# VLA-4D: Embedding 4D Awareness into Vision-Language-Action Models for SpatioTemporally Coherent Robotic Manipulation

Hanyu Zhou[1], Chuanhao Ma[2], Gim Hee Lee[1]

[1] School of Computing, National University of Singapore

[2] School of Artificial Intelligence and Automation, Huazhong University of Science and Technology

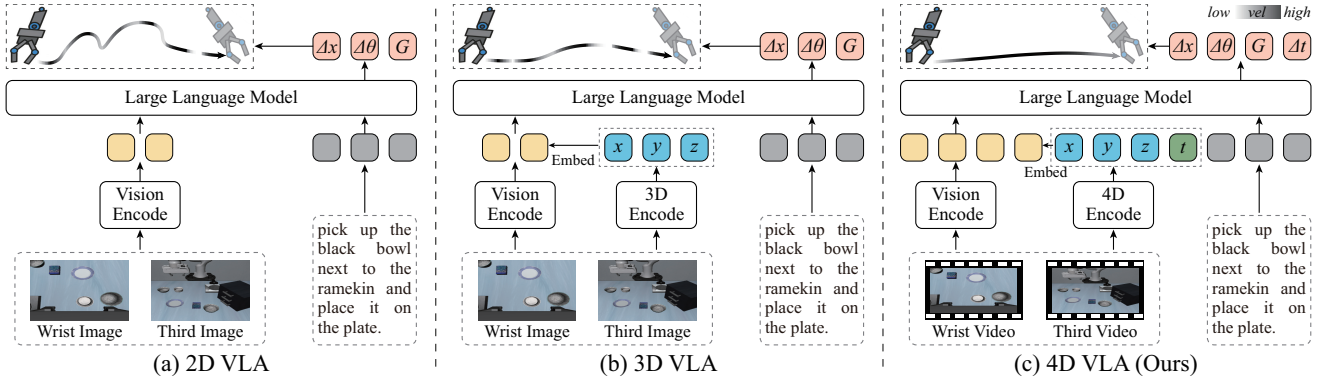{hy.zhou, gimhee.lee}@nus.edu.sg

Figure 1. Illustration of various VLA paradigms for robotic manipulation. (a) 2D VLAs encode visual and language modalities into the LLM to predict actions, which remain spatiotemporally discontinuous. (b) 3D VLAs embed 3D positions into visual representations to improve spatial precision and smoothness of actions, but lack temporal coherence. (c) Our VLA-4D integrates both 3D positions and 1D time into visual representations, and extends the action representation into the spatiotemporal domain to achieve spatiotemporal coherence.

## Abstract

*Vision-language-action (VLA) models show potential for general robotic tasks, but remain challenging in spatiotemporally coherent manipulation, which requires fine-grained representations. Typically, existing methods embed 3D positions into visual representations to enhance the spatial precision of actions. However, these methods struggle to achieve temporally coherent control over action execution. In this work, we propose **VLA-4D**, a general VLA model with 4D awareness for spatiotemporally coherent robotic manipulation. Our model is guided by two key designs: 1) **4D-aware visual representation**. We extract visual features, embed 1D time into 3D positions for 4D embeddings, and fuse them into a unified visual representation via a cross-attention mechanism. 2) **Spatiotemporal action representation.** We extend conventional spatial action representations with temporal information to enable the spatiotemporal planning, and align the multimodal representations into the LLM for spatiotemporal action prediction. Within this unified framework, the designed visual and action representations jointly make robotic manipulation spatially-smooth and temporally-coherent. In addition, we extend the VLA dataset with temporal action annotations for fine-tuning our model. Extensive experiments have been conducted to verify the superiority of our method across different tasks of robotic manipulation.*

## 1. Introduction

With the success of vision-language models (VLMs) in scene reasoning [1–3], recent research explores their use in downstream robotics for manipulation [4–6] and navigation [7, 8]. As shown in Fig. 1 (a), this has led to a specialized paradigm of vision-language-action (VLA) models that connect visual reasoning to action planning. Despite strong results on general robotic tasks [5, 9–11], VLAs still struggle with spatiotemporally coherent manipulation that demands fine-grained representations. The key issues are: single-image inputs induce coarse visual reasoning, and 2D–3D (image-robot) coordinate system mismatches degrade action precision. Therefore, our goal is to strengthen fine-grained representations for visual reasoning and action planning.

Most existing VLA models [12–14] embed 3D positional cues into visual features to strengthen spatial reasoning, which improves action precision and smoothness as shown in Fig. 1 (b). However, action planning is inherently spatiotemporal and demands continuity over time. Consequently, these 3D VLAs often struggle with fine-grained temporal control that leads to behaviors such as idle pauses or jitter. Recent works [15, 16] add temporal signals such as frame indices and fuse them with 3D position embeddings in the visual stream. Although this helps temporal reasoning such as state precedence, it does not directly enforce temporally coherent action plans. We argue that coherent manipulation requires jointly enhancing spatiotemporal perception in both the visual and action representations within a VLA.

To address these issues, we propose to embed distinct spatiotemporal information into both the visual representation for reasoning and the action representation for planning. For visual representation, we suggest that 3D positional information enhances understanding of scene geometry and spatial localization of subsequent actions, and 1D temporal information further improves the perception of dynamic patterns and temporal states of actions. This indicates that 4D spatiotemporal cues play a crucial role in promoting both visual reasoning and action planning. For action representation, we observe that spatial control serves as the foundation for fine-grained planning, and temporal control is indispensable for achieving a coherent action process. This motivates us to augment conventional spatial actions with temporal information for coherent spatiotemporal operations. Consequently, these enhanced visual and action representations jointly unlock the potential of VLA models for spatiotemporally coherent robotic manipulation.

In this work, we propose **VLA-4D**, a general vision-language-action model with 4D awareness for spatiotemporally coherent robotic manipulation. As illustrated in Fig. 1 (c), our VLA-4D relies on a 4D-aware visual representation and a spatiotemporal action representation. For the 4D-aware visual representation, we first encode the input video sequence into visual features and geometric features. Within the geometric space, we encode 3D positions and 1D time into 4D spatiotemporal embeddings. We then fuse the 4D embeddings into the visual features via a cross-attention mechanism. For the spatiotemporal action representation, we extend the conventional spatial control parameters by introducing additional varying temporal variables to enable fine-grained spatiotemporal planning of robotic actions. By performing multimodal alignment, the LLM produces spatiotemporal actions for the robot. Under this unified framework, the 4D-aware visual representation and the spatiotemporal action representation jointly ensure the spatial smoothness and temporal coherence for robotic manipulation. In addition, we augment the LIBERO [17] dataset with temporal action annotations to fine-tune our model for more

effective spatiotemporal robotic manipulation. Our main contributions are summarized as follows:

- We propose **VLA-4D**, a general 4D vision-language-action model for spatiotemporally coherent robotic manipulation. Our model embeds spatiotemporal information into both visual and action representations for fine-grained visual reasoning and action planning.
- We design an explicit 4D-aware visual representation. A cross-attention mechanism fuses 3D positions and 1D time into visual features for stronger fine-grained spatiotemporal reasoning of our model.
- We construct a spatiotemporal action representation. Temporal control variables are incorporated into conventional spatial actions to improve the spatial smoothness and temporal coherence of robotic operations.
- We extend the VLA dataset with temporal action annotations to fine-tune our model. Extensive experiments demonstrate that our framework achieves state-of-the-art performance across various robotic manipulation tasks.

## 2. Related Work

**Vision-Language-Action Models.** Vision-language models [1, 3, 18–20] have achieved remarkable success in scene reasoning, inspiring a new class of end-to-end models for embodied intelligence, *i.e.*, vision-language-action models. VLAs [5, 9, 21–25] often encode actions as sequences of discrete, text-like tokens. These tokens are co-embedded with visual and linguistic inputs in a large language model [26, 27] for a motion-planning paradigm grounded in scene reasoning for robotic manipulation. Although these VLA models indeed show great potential for general robotic tasks, they still face challenges in achieving spatiotemporally coherent manipulation. Robotic operations demand precise and continuous action parameters. In contrast, the knowledge learned directly from images and text is often coarse-grained. This mismatch makes it difficult to plan motions that are both accurate and temporally coherent. In this work, we focus on enhancing the fine-grained representation for visual reasoning and action planning.

**3D Vision-Language-Action Models.** For VLA models, the key to improving the action precision lies in improving the fine-grained visual reasoning. Typically, existing VLA models [12–14, 28, 29] embed 3D positional information into visual representations to strengthen spatial reasoning, thus optimizing the spatial control parameters of actions. Although these methods may achieve a relatively optimal spatial trajectory, they are still limited in achieving fine-grained reasoning and control along the temporal dimension, such as sequential state chaos and potential operational idling. We thus propose to embed 3D positions and 1D time into the visual representation to strengthen spatiotemporal reasoning and improve action precision in robotic manipulation.
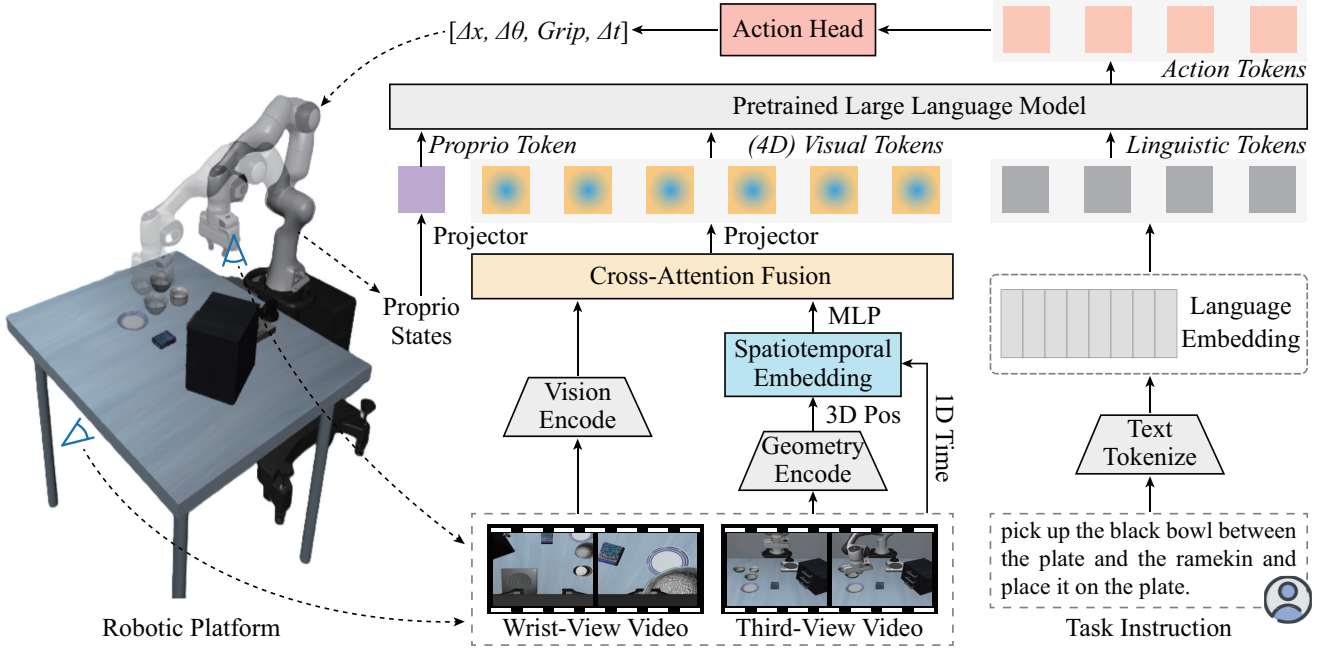
Figure 2. Our VLA-4D consists of two stages: 1) **4D-aware visual representation.** Encode 3D positions and 1D time into 4D spatiotemporal embeddings, and fuse them into visual features via a cross-attention mechanism. 2) **Spatiotemporal action representation.** Extend action parameters into the spatiotemporal domain, and align multimodal representations into the LLM for robotic action prediction.

**4D Vision-Language-Action Models.** The research line most closely related to our work is 4D VLA models [15, 16, 30]. These 4D models integrate spatial positional and temporal information into visual representations to further improve the prediction precision of robotic actions. For example, Niu *et al.* [15] learn 4D point trajectories from videos to model visual representations and transfer the learned distribution to robotic control for action prediction via an autoregressive model. Zhang *et al.* [16] encode frame indices from past videos and fuse them with 3D positional embeddings in the visual stream. This mitigates temporal state ambiguity during robotic execution. Although these methods strengthen fine-grained spatiotemporal visual reasoning, they still struggle to explicitly improve the precision of temporal action planning and the overall coherence of operation execution. In contrast, we posit that it is essential to enhance the 4D awareness of visual representations, and to represent the spatiotemporal control parameters of robotic actions for spatiotemporally coherent robotic manipulation.

## 3. Our VLA-4D

**Overview.** Fig. 2 shows the architecture of our VLA-4D with two key stages: 1) *4D-Aware Visual Representation,* where visual features are enhanced by 4D spatiotemporal embeddings with cross-attention fusion. 2) *SpatioTemporal Action Representation,* where action parameters are extended to the spatiotemporal domain and aligned with multimodal representations for task optimization. These components

enable VLAs to achieve finer-grained spatiotemporal visual reasoning and action planning for robotic manipulation.

**Our Framework.** Given a video sequence $I$ and instruction texts, we employ a vision encoder to extract visual features $f_v$ and a text tokenizer to extract linguistic tokens $\tau_l$.

**1)** *4D-Aware Visual Representation (cf. Sec. 3.1).* This stage enhances 4D perception of visual reasoning. We first encode the video sequence into a geometric latent space to obtain 3D positions $p_{3D}$. We then encode 1D time $t$ and 3D positions $p_{3D}$ into 4D learnable embeddings $f_{4D}$, which are fused with visual features through a cross-attention mechanism:

$$f_{4D} = \text{STE}(p_{3D}, t), \quad f_v^{4D} = \text{CAtt}(f_v, f_{4D}), \quad (1)$$

where STE denotes the spatiotemporal embedding operation. The unified visual representation $f_v^{4D}$ can perceive the 4D semantic and geometry of the scene for visual reasoning.

**2)** *SpatioTemporal Action Representation (cf. Sec. 3.2).* This stage extends action planning into the spatiotemporal dimension. We formulate the robotic action representation along the spatiotemporal dimension: $A = [X, T]$, where $X, T$ denote spatiotemporal control variables. The fused visual features $f_v^{4D}$ from the previous stage and the proprioceptive states $f_p$ are projected into language embedding space: $[\tau_v^{4D}, \tau_p] = \text{Proj}(f_v^{4D}, f_p)$, where $\tau_v^{4D}, \tau_p$ are the visual and proprioceptive tokens aligned with linguistic tokens $\tau_l$. The LLM with action head then learns the mapping from the aligned multimodal tokens to actions: $A = \mathcal{LLM}(\tau_v^{4D}, \tau_p, \tau_l)$.
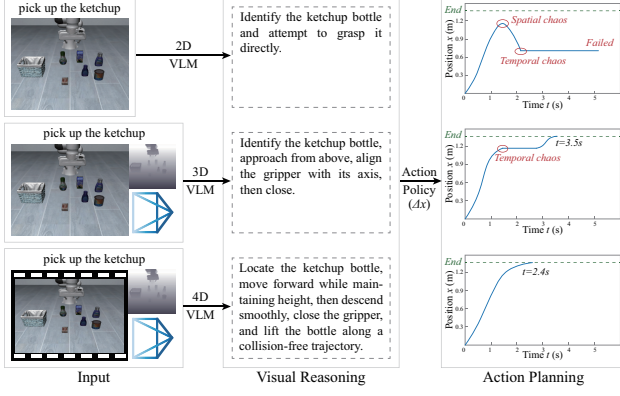
3

Figure 3. Effect of different visual representations. 3D spatial information enhances the understanding of scene geometry and subsequent action localization, while 1D temporal information further ensures the dynamic perception and temporal action state.

**Remarks.** The output $A$ denotes the fine-grained spatiotemporal action parameters. Our unified framework enhances the 4D awareness of visual reasoning and enables spatiotemporally coherent action planning for robotic manipulation. The following sections detail the design of each stage.

### 3.1. 4D-Aware Visual Representation

Existing VLA models utilize images and instructions to directly learn visual reasoning to enable action planning for robotic manipulation. However, the visual knowledge remains relatively coarse-grained and carries the risk of reducing the precision of the predicted action control. In this section, we explore how to enhance fine-grained visual representations to achieve precise robotic action planning.

**4D SpatioTemporal Embedding.** Visual and action representations within VLAs exhibit certain discrepancies in the spatial and temporal dimensions. There is a coordinate system mismatch in the spatial dimension, where the image is in a 2D coordinate system and the action is in a 3D world (or robot) coordinate system. Furthermore, there is also a disagreement in the temporal dimension, where the image is a static snapshot of a scene and the action forms a sequence of 3D trajectories over time. The discrepancies in the spatial-temporal domains led us to consider the unification of the coordinate systems with temporal cues. To this end, we analyze the effects of visual inputs corresponding to different coordinate systems (2D pixel space *v.s.* 3D world space) and temporal settings on visual reasoning and action planning. As shown in Fig. 3, 3D coordinate information enhances the understanding of scene geometry and spatial localization of subsequent actions, while 1D temporal information promotes the dynamic perception and temporal action state. Inspired by this observation, we conjecture that visual reasoning endowed with 4D awareness can facilitate fine-grained action planning. Given third-view and

wrist-view video sequences, we adopt VGGT [31] as the geometry encoder to extract the camera pose $P$ and depth $D$ at a timestamp $t$. Combined with intrinsic parameter $K$, we transform 2D pixel coordinate $p_{2D}$ to world (or robot) coordinate system via geometric projection [32, 33]:

$$p_{3D} = P^{-1}(DK^{-1}p_{2D}). \quad (2)$$

After traversing all timestamps for geometric extraction, we propose a spatiotemporal embedding operation STE to integrate 3D positions and 1D time. Specifically, we introduce a Fourier-based encoding strategy [34] to convert positions and timestamps into learnable patterns, and map them into a 4D representation through a linear layer:

$$\psi(x) = 1/\sqrt{d}\,[cos(xW_r^\top) \,||\, sin(xW_r^\top)],$$
$$f_{4D} = w_p \cdot [\psi(p_{3D}) \,||\, \psi(t)], \quad (3)$$

where $d$ denotes the dimension and $W_r$ is the learnable parameter of the Fourier feature. The resulting $f_{4D}$ represents the 4D spatiotemporal geometry of the scene.

**Cross-Attention Fusion.** In our model, we employ a ViT variant [35] as the vision encoder to extract high-level visual features with semantic information. However, such standalone features cannot be localized to the world coordinate system or aligned with the temporal sequence state required for robotic manipulation. This motivates us to further embed 4D awareness into the visual features. We first introduce a Multi-Layer Perceptron (MLP) to make the dimension of the 4D spatiotemporal embeddings the same as the dimension of the visual features: $\hat{f}_{4D} = \text{MLP}(f_{4D})$. Next, we fuse the 4D spatiotemporal embeddings into the visual features via a cross-attention mechanism [36, 37]:

$$q = w_q f_v,\ k = w_k \hat{f}_{4D},\ v = w_v \hat{f}_{4D},$$
$$f_v^{4D} = f_v + \text{softmax}(qk^\top\sqrt{d})v, \quad (4)$$

where $w$ is a learnable weight. The fused result $f_v^{4D}$ denotes the unified visual representation that can capture both the 4D semantic and geometric characteristics of the scene for spatiotemporal visual reasoning.

### 3.2. SpatioTemporal Action Representation

Ideally, robotic motion planning should be a spatially smooth and temporally coherent dynamic process. Although 4D-aware visual representations facilitate fine-grained action prediction, they are not sufficient to explicitly ensure the overall coherence of robotic operations. Therefore, our goal is to build the fine-grained spatiotemporal action representation for coherent robotic manipulation.

**Spatiotemporal Action Definition.** Conventional action representations [5, 11] focus mainly on spatial control parameters such as translation for fine-grained robotic operations. However, such representations ignore the specific execution duration and enable the robot to learn only spatial movements within an action step. This omission increases
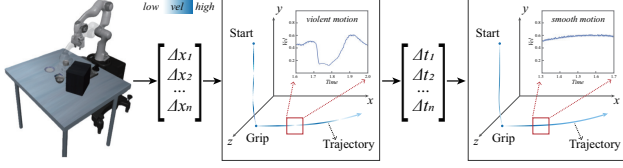
Figure 4. Illustration of spatiotemporal action representation. Spatial parameters enable fine-grained action planning, while temporal parameters further improve the action coherence during execution.

the risk of premature termination or delayed responses in robotic tasks, which leads to a discontinuous control process. As illustrated in Fig. 4, we analyze the effects of spatial and temporal control parameters on the robotic operation process. Note that the variable time here is predefined according to the parameters of the robot and task configurations. We can observe that spatial actions form the foundation of fine-grained robotic operations, and the incorporation of reasonable temporal parameters further promotes the continuity and coherence of the overall action execution process. Driven by this motivation, we augment conventional spatial action representations with temporal information. Specifically, the conventional spatial action is defined as $X = [\Delta x, \Delta \theta, Grip]$, where $\Delta x$ denotes the translational displacement of the end-effector, $\Delta \theta$ represents its rotational change in orientation, and $Grip$ specifies the gripper control signal indicating the open or close state. We further formulate a temporal action representation: $T = \Delta t$, which is a time variable for step-level action control jointly determined by the scene content *corresponding to vision in VLA*, the operation task *corresponding to language in VLA*, and the robotic embodiment *corresponding to proprioceptive state feedback*. In this way, we extend the spatial action into a spatiotemporal action: $A = [\Delta x, \Delta \theta, Grip, \Delta t]$.

**Multimodal Alignment and Optimization.** Similar to conventional VLMs [1, 38–40], our VLA model produces spatiotemporal actions through two key components: multimodal representation alignment and task optimization. For multimodal alignment, we align vision, language, and additional proprioceptive states into a unified learnable space. Considering that the input to a large language model requires text-like tokens, we first introduce a multi-layer perceptron as the projection function $\mathrm{Proj}(\cdot)$ to map the 4D-aware visual features and proprioceptive states into the language embedding space: $\tau_v^{4D} = \mathrm{MLP}(f_v^{4D})$, $\tau_p = \mathrm{MLP}(f_p)$, where $\tau_v^{4D}$ and $\tau_p$ denote the visual and proprioceptive tokens. Next, we tokenize the input instruction into the same language embedding space to obtain the linguistic tokens $\tau_l$.

For task optimization, our goal is to learn a model that maps the aligned multimodal representations to spatiotemporal actions. We first introduce a pretrained large language model $\mathcal{T}(\cdot)$ and append an MLP-based action head $\mathcal{H}(\cdot)$. We then concatenate the visual, linguistic, and proprioceptive

tokens, and feed them into the model for action prediction:

$$[\Delta x, \ \Delta \theta, \ Grip, \ \Delta t] = \mathcal{H}(\mathcal{T}([\tau_v^{4D}, \ \tau_p, \ \tau_l])). \quad (5)$$

To train the model for precise spatiotemporal action prediction, we further apply an L1-norm loss [41] on the predicted spatiotemporal action variables:

$$\mathcal{L}_{action} = \sum (|\Delta x - \tilde{\Delta x}|_1 + |\Delta \theta - \tilde{\Delta \theta}|_1 \\ + |Grip - \tilde{Grip}|_1 + |\Delta t - \tilde{\Delta t}|_1), \quad (6)$$

where variables with a tilde ( ˜ ) denote ground-truth action values. During inference, our VLA-4D can directly predict future spatiotemporal actions based on historical videos, proprioceptive states, and task-specific textual instructions for fine-grained and coherent robotic manipulation.

## 4. Dataset and Training Pipeline

### 4.1. Robotic Dataset

Many robotic datasets [17, 43–45] have been proposed to evaluate the performance of VLAs. However, there is currently no dataset specifically designed for spatiotemporally coherent robotic manipulation in VLAs. To address this limitation, we select a representative robotic dataset [17] and extend its input modalities and action annotations to fine-tune our model for evaluating robotic tasks.

**LIBERO.** LIBERO [17] is a simulation suite with four benchmark settings designed to advance lifelong learning in robotic manipulation, including spatial reasoning (LIBERO-Spatial), object understanding (LIBERO-Object), task goal (LIBERO-Goal) and long-horizon planning (LIBERO-Long). Note that LIBERO-Average denotes the whole dataset for task evaluation. For the input modalities, we render human-designed trajectories in a simulated environment at a fixed sampling frequency to obtain camera parameters, multi-view videos with timestamps, depths, and proprioceptive states. For the action annotations, without altering the original spatial action parameters, we manually select action chunks that exhibit a consistent motion trend and convert their step counts to variable temporal action annotations based on the sampling frequency. After data cleaning and selection, the final dataset contains 40 subtasks with a total of 150k paired vision–language–action samples.

### 4.2. Training Pipeline

Our VLA-4D first loads the weights of several pretrained models [31, 35] to initialize its multimodal understanding capability, and then performs fine-tuning for robotic tasks. To fully enhance the learning capability of visual representations and improve the performance of action representations, we divide the training process into two stages:

**Stage 1: 4D Vision-Language Alignment.** This stage trains the 4D visual representations for reasoning within our VLM architecture. We select several 3D and 4D vision-language

Table 1. Quantitative results of VLAs for fine-tuned robotic manipulation tasks on the LIBERO benchmark.

| Methods | | Spatial | | Object | | Goal | | Long | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Succ. rate(%)↑ | Time(s)↓ | Succ. rate(%)↑ | Time(s)↓ | Succ. rate(%)↑ | Time(s)↓ | Succ. rate(%)↑ | Time(s)↓ | Succ. rate(%)↑ | Time(s)↓ |
| 2D | OpenVLA [5] | $84.7 \pm 0.9$ | 5.5 | $88.4 \pm 0.8$ | 7.5 | $79.2 \pm 1.0$ | 6.1 | $53.7 \pm 1.3$ | 13.1 | $76.5 \pm 0.6$ | 8.1 |
| | Octo [21] | $78.9 \pm 1.0$ | 5.7 | $85.7 \pm 0.9$ | 6.9 | $84.6 \pm 0.9$ | 6.3 | $51.1 \pm 1.3$ | 9.3 | $75.1 \pm 0.6$ | 7.1 |
| | DiffusionPolicy [11] | $78.3 \pm 1.1$ | 6.4 | $92.5 \pm 0.7$ | 7.8 | $68.3 \pm 1.2$ | 6.4 | $50.5 \pm 1.3$ | 15.2 | $72.4 \pm 0.7$ | 8.7 |
| | CogACT [42] | $87.5 \pm 0.9$ | 5.4 | $90.2 \pm 1.1$ | 6.8 | $78.4 \pm 0.8$ | 5.9 | $53.2 \pm 1.2$ | 10.7 | $76.5 \pm 0.9$ | 7.0 |
| 3D | TraceVLA [12] | $84.6 \pm 0.2$ | – | $85.2 \pm 0.4$ | – | $75.1 \pm 0.3$ | – | $54.1 \pm 1.0$ | – | $74.8 \pm 0.4$ | – |
| | SpatialVLA [13] | $88.2 \pm 0.5$ | 5.3 | $89.9 \pm 0.7$ | 6.4 | $78.6 \pm 0.6$ | 5.9 | $55.5 \pm 1.0$ | 8.9 | $78.1 \pm 0.7$ | 6.8 |
| 4D | 4D-VLA [16] | $88.9 \pm 0.5$ | – | $95.2 \pm 0.3$ | – | $90.9 \pm 0.4$ | – | $79.1 \pm 1.2$ | – | $88.6 \pm 0.3$ | – |
| | VLA-4D (Ours) | $\mathbf{97.9 \pm 0.2}$ | $\mathbf{4.1}$ | $\mathbf{98.6 \pm 0.3}$ | $\mathbf{5.6}$ | $\mathbf{97.8 \pm 0.3}$ | $\mathbf{4.6}$ | $\mathbf{94.8 \pm 0.8}$ | $\mathbf{6.9}$ | $\mathbf{97.4 \pm 0.3}$ | $\mathbf{5.8}$ |



pick up the black bowl next to cookie box and place it on plate

turn on the stove and put the moka pot on it

open the middle drawer of the cabinet

put the bowl on the stove

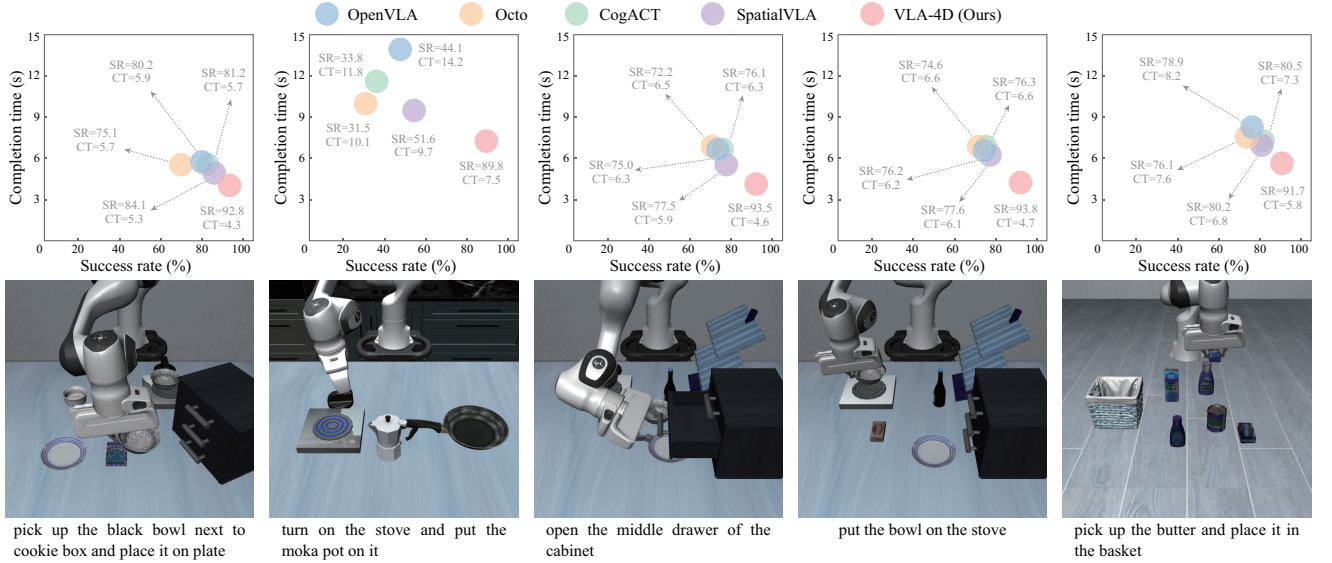pick up the butter and place it in the basket

Figure 5. Quantitative comparison of VLAs on zero-shot robotic manipulation tasks.

datasets including Scan2Cap [46], ScanQA [47], ScanRef [48], Multi3DRefer [49], Chat4D [20] to train the VLM component of our model using only the LLM loss [26, 50]. This ensures that the visual representations acquire strong 4D spatiotemporal perception and can interact with language. At this stage, we update the weights of cross-attention, spatiotemporal embedding, projector, and optimize the LLM component using LoRA adapters [51] while freezing the vision encoder, geometry encoder, and action head.

**Stage 2: Robotic Task Fine-Tuning.** This stage trains the spatiotemporal action representations for planning within our entire VLA architecture. We employ the modified LIBERO dataset with spatiotemporal action annotations as the training settings to fine-tune our model using the loss $\mathcal{L}_{action}$. This enables our VLA-4D to produce a series of spatiotemporal actions conditioned on task instructions and 4D vision inputs for diverse robotic manipulation tasks. At this stage, we update the weights of the action head and projector, optimize cross-attention and LLM components using LoRA adapters [51] while freezing the remaining modules.

## 5. Experiments

### 5.1. Experiment Setup

**Implementation Details.** Our VLA-4D model utilizes the pre-trained weights of Qwen2.5-VL-7B [35] as the VLM (vision encoder + LLM) backbone, and VGGT [31] as the geometry encoder. The cross-attention fusion module is a Transformer-based network architecture. The whole model is trained on 8 RTX 6000 Ada GPUs using AdamW as the optimizer. In training stage 1, we set the learning rate to $1.0e - 4$ with a batch size of 16. In training stage 2, we set the learning rate to $5.0e - 5$ with a batch size of 24.

**Comparison Methods.** We compare our model with several typical VLAs: OpenVLA [5], Octo [21], CogACT [42], DiffusionPolicy [11], TraceVLA [12], SpatialVLA [13], and 4D-VLA [16]. The first four methods are 2D VLA models using single images as inputs, the next two methods are 3D VLA models using images and 3D positions as inputs, and the last method is a 4D VLA model using images, 3D positions, and 1D time as inputs.
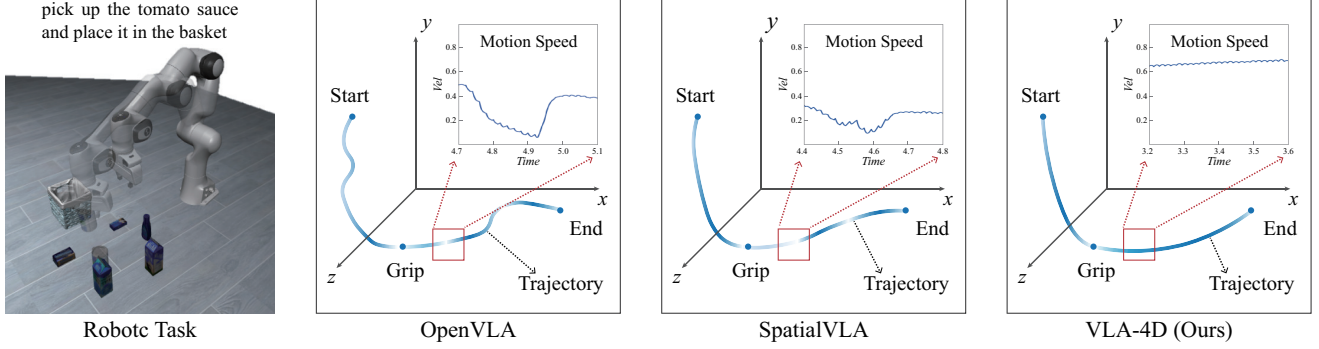
Figure 6. Visual comparison of VLAs on spatiotemporal action planning.

Table 2. Effect of visual representation modules.

| Spatial embed | Temporal embed | Feature fusion | LIBERO-Spatial | | LIBERO-Goal | |
|---|---|---|---|---|---|---|
| | | | Succ(%)↑ | Time(s)↓ | Succ(%)↑ | Time(s)↓ |
| × | × | × | $89.4 \pm 0.6$ | 5.7 | $90.1 \pm 0.7$ | 6.3 |
| ✓ | × | × | $92.2 \pm 0.4$ | 5.1 | $94.3 \pm 0.5$ | 5.6 |
| ✓ | ✓ | × | $96.5 \pm 0.3$ | 4.4 | $95.7 \pm 0.4$ | 4.9 |
| ✓ | ✓ | ✓ | $\mathbf{97.9 \pm 0.2}$ | $\mathbf{4.1}$ | $\mathbf{97.8 \pm 0.3}$ | $\mathbf{4.6}$ |

Table 3. Effect of action representation components.

| Action representation | LIBERO-Spatial | | LIBERO-Goal | |
|---|---|---|---|---|
| | Succ(%)↑ | Time(s)↓ | Succ(%)↑ | Time(s)↓ |
| Spatial param. | $96.8 \pm 0.3$ | 5.0 | $97.1 \pm 0.3$ | 5.7 |
| Spatial + Temporal param. | $\mathbf{97.9 \pm 0.2}$ | $\mathbf{4.1}$ | $\mathbf{97.8 \pm 0.3}$ | $\mathbf{4.6}$ |

**Evaluation Metric.** We compare all competing methods on the LIBERO dataset. We set two types of evaluation settings. The first directly fine-tunes competing models and evaluates them across all robotic subtasks, and the second evaluates the zero-shot generalization performance of competing models on a selected subset of subtasks. Moreover, we adopt task success rate (SR) and completion time (CT) as primary evaluation metrics. Note that the experiments in Sec. 5.3 and Sec. 5.4 are evaluated on the LIBERO-Spatial and LIBERO-Goal benchmarks.

## 5.2. Comparison with State-of-the-Art Models

**Comparison on Fine-Tuned Tasks.** In Tab. 1, we compare the fine-tuning performance of competing methods on the LIBERO benchmark. We draw three key conclusions from the results. First, 2D models exhibit relatively inferior performance on complex real-world robotic manipulation tasks while 3D and 4D models achieve significantly better results. Second, 4D models consistently outperform all other methods. The main reason is 4D models add 3D positions to align visual and action coordinate frames, and 1D time to resolve temporal ambiguity in planning. This joint space–time encoding stabilizes execution and reduces spurious or unstable actions. Third, our model achieves a higher task success rate than other competing models and requires the shortest task completion time. Overall, our VLA-4D provides an effective paradigm for spatiotemporal robotic manipulation.

**Comparison on Zero-Shot Tasks.** In Fig. 5, we compare the generalization performance of different methods on zero-shot tasks. The results show that our model achieves substan-

tially higher success rates and shorter task completion times than the competing models across multiple zero-shot tasks. These results demonstrate that incorporating spatiotemporal representations into both the visual and action modalities enables our model to maintain strong generalization performance, even on unseen robotic tasks.

**Comparison on SpatioTemporal Planning.** In Fig. 6, we compare the predicted spatiotemporal action trajectories produced by representative 2D, 3D, and 4D models: OpenVLA [5], SpatialVLA [13] and our VLA-4D. The results show that the 2D model exhibits substantial redundant global motion and pronounced oscillations in local motion speed. The 3D model generates much smoother global trajectories, but the local motion speed still fluctuates noticeably. In contrast, the 4D model produces both smooth global trajectories and stable local motion speeds. These results demonstrate that our VLA-4D provides a new and reliable paradigm for achieving spatiotemporally coherent robotic manipulation.

## 5.3. Ablation Study

**Effect of Visual Representation Modules.** In Tab. 2, we validate the effectiveness of spatial embedding, temporal embedding, and feature fusion modules in the visual representation. The results show that both the spatial and temporal embeddings are key to significantly improving the overall performance of our model by a large margin. Furthermore, the feature fusion further enhances the upper limit of the performance of our model on robotic manipulation tasks.

**Effect of Action Representation Components.** In Tab. 3, we compare the effects of spatial and temporal components in the action representation. When only the spatial action representation is introduced, our model already achieves

Table 4. Impact of various input modalities on robotic manipulation.

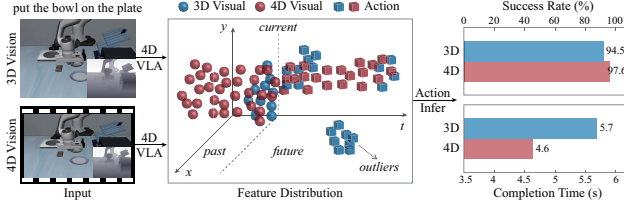| Vision data | 4D cues | Proprio. | LIBERO-Spatial | | LIBERO-Goal | |
|---|---|---|---|---|---|---|
| | | | Succ(%)↑ | Time(s)↓ | Succ(%)↑ | Time(s)↓ |
| Image | × | × | 85.9 ± 0.6 | 5.9 | 88.0 ± 0.8 | 6.5 |
| Video | × | × | 89.2 ± 0.6 | 5.7 | 90.1 ± 0.7 | 6.3 |
| Video | √ | × | 97.1 ± 0.2 | 4.1 | 97.3 ± 0.4 | 4.6 |
| Video | √ | √ | **97.9 ± 0.2** | **4.1** | **97.8 ± 0.3** | **4.6** |



Figure 7. Influence of visual representations on actions. Compared with 3D features, 4D visual features allow the VLA to exploit past visual cues, producing accurate and coherent actions.

Table 5. Discussion on spatiotemporal encoding.

| 4D cues | LIBERO-Spatial | | LIBERO-Goal | |
|---|---|---|---|---|
| | Succ(%)↑ | Time(s)↓ | Succ(%)↑ | Time(s)↓ |
| w/o Spatiotemporal encoding | 96.3 ± 0.4 | 4.3 | 96.1 ± 0.3 | 4.9 |
| w/ Spatiotemporal encoding | **97.9 ± 0.2** | **4.1** | **97.8 ± 0.3** | **4.6** |

Table 6. Impact of various visual feature fusion strategies.

| Fusion strategy | LIBERO-Spatial | | LIBERO-Goal | |
|---|---|---|---|---|
| | Succ(%)↑ | Time(s)↓ | Succ(%)↑ | Time(s)↓ |
| Concatenation | 94.2 ± 0.4 | 4.5 | 94.9 ± 0.4 | 4.8 |
| Weighting | 95.8 ± 0.3 | 4.3 | 96.5 ± 0.3 | 4.8 |
| Attention | **97.9 ± 0.2** | **4.1** | **97.8 ± 0.3** | **4.6** |

Table 7. Comparison between different training strategies.

| Training strategy | LIBERO-Spatial | | LIBERO-Goal | |
|---|---|---|---|---|
| | Succ(%)↑ | Time(s)↓ | Succ(%)↑ | Time(s)↓ |
| Only Stage 2 | 91.2 ± 0.5 | 4.9 | 90.7 ± 0.6 | 5.3 |
| Stage 1 + Stage 2 | **97.9 ± 0.2** | **4.1** | **97.8 ± 0.3** | **4.6** |

a relatively high task success rate, but with a higher task completion time. After further incorporation of temporal representation, the success rate improves slightly, but the task completion time decreases significantly. These results indicate that the spatial parameters determine the effectiveness of action planning, and the temporal parameters further enhance its efficiency and superiority.

**Influence of Input Modality.** In Tab. 4, we conduct an ablation study on the impacts of different input modalities: vision data (*e.g.*, images and videos), 4D cues, and proprioceptive states. When only image-based vision input is used, our model achieves a moderate task success rate. When switching to video-based vision inputs, the success rate improves slightly. With the further introduction of 4D cues, both the task success rate and the completion time improve significantly. When incorporating proprioceptive states, our model yields a minor but consistent overall performance gain. These results suggest that video-based vision data ensure a preliminary baseline capability, 4D cues substantially enhance model performance, and proprioceptive states provide an additional complementary improvement.

### 5.4. Discussion

**How 4D Vision Affects Spatiotemporal Action?** In Fig. 7, we analyze the interaction mechanism between 4D-aware visual representations and spatiotemporal action representations through feature distribution visualization. Specifically, we compare the 3D visual representations (derived from image+depth input) and their corresponding action representations with the 4D visual representations (derived from video+depth input) and their corresponding action representations under the same task settings. First, the 3D visual feature distribution lacks the temporal dimension. As a re-

sult, the corresponding action features are spatially clustered but temporally scattered. This leads to robotic operations that achieve a reasonable success rate but incur a longer completion time. Second, the 4D visual feature distribution forms a continuous spatiotemporal manifold, and the corresponding action features are spatiotemporally clustered and continuous. Consequently, the robot achieves the highest task success rate with a reduced completion time. Within our VLA-4D framework, spatial knowledge embedded in 4D visual features enhances spatial localization of robotic actions, and temporal knowledge improves their temporal efficiency during manipulation.

**Role of Spatiotemporal Encoding.** We study the importance of spatiotemporal encoding on positions and timestamps for our model. As demonstrated in Tab. 5, the model equipped with spatiotemporal encoding achieves better performance on robotic manipulation tasks compared to the one without it. This is because spatial positions and temporal information operate on different scales, which can lead to gradient conflicts during training. In contrast, a dedicated encoding strategy can map spatiotemporal information into a unified feature space to ensure stable convergence and prevention of feature space collapse.

**Choice of Visual Feature Fusion Strategy.** Tab. 6 compares the impact of different visual feature fusion strategies: concatenation, weighting, and attention. The results show that attention-based fusion outperforms other fusion strategies. This is because concatenation and weighting rely on global unified fusion with fixed weights. In contrast, attention-based fusion can dynamically adjust the fusion weights of visual features according to 4D spatiotemporal embeddings.

This allows our model to focus more on meaningful 4D-aware visual features for action prediction.

**Impact of Training Strategy.** In Tab. 7, we study the effects of different training strategies on our model. Compared to using only stage 2 for fine-tuning, the multi-stage training strategy that combines stage 1 and stage 2 leads to a substantial improvement in model performance on robotic manipulation tasks. The underlying reason is that stage 1 training enhances the spatiotemporal visual reasoning ability through large-scale 3D and 4D vision-language pretraining, facilitating more accurate spatiotemporal action prediction for coherent robotic planning.

**Limitation.** Our VLA-4D can predict fine-grained action parameters for spatiotemporally coherent robotic manipulation, but still faces challenges when deployed in unseen real-world environments. The main reason is that robotic operations in unseen environments may introduce action errors due to uncontrollable factors such as mechanical wear and calibration drift, which in turn reduce the overall efficiency of manipulation. In the future, we plan to incorporate reinforcement learning [52] to correct the errors online in predicted spatiotemporal actions, enabling more robust and adaptive robotic planning in real-world environments.

# 6. Conclusion

In this work, we propose VLA-4D, a general 4D vision-language-action model for spatiotemporally coherent robotic manipulation. We design a 4D-aware visual representation that fuses 3D positions and 1D time into visual features for fine-grained spatiotemporal reasoning. Moreover, we construct a spatiotemporal action representation that incorporates temporal variables into spatial parameters for fine-grained spatiotemporal planning. By performing multimodal alignment, our model can produce precise and coherent actions for the robot. To support training, we extend the representative VLA dataset with temporal action annotations to fine-tune our model for spatiotemporal robotic manipulation. Extensive experiments and ablations on benchmark datasets verify the superiority of our method.

# References

[1] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. 1, 2, 5

[2] Mandy Chen, Adams Wei Yu, Hamid Palangi, Paul Smolensky, Yinfei Yang, Xiaowei Yuan, Kathy Meier-Hellstern, Jianfeng Gao, Ed Chi, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2303.07892*, 2023.

[3] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Arthur Mensch, Katie Millican, David Moore, Michael Needham, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23732, 2022. 1, 2

[4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 1

[5] 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. 1, 2, 4, 6, 7

[6] Moo Jin Kim, Chelsea Finn, and Percy Liang. Fine-tuning vision-language-action models: Optimizing speed and success. *arXiv preprint arXiv:2502.19645*, 2025. 1

[7] Rui Liu, Wenguan Wang, and Yi Yang. Volumetric environment representation for vision-language navigation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16317–16328, 2024. 1

[8] Liuyi Wang, Zongtao He, Ronghao Dang, Mengjiao Shen, Chengju Liu, and Qijun Chen. Vision-and-language navigation via causal learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13139–13150, 2024. 1

[9] Kevin Black, Noah Brown, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo 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. 1, 2

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

[11] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 44(10-11):1684–1704, 2025. 1, 4, 6

[12] Ruijie Zheng, Yongyuan Liang, Shuaiyi Huang, Jianfeng Gao, Hal Daumé III, Andrey Kolobov, Furong Huang, and Jianwei Yang. Tracevla: Visual trace prompting enhances spatial-temporal awareness for generalist robotic policies. *arXiv preprint arXiv:2412.10345*, 2024. 2, 6

[13] 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. 6, 7

[14] Haoyu Zhen, Xiaowen Qiu, Peihao Chen, Jincheng Yang, Xin Yan, Yilun Du, Yining Hong, and Chuang Gan. 3d-vla: A 3d vision-language-action generative world model. *arXiv preprint arXiv:2403.09631*, 2024. 2

[15] Dantong Niu, Yuvan Sharma, Haoru Xue, Giscard Biamby, Junyi Zhang, Ziteng Ji, Trevor Darrell, and Roei Herzig. Pre-training auto-regressive robotic models with 4d representations. *arXiv preprint arXiv:2502.13142*, 2025. 2, 3

[16] Jiahui Zhang, Yurui Chen, Yueming Xu, Ze Huang, Yanpeng Zhou, Yu-Jie Yuan, Xinyue Cai, Guowei Huang, Xingyue Quan, Hang Xu, et al. 4d-vla: Spatiotemporal vision-language-action pretraining with cross-scene calibration. *arXiv preprint arXiv:2506.22242*, 2025. 2, 3, 6

[17] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36:44776–44791, 2023. 2, 5

[18] Haotian Liu, Pengchuan Zhang, Jianwei Yang, and Lei Zhang. Improved visual instruction tuning for image-text alignment. *arXiv preprint*, abs/2310.00067, 2023. 2

[19] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 19730–19742. PMLR, 2023.

[20] Hanyu Zhou and Gim Hee Lee. Llava-4d: Embedding spatiotemporal prompt into lmms for 4d scene understanding. *arXiv preprint arXiv:2505.12253*, 2025. 2, 6

[21] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. 2, 6

[22] Zewei Zhou, Tianhui Cai, Seth Z Zhao, Yun Zhang, Zhiyu Huang, Bolei Zhou, and Jiaqi Ma. Autovla: A vision-language-action model for end-to-end autonomous driving with adaptive reasoning and reinforcement fine-tuning. *arXiv preprint arXiv:2506.13757*, 2025.

[23] Haozhan Li, Yuxin Zuo, Jiale Yu, Yuhao Zhang, Zhaohui Yang, Kaiyan Zhang, Xuekai Zhu, Yuchen Zhang, Tianxing Chen, Ganqu Cui, et al. Simplevla-rl: Scaling vla training via reinforcement learning. *arXiv preprint arXiv:2509.09674*, 2025.

[24] 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.

[25] Junjie Wen, Yichen Zhu, Jinming Li, Minjie Zhu, Zhibin Tang, Kun Wu, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin Shen, et al. Tinyvla: Towards fast, data-efficient vision-language-action models for robotic manipulation. *IEEE Robotics and Automation Letters*, 2025. 2

[26] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 2, 6

[27] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 2

[28] Xiaoqi Li, Liang Heng, Jiaming Liu, Yan Shen, Chenyang Gu, Zhuoyang Liu, Hao Chen, Nuowei Han, Renrui Zhang, Hao Tang, et al. 3ds-vla: A 3d spatial-aware vision language action model for robust multi-task manipulation. In *9th Annual Conference on Robot Learning*, 2025. 2

[29] Helong Huang, Min Cen, Kai Tan, Xingyue Quan, Guowei Huang, and Hong Zhang. Graphcot-vla: A 3d spatial-aware reasoning vision-language-action model for robotic manipulation with ambiguous instructions. *arXiv preprint arXiv:2508.07650*, 2025. 2

[30] Chengkai Hou, Yanjie Ze, Yankai Fu, Zeyu Gao, Songbo Hu, Yue Yu, Shanghang Zhang, and Huazhe Xu. 4d visual pre-training for robot learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8451–8461, 2025. 3

[31] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 5294–5306, 2025. 4, 5, 6

[32] Yuliang Zou, Zelun Luo, and Jia-Bin Huang. Df-net: Unsupervised joint learning of depth and flow using cross-task consistency. In *Eur. Conf. Comput. Vis.*, pages 1–18. Springer, 2018. 4

[33] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1851–1858, 2017. 4

[34] Yang Li, Si Si, Gang Li, Cho-Jui Hsieh, and Samy Bengio. Learnable fourier features for multi-dimensional spatial positional encoding. *Advances in Neural Information Processing Systems*, 34:15816–15829, 2021. 4

[35] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. 4, 5, 6

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Adv. Neural Inform. Process. Syst.*, 30, 2017. 4

[37] Chun-Fu Richard Chen, Quanfu Fan, and Rameswar Panda. Crossvit: Cross-attention multi-scale vision transformer for image classification. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 357–366, 2021. 4

[38] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 12888–12900. PMLR, 2022. 5

[39] Hanyu Zhou and Gim Hee Lee. Uni4d-llm: A unified spatiotemporal-aware vlm for 4d understanding and generation. *arXiv preprint arXiv:2509.23828*, 2025.

[40] Hanyu Zhou and Gim Hee Lee. Llafea: Frame-event complementary fusion for fine-grained spatiotemporal understanding in lmms. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22294–22304, 2025. 5

[41] Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics: Methodology and distribution*, pages 492–518. Springer, 1992. 5

[42] 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. 6

[43] Abby O'Neill, Abdul Rehman, Abhiram Maddukuri, Abhishek Gupta, Abhishek Padalkar, Abraham Lee, Acorn Pooley, Agrim Gupta, Ajay Mandlekar, Ajinkya Jain, et al. Open x-embodiment: Robotic learning datasets and rt-x models: Open x-embodiment collaboration 0. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6892–6903. IEEE, 2024. 5

[44] 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. *arXiv preprint arXiv:2403.12945*, 2024.

[45] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023. 5

[46] Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X Chang. Scan2cap: Context-aware dense captioning in rgb-d scans. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3193–3203, 2021. 6

[47] Daichi Azuma, Taiki Miyanishi, Shuhei Kurita, and Motoaki Kawanabe. Scanqa: 3d question answering for spatial scene understanding. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 19129–19139, 2022. 6

[48] Dave Zhenyu Chen, Angel X Chang, and Matthias Nießner. Scanrefer: 3d object localization in rgb-d scans using natural language. In *European conference on computer vision*, pages 202–221. Springer, 2020. 6

[49] Yiming Zhang, ZeMing Gong, and Angel X Chang. Multi3drefer: Grounding text description to multiple 3d objects. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15225–15236, 2023. 6

[50] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018. 6

[51] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 6

[52] Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998. 9