for the tasks.

Appendix

A Tasks

B Prompts used in model variants

B.1 OG-VLA Prompt and Response

Prompt: Task: “the bottle should be twenty centimeters from the ground.”. Where should the robot move next? Show the gripper’s next pose as translation and rotation heatmaps on the input orthographic views. Translation should be represented as red heatmap on all 4 views. Follow the provided instruction to compute correct translation points in the images. Rotation should be represented as yellow, blue, and green heatmaps for the front, top, and left views, corresponding to the x, z, and y axes respectively.

Response: The next gripper pose for the given task ‘the bottle should be twenty centimeters from the ground.’ for timestep 1 is shown in the generated images.

B.2 Text Action model Prompt and Response

Prompt: Task: “pull the top drawer 50% open”. Where should the robot move next? Format the robot’s gripper action as a relative 3D coordinate, an Euler rotation, and a binary gripper open/close state. All numbers are floats with two decimal places, each in relative coordinates.

Response: pos: [0.54, 0.42, 0.62], rot: [-1.57, -0.0, -1.57]

C Detailed Ablation Results and Analysis

C.1 Action Prediction Ablations

We experimented with two common alternative architectural choices for action generation in VLAs, both of which failed to learn a working policy under similar training and evaluation settings as for OG-VLA.

1) Text Action: is an architecture variant where the LLM also produces gripper’s next pose in the form of text akin to prior VLA models [41]. In general, the raw output sequence from the LLM can be of slightly different format and contain additional text (e.g. text like *the next robot action is:*), so long as it contains the information above. We apply regex parsing to extract from the text the gripper position p^{text} , orientation ω^{text} . The prompt used for this ablation is shown in Appendix B.2.

Predicting actions as text without any visual reasoning on the output end of the model results in an imprecise action prediction model. Li et al. [41] have shown this architecture to work when training with large-scale robot datasets. However, we find that with a small robotics dataset, it is hard for VLAs to learn precise control via direct text token prediction.

2) Action Tokens: is an architecture variant where the LLM produces special action tokens added to the LLM vocabulary, such as $[trans0]$ or $[rot0]$ for translation and rotation modalities, akin to that for the image modality. This architecture is closer to works that perform action tokenization [3]. We decode the hidden state vectors for these tokens using additional MLP decoders to predict 3D translation and rotation vectors.

The model still struggles to predict sufficiently precise actions to perform the tasks. This may be due to insufficient visual reasoning available for action decoding, as the decoders use the hidden states corresponding to the special token produced by the LLM. We also observed degeneration of LLM’s ability to produce coherent and grammatical text as the training progressed for this model design. This may have been caused by the training of new tokens and additional decoders failing to preserve original reasoning capabilities of the model, an issue that might be addressed by co-training with the pretraining datasets. We do not observe degeneration when finetuning the model with original architectural components in OG-VLA without any co-training with other datasets. We leave the study with co-training for above architectural variants to future works.

C.2 Image Generation Mode

Model	Pickup Object	Reorient Object	Open Drawer	Close Drawer	Open Cabinet	Close Cabinet	Pour Water	Transfer Water	Overall
(1) No Reconstruction	76.7 ± 6.2	10.0 ± 10.8	3.3 ± 2.4	20.0 ± 0.0	1.7 ± 2.4	10.0 ± 7.1	28.3 ± 11.8	13.3 ± 12.5	20.4 ± 4.7
-Novel Object	56.7 ± 6.2	10.0 ± 0.0	3.3 ± 2.4	23.3 ± 15.5	0.0 ± 0.0	10.0 ± 8.2	13.3 ± 8.5	6.7 ± 2.4	15.4 ± 5.1
-Novel Scene	51.7 ± 8.5	6.7 ± 6.2	10.0 ± 4.1	10.0 ± 8.2	0.0 ± 0.0	5.0 ± 4.1	11.7 ± 6.2	6.7 ± 6.2	12.7 ± 3.0
-Novel State	0.0 ± 0.0	6.7 ± 6.2	6.7 ± 2.4	5.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	13.3 ± 8.5	0.0 ± 0.0	4.0 ± 2.1
(2) Reconstruction	86.7 ± 2.4	15.0 ± 7.1	38.3 ± 2.4	51.7 ± 2.4	0.0 ± 0.0	16.7 ± 2.4	25.0 ± 4.1	16.7 ± 6.2	31.2 ± 2.3
-Novel Object	85.0 ± 4.1	0.0 ± 0.0	1.7 ± 2.4	55.0 ± 10.8	0.0 ± 0.0	5.0 ± 4.1	18.3 ± 2.4	6.7 ± 6.2	21.5 ± 0.3
-Novel Scene	73.3 ± 2.4	1.7 ± 2.4	26.7 ± 9.4	36.7 ± 4.7	1.7 ± 2.4	1.7 ± 2.4	16.7 ± 9.4	8.3 ± 2.4	20.8 ± 1.1
-Novel State	0.0 ± 0.0	13.3 ± 6.2	13.3 ± 2.4	20.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	8.3 ± 6.2	13.3 ± 2.4	8.5 ± 1.6
(3) Faded Reconstruction	93.3 ± 2.4	5.0 ± 4.1	23.3 ± 2.4	36.7 ± 8.5	0.0 ± 0.0	5.0 ± 4.1	30.0 ± 0.0	20.0 ± 10.8	26.7 ± 2.3
-Novel Object	75.0 ± 4.1	18.3 ± 2.4	0.0 ± 0.0	55.0 ± 4.1	0.0 ± 0.0	8.3 ± 4.7	36.7 ± 11.8	5.0 ± 4.1	24.8 ± 2.3
-Novel Scene	76.7 ± 6.2	11.7 ± 6.2	28.3 ± 4.7	28.3 ± 2.4	5.0 ± 4.1	8.3 ± 2.4	20.0 ± 8.2	11.7 ± 6.2	23.8 ± 3.1
-Novel State	0.0 ± 0.0	18.3 ± 10.3	6.7 ± 2.4	10.0 ± 4.1	1.7 ± 2.4	0.0 ± 0.0	1.7 ± 2.4	11.7 ± 6.2	6.2 ± 0.5

Table 5: Success rates for image generation modes for each task. We show that generating actions with reconstruction or with faded reconstructions work better than that without reconstruction.

We study three image generation modes for action prediction: (1) Generation without reconstruction: an all black image background (2) Generation with reconstruction: an RGB image that is a reconstruction of the input image, and (3) Generation with faded reconstruction: a shifted RGB image between the range $[0, 128]$ that is a reconstruction of the input image. For each image mode, the action Gaussian distributions are overlaid on these backgrounds. These three choices carry trade-offs. The first method is the purest way to represent actions, but it appears challenging for image generators pre-trained on generating color images to learn to generate uniform black backgrounds. The second method does not require un-learning generation of color images, but burdens the generator with the additional reconstruction task, which has the potential to take model capacity away from action generation. On the flipside, it has the potential for some positive cross-task transfer, and better scene understanding and visual reasoning. The third method is a middle ground between the other two methods, which we include in our study. It eases action decoding by keeping values in the $(128, 255]$ range exclusively for action annotations. For methods (2) and (3), we apply an additional filtering step to identify the Gaussian and recover a grayscale heatmap.

We report success rates across tasks and test splits in Table 5. Generation without reconstruction of the scene performs the worst of all generation modes. This may be due to forcing the IMAGEGENERATOR to unlearn its prior of generating natural images and forcing it to only predict Gaussian distributions. We also observe instability in training this version, as sometimes the IMAGEGENERATOR

ATOR would collapse to produce completely black or noisy images and stop generating actions on them. Generation with reconstruction and that with faded reconstruction did not show consistent difference over all test splits. Generation with reconstruction learns a policy that’s better by 4.4% on Novel Pose split and 2.3% Novel State split. Generation with faded reconstruction performs better on Novel Object and Novel State splits by 3.3% and 3%. Therefore, we conclude that these model variants have similar performance, and generation with faded reconstruction did not work better as we had hypothesized due to its balance between forcing the model to focus less on scene reconstruction and more on action prediction. Therefore, we report generation with reconstruction as our final method.

C.3 Model Ablations

Model	Pickup Object	Reorient Object	Open Drawer	Close Drawer	Open Cabinet	Close Cabinet	Pour Water	Transfer Water	Overall
OG-VLA	86.7 ± 2.4	15.0 ± 7.1	38.3 ± 2.4	51.7 ± 2.4	0.0 ± 0.0	16.7 ± 2.4	25.0 ± 4.1	16.7 ± 6.2	31.2 ± 2.3
-Novel Object	85.0 ± 4.1	0.0 ± 0.0	1.7 ± 2.4	55.0 ± 10.8	0.0 ± 0.0	5.0 ± 4.1	18.3 ± 2.4	6.7 ± 6.2	21.5 ± 0.3
-Novel Scene	73.3 ± 2.4	1.7 ± 2.4	26.7 ± 9.4	36.7 ± 4.7	1.7 ± 2.4	1.7 ± 2.4	16.7 ± 9.4	8.3 ± 2.4	20.8 ± 1.1
-Novel State	0.0 ± 0.0	13.3 ± 6.2	13.3 ± 2.4	20.0 ± 0.0	0.0 ± 0.0	0.0 ± 0.0	8.3 ± 6.2	13.3 ± 2.4	8.5 ± 1.6
+Tiled Views	75.0 ± 7.1	6.7 ± 4.7	33.3 ± 2.4	28.3 ± 6.2	1.7 ± 2.4	20.0 ± 4.1	15.0 ± 10.8	18.3 ± 6.2	24.8 ± 0.3
-Novel Object	58.3 ± 8.5	15.0 ± 7.1	1.7 ± 2.4	30.0 ± 0.0	0.0 ± 0.0	10.0 ± 4.1	40.0 ± 4.1	6.7 ± 6.2	20.2 ± 1.3
-Novel Scene	60.0 ± 7.1	15.0 ± 10.8	45.0 ± 7.1	21.7 ± 6.2	5.0 ± 4.1	5.0 ± 4.1	26.7 ± 6.2	1.7 ± 2.4	22.5 ± 0.5
-Novel State	0.0 ± 0.0	10.0 ± 4.1	10.0 ± 7.1	18.3 ± 4.7	1.7 ± 2.4	1.7 ± 2.4	1.7 ± 2.4	16.7 ± 11.8	7.5 ± 2.6
-LLM	86.7 ± 2.4	5.0 ± 7.1	6.7 ± 4.7	33.3 ± 4.7	0.0 ± 0.0	6.7 ± 4.7	15.0 ± 0.0	6.7 ± 2.4	20.0 ± 1.5
-Novel Object	68.3 ± 6.2	1.7 ± 2.4	6.7 ± 9.4	40.0 ± 4.1	0.0 ± 0.0	3.3 ± 2.4	10.0 ± 4.1	8.3 ± 2.4	17.3 ± 1.6
-Novel Scene	71.7 ± 6.2	8.3 ± 6.2	18.3 ± 2.4	21.7 ± 8.5	3.3 ± 2.4	0.0 ± 0.0	15.0 ± 4.1	5.0 ± 4.1	17.9 ± 3.3
-Novel State	0.0 ± 0.0	0.0 ± 0.0	6.7 ± 2.4	16.7 ± 2.4	0.0 ± 0.0	1.7 ± 2.4	3.3 ± 2.4	10.0 ± 4.1	4.8 ± 0.8
+Tiled Views -LLM	71.7 ± 10.3	1.7 ± 2.4	13.3 ± 8.5	16.7 ± 4.7	0.0 ± 0.0	8.3 ± 2.4	15.0 ± 8.2	10.0 ± 0.0	17.1 ± 3.4
-Novel Object	56.7 ± 8.5	8.3 ± 2.4	1.7 ± 2.4	16.7 ± 8.5	0.0 ± 0.0	1.7 ± 2.4	15.0 ± 4.1	6.7 ± 2.4	13.3 ± 1.6
-Novel Scene	61.7 ± 6.2	5.0 ± 4.1	20.0 ± 4.1	11.7 ± 6.2	0.0 ± 0.0	3.3 ± 4.7	10.0 ± 0.0	10.0 ± 0.0	15.2 ± 1.2
-Novel State	0.0 ± 0.0	6.7 ± 2.4	10.0 ± 0.0	30.0 ± 4.1	1.7 ± 2.4	0.0 ± 0.0	6.7 ± 6.2	3.3 ± 4.7	7.3 ± 0.8
-Instruction to IG	71.7 ± 4.7	8.3 ± 2.4	20.0 ± 10.8	40.0 ± 12.2	1.7 ± 2.4	11.7 ± 9.4	5.0 ± 4.1	15.0 ± 4.1	21.7 ± 2.6
-Novel Object	66.7 ± 6.2	0.0 ± 0.0	1.7 ± 2.4	45.0 ± 4.1	0.0 ± 0.0	8.3 ± 2.4	20.0 ± 4.1	1.7 ± 2.4	17.9 ± 0.8
-Novel Scene	50.0 ± 8.2	11.7 ± 2.4	25.0 ± 0.0	26.7 ± 2.4	10.0 ± 0.0	11.7 ± 2.4	13.3 ± 4.7	5.0 ± 4.1	19.2 ± 2.1
-Novel State	0.0 ± 0.0	3.3 ± 4.7	8.3 ± 6.2	13.3 ± 8.5	0.0 ± 0.0	0.0 ± 0.0	1.7 ± 2.4	8.3 ± 2.4	4.4 ± 1.8

Table 6: Model ablation results for each task. We ablate components of OG-VLA to justify our design choices such as using tiled views, and the contribution of the LLM in the pipeline.

We present ablations on OG-VLA’s design choices in Table 6. The first result shows the effect of tiling the 4 orthographic views instead of feeding them to the IMAGEGENERATOR in batch of 4 as also explored in Genima [25]. For tiling, we stack the 4 views as in 2D array, following prior work. Tiling did not improve the performance for our model as opposed to results reported in the prior works. This may have occurred because the prior work did not have an LLM in the pipeline, so tiling becomes necessary for modeling interactions between views to generate multi-view consistent predictions. However, due to the LLM in OG-VLA conditioning generation in all views, there’s sufficient interaction and reasoning between views through the predicted image tokens (t_a^i , $i \in \{1, \dots, 4\}$).

The second experiment ablates the LLM to study its contribution to the overall performance. We remove the LLM from the system, directly passing image and text representations to the IMAGEGENERATOR. This result shows a drop in performance, highlighting the importance of the LLM in our pipeline.

Next, we study adding tiling to the previous LLM ablation to ensure that interaction between views, which was earlier happening through the LLM, can now take place in the IMAGEGENERATOR. Tiling further leads to a drop in performance, perhaps because of the reduction in the number of tokens used to represent each view. Without tiling, each image is represented by 256 patch tokens, but with tiling all 4 images are represented by 256 tokens in total, potentially creating too tight a representational bottleneck. This version is also similar to Genima with only an IMAGEGENERATOR and tiled camera input views in the pipeline.

Finally, we study the effect of bypassing the prompt and instruction into the IMAGEGENERATOR, shown in the last section of Table 6. The performance drop suggests that some instruction information may have been lost in image tokens output from the LLM, therefore it is beneficial to provide that information separately to the IMAGEGENERATOR for more accurate action prediction.

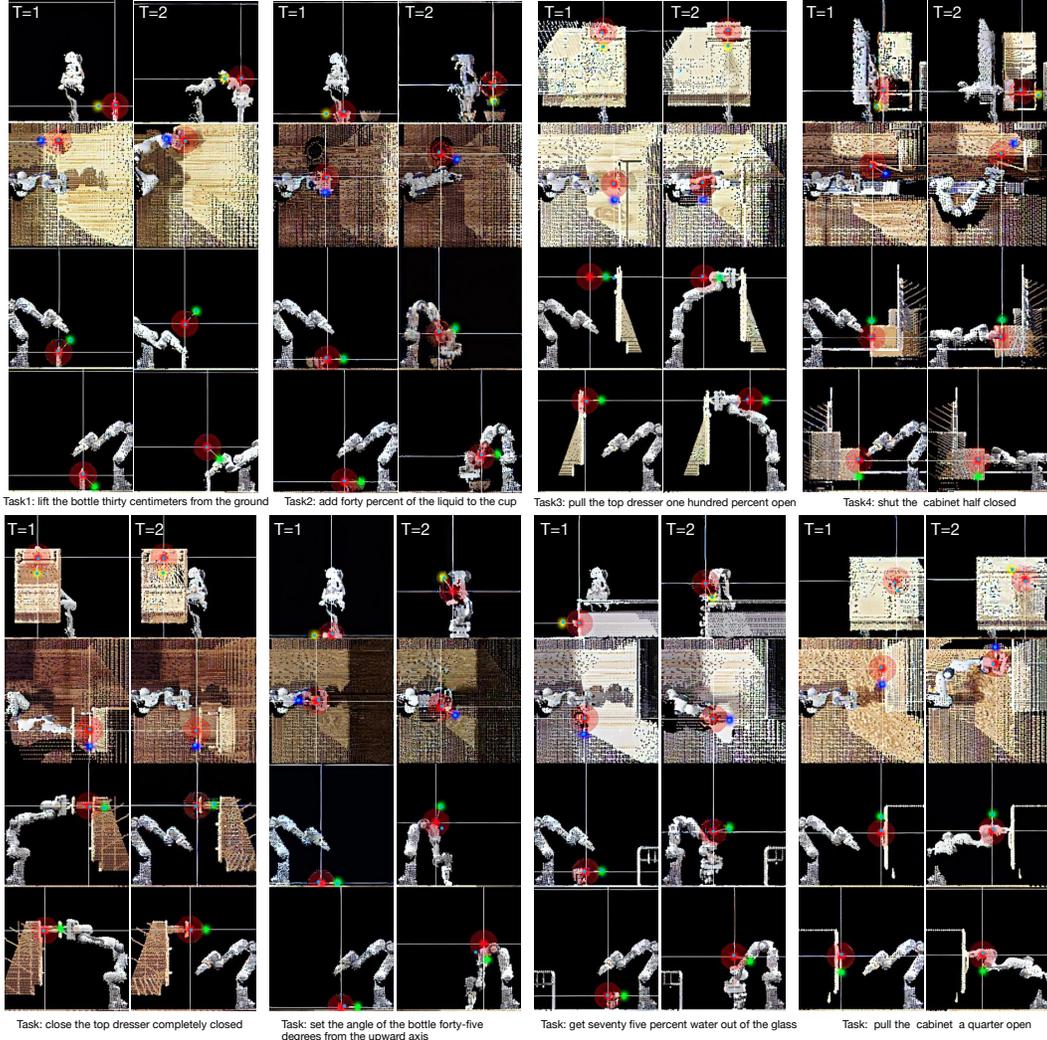


Figure 5: Example gripper position and rotation outputs from OG-VLA for eight different tasks. The rows are the different views: front, top, left, and right. For each task, the two columns are two timesteps required to solve the task. The red Gaussian is the predicted position. The yellow, blue, and green Gaussians are predicted rotation angles along x , z , and y -axis respectively. The blue dot is our model’s output gripper position, back-projected to each view. The dots on rotation Gaussians are showing the extracted pixel for computing the rotation angle in reference to the horizontal right axis in each view.

C.4 Qualitative results

We present qualitative results in Figure 5 of the output of OG-VLA, showing the generated actions, inferred gripper positions and rotations. We show predictions for the tasks: Pickup Object, Transfer Water, Open Drawer, and Close Cabinet. We add remaining task prediction examples in the Appendix. These visualizations show that OG-VLA can do complex translation and rotation tasks, that require both object and free space reasoning. We observe that the translation predictions are quite consistent for most cases, except for Task3 at T=1. The prediction in the Left view is incorrect, however, due to correct predictions in other views, the most likely 3D point (blue dot) is still correctly extracted at the handle of the drawer using the optimization in Equation 2. In Task4, we show

a failure case where the task is to half-close the cabinet, however the predicted point in T=2 is less than the half-way point.

C.5 Ablation with ground truth translation and rotation

	Pickup Object	Reorient Object	Open Drawer	Close Drawer	Open Cabinet	Close Cabinet	Pour Water	Transfer Water	Overall
Training set evaluation	76.7 \pm 4.7	11.7 \pm 2.4	35.0 \pm 0.0	58.3 \pm 4.7	0.0 \pm 0.0	10.0 \pm 4.1	18.3 \pm 6.2	20.0 \pm 4.1	28.8 \pm 0.5
+ Ground Truth Translation	91.7 \pm 2.4	21.7 \pm 14.3	76.7 \pm 6.2	81.7 \pm 2.4	15.0 \pm 7.1	21.7 \pm 4.7	28.3 \pm 2.4	33.3 \pm 2.4	46.2 \pm 0.9
+ Ground Truth Rotation	96.7 \pm 4.7	60.0 \pm 4.1	43.3 \pm 6.2	75.0 \pm 8.2	5.0 \pm 0.0	40.0 \pm 7.1	40.0 \pm 8.2	11.7 \pm 2.4	46.5 \pm 3.8

Table 7: Ablated evaluation of the model’s translation and rotation prediction capabilities on a sampled training set of 20 episodes, similar to the test sets. In each evaluation, we individually ablate either the translation or rotation predictions to assess the model’s prediction capabilities for each.

We present another set of ablations in Table 7 for studying our model’s translation and rotation prediction abilities in an ablated setting on a sampled training set. First, we note that the overall performance of our model on Training and Novel Pose split is quite similar, which indicates that the model has not overfit and is generalizing well w.r.t the training performance. We observe that there is huge scope for improvement, both in predicting more precise translations and rotations from the last two rows of the table. We believe that this gap can be closed with training techniques like co-training with other datasets for better reasoning [42] and leveraging existing large robotics datasets like that in other VLA models [3]. We believe leveraging these datasets with a 3D-aware VLA framework can unlock better learning signals from these large-scale datasets and significantly improve results for our proposed VLA architecture.

D Real world details

We calibrate the camera using MoveIt hand-eye calibration using an ArUco board ¹. We record the trajectory at 30 Hz frame rate. For training OG-VLA, we only use the first 1/5-th of the trajectory before each annotated keyframe action in the trajectory, for augmenting states prior to each keyframe.

Figure 6 shows the different training and test scenes. We provide evaluation videos of successful and failed trajectories for these test scenes on our website.



(a) Training demonstration objects (b) Novel objects used at test time (c) Novel scene distractors used at test time

Figure 6: Real-world setup: objects used at training and test-time. For novel scene, we additionally vary the lighting and background by switching on/off external lighting source, and removing curtains.

¹https://github.com/moveit/moveit_calibration