

# MotionTrans: Human VR Data Enable Motion-Level Learning for Robotic Manipulation Policies

Chengbo Yuan<sup>1,2</sup>, Rui Zhou<sup>\*5</sup>, Mengzhen Liu<sup>\*3</sup>, Yingdong Hu<sup>1,2</sup>, Shengjie Wang<sup>1,2</sup>

Li Yi<sup>1,2</sup>, Chuan Wen<sup>4</sup>, Shanghang Zhang<sup>3</sup>, Yang Gao<sup>1,2†</sup>

<sup>1</sup>Institute for Interdisciplinary Information Sciences, Tsinghua University <sup>2</sup>Shanghai Qi Zhi Institute

<sup>3</sup>State Key Laboratory of Multimedia Information Processing, School of Computer Science, Peking University

<sup>4</sup>Shanghai Jiao Tong University <sup>5</sup>Wuhan University

\* Indicates equal contribution. † The corresponding author.

<https://motiontrans.github.io/>

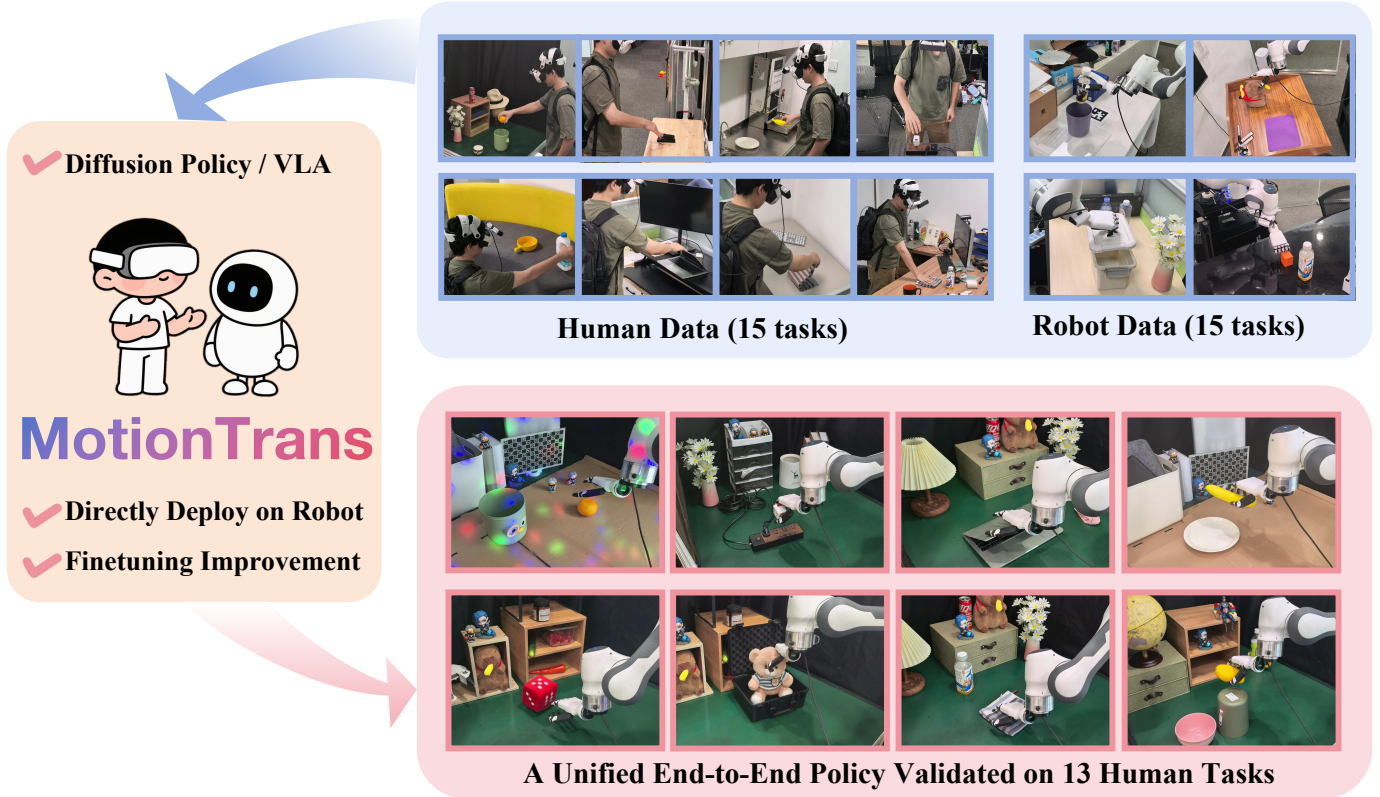


Fig. 1: We propose *MotionTrans*, a framework that enables **motion-level** learning from VR-collected human data. By cotraining on 15 human tasks and 15 robot tasks, we empower end-to-end robotic manipulation policies to directly perform tasks in human data on real robot hardware. Our framework also improves finetuning performance when a few robot demonstrations are available for these tasks.

**Abstract**—Scaling real robot data is a key bottleneck in imitation learning, leading to the use of auxiliary data for policy training. While other aspects of robotic manipulation such as image or language understanding may be learned from internet-based datasets, acquiring motion knowledge remains challenging. Human data, with its rich diversity of manipulation behaviors, offers a valuable resource for this purpose. While previous works show that using human data can bring benefits, such as improving robustness and training efficiency, it remains unclear whether it can realize its greatest advantage: *enabling robot policies to directly learn new motions for task completion*. In this paper, we systematically explore this potential through multi-task human-robot cotraining. We introduce *MotionTrans*, a framework that includes a data collection system, a human data transformation

pipeline, and a weighted cotraining strategy. By cotraining 30 human-robot tasks simultaneously, we directly transfer motions of 13 tasks from human data to deployable end-to-end robot policies. Notably, 9 tasks achieve non-trivial success rates in zero-shot manner. *MotionTrans* also significantly enhances pretraining-finetuning performance (+40% success rate). Through ablation study, we also identify key factors for successful motion learning: cotraining with robot data and broad task-related motion coverage. These findings unlock the potential of motion-level learning from human data, offering insights into its effective use for training robotic manipulation policies. All data, code, and model weights are open-sourced <https://motiontrans.github.io/>.

## I. INTRODUCTION

Learning robotic manipulation policies from teleoperated demonstrations has progressed rapidly in recent years [11, 12, 7]. However, collecting large-scale robot datasets remains costly and labor-intensive [26, 49], creating a significant bottleneck for further improvement of manipulation abilities. To address data scarcity, researchers have turned to auxiliary sources, such as images or language [23, 77] to help policy training. While internet data provides abundant vision-language knowledge to aid policy learning [21], acquiring motion knowledge remains a significant challenge.

Human data [54, 18] represents a particularly promising source to solve this: it is abundant, easy to collect, and rich in diverse manipulation behaviors [18]. Previous works have leveraged human demonstrations to extract task-aware representations, such as affordances [3] or keypoint flows [74], to support motion transfer. However, the introduction of intermediate representation hinders integration with mainstream end-to-end policies. More recently, with advances in wearable sensing, researchers begin to explore the use of human motion data (with hand poses recorded from VR device) directly for robot policy cotraining or pretraining [25, 54, 70, 44, 6]. These approaches have shown benefits for visual grounding [44], robustness [70] and training efficiency [6]. However, it is still uncertain whether it can fully realize its greatest advantage: *allowing robot policies to directly acquire new motions for task completion.*

In this paper, we investigate this question by introducing *MotionTrans*, a framework designed to **directly learn 13 robot-executable motions from human data for a unified, end-to-end robot policy**. This is achieved through multi-task human-robot cotraining. We develop a VR-based teleoperation system and data collection pipeline to construct the *MotionTrans Dataset*, which includes 3,213 demonstrations across 15 human tasks and 15 robot tasks from more than 10 scenes. We further propose a transformation procedure that maps human demonstrations into the robot’s observation–action space, making them compatible with mainstream end-to-end policies such as Diffusion Policy [12] or the Vision-Language-Action model ( $\pi_0$ -VLA) [7]. Finally, we adopt a weighted cotraining strategy that jointly optimizes over both human and robot tasks. We name the entire framework *MotionTrans* because it enables motion transfer from human data to deployable robot policies.

We first evaluate the zero-shot performance on all human tasks. This means that we directly deploy policies to robot without collecting any robot data for these tasks. Results show that Diffusion Policy [12] and  $\pi_0$ -VLA model [7] achieve non-trivial success rates for 9 tasks in total. Even in unsuccessful cases, they exhibit meaningful motion for task completion, such as reaching target objects. We also find that, when few robot demonstrations of these human tasks are available for finetuning, pretraining on the *MotionTrans Dataset* leads to an average 40% boost in success rate on these tasks. Further analysis indicates that the effectiveness of motion transfer depends on the presence of both robot demonstrations and sufficient task-related motion

coverage during training. Together, these findings highlight the possibility for motion-level learning from human data, and provide a clear framework and principles for achieving this. Our contributions can be summarized as:

- *MotionTrans*, a framework for end-to-end human-to-robot motion transfer, including data collection system, a pipeline to transform human data into robot format, and a weighted human-robot cotraining strategy.
- *MotionTrans Dataset*, containing 3,213 demonstrations for 15 human tasks and 15 robot tasks across 10+ scenes.
- ***MotionTrans* enables explicit human motions transfer for end-to-end robot policies, even for zero-shot settings** (directly learn 13 tasks from human data).
- Key factors for successful motion transfer: robot data cotraining and sufficient task-related motion coverage.

## II. RELATED WORK

### A. Imitation Learning for Robot Manipulation

Imitation learning [32, 4, 58, 34] has made significant progress in recent years. By learning motion from training data [11, 10], imitation policies can effectively perform a wide range of manipulation tasks [12, 76], including challenging multi-task settings [80, 38, 7, 6, 35]. In this paper, we focus on two widely-used architectures for imitation learning: Diffusion Policy [12] and the  $\pi_0$  Vision-Language-Action Model ( $\pi_0$ -VLA) [7]. However, the scalability of training data remains a major challenge, due to the high cost of collecting real-robot data [49, 26, 67]. This has led to the use of auxiliary data [23, 77, 36] for policy training. Despite ability such as image or language understanding in robotic manipulation could improve from internet-based pretraining [21, 33], acquiring motion knowledge remains difficult. Human data [18, 41, 73, 14], with its abundant and diverse manipulation behaviors, provides a valuable supplement for this.

### B. Task-Aware Representation Learning from Human

Early works have leveraged task-aware representations for human-to-robot knowledge transfer. Self-supervised learning has been used for implicit task-aware representations [46, 24, 45, 71, 9] learning, while representations like affordances [3, 28, 57], object poses [19], videos [5, 51], and motion flows [74, 66, 68, 55] support motion-aware representation learning. Some approaches use wrist trajectories as prompts for one-shot human-to-robot skill transfer [27, 78, 79, 60, 50]. EgoZero [39] predicts wrist poses from smart glasses, but relies on keypoint-based representations [64] for policy observations. The use of intermediate representations in these methods limits their integration with mainstream end-to-end visuomotor policy learning [12, 7], restricting their future applicability.

### C. End-to-End Policy Learning with Posed Human Data

Human motion data can be captured through hand-held SLAM-based device [13, 69], but often limited to only wrist camera sensing [61]. Recent advancements in wearable sensing [15, 10, 54] now allow easy collection of posed human data (with hand keypoints, wrist poses information etc.) through

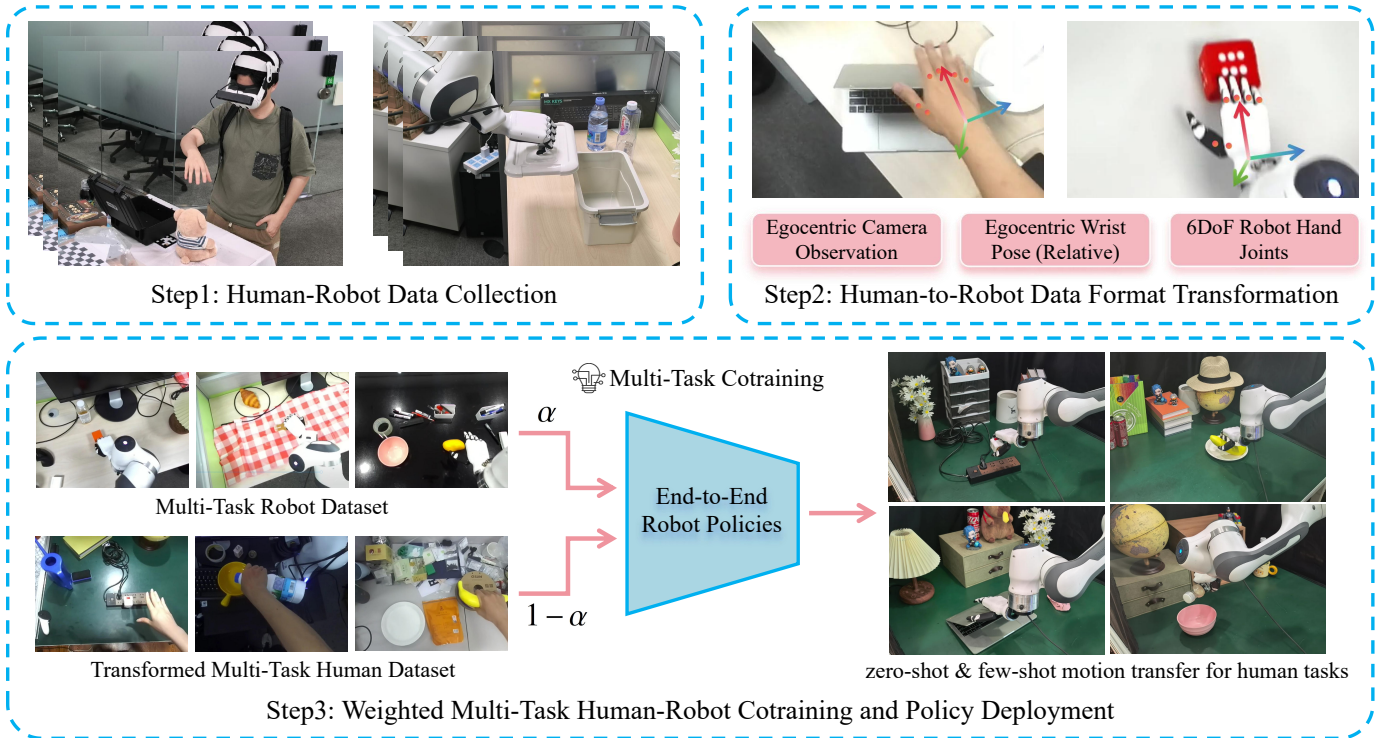


Fig. 2: Illustration of our proposed *MotionTrans* framework, which consists of a human-robot data collection system, a pipeline for transforming human data into robot format, and a weighted human-robot multi-task cotraining strategy. After training, we enable the direct deployment of the trained policies to perform tasks in human datasets on real robots.

VR devices [18]. This data provide action label for prediction, supporting end-to-end policy learning [30]. Some studies cotrain human and robot data [25, 54, 47, 29, 61, 40], while others first pretrain with human data and then finetune with robot demonstrations [70, 44, 6]. These works have shown policy improvements in visual grounding [44], robustness [54, 70], and training efficiency [6, 25]. However, whether it can achieve direct transfer of motions from human to robot remains unclear [39]. To the best of our knowledge, our paper is the first to systematically verify motion-level end-to-end learning from human data.

### III. MOTIONTRANS

In this section, we present our proposed *MotionTrans* framework (Figure 2). The core idea is to first transform human data to robot data format, and then jointly learn from human and robot data within the robot observation-action space. By training policies in robot space, we can directly deploy policies to perform tasks from human data on real-world robots, i.e., enabling explicit human-to-robot motion transfer. We first introduce the motion transfer problem and define the observation-action space of the policy (Section III-A). To facilitate human-robot data cotraining, we develop data collection systems for both human and robot data (Section III-B). We then propose a pipeline to convert human data into robot format (Section III-C). This ensures compatibility with mainstream robot policies, enabling subsequent end-to-end cotraining. Finally, we choose the architecture of robot policies and apply human-robot multi-task cotraining (Section III-D).

#### A. Problem Definition

Our goal is to enable explicit human-to-robot motion transfer. Considering the embodiment gap between human and robot [25], we explore this problem within a multi-task human-robot cotraining framework, where robot data for certain tasks are available to help motions in human data adapt to the robot. Specifically, we aim to train a policy  $P_{\text{policy}}$  on  $D = D_{\text{robot}} \cup D_{\text{human}}$ , where  $D_{\text{robot}} = \{D_{\text{robot}}^i \mid i = 1, \dots, N_{\text{robot}}\}$  is the robot dataset, and  $D_{\text{human}} = \{D_{\text{human}}^i \mid i = 1, \dots, N_{\text{human}}\}$  is the human dataset. Each  $D^i$  represents a sub-dataset corresponding to a specific task, and the task sets of the human and robot data are **non-overlapping**. After training, we deploy  $P_{\text{policy}}$  on a real-world robot and evaluate its performance on **tasks from  $D_{\text{human}}$**  to assess the effectiveness of motion transfer. This is defined as the **zero-shot** setting, since the evaluation tasks contain no corresponding robot data for training. We also evaluate the performance of **few-shot finetuning** setting, where a small number of robot demonstrations for the tasks from  $D_{\text{human}}$  are available to further finetune  $P_{\text{policy}}$ .

We define the input and output of our policies within the robot observation-action space  $S = (I_t, P_t, A_t)$ . At each timestamp  $t$ , the policy receives an egocentric RGB image  $I_t \in \mathbb{R}^{H \times W \times 3}$  and proprioceptive states  $P_t \in \mathbb{R}^{T_P \times D}$ , where  $T_P$  is the history length and  $D$  is the state dimension. For simplicity, this work focuses on single-arm tasks (Figure 4), thus  $D$  corresponds to the concatenation of one robot wrist pose and one robot hand joint state (Figure 3(c)). The policy outputs an action



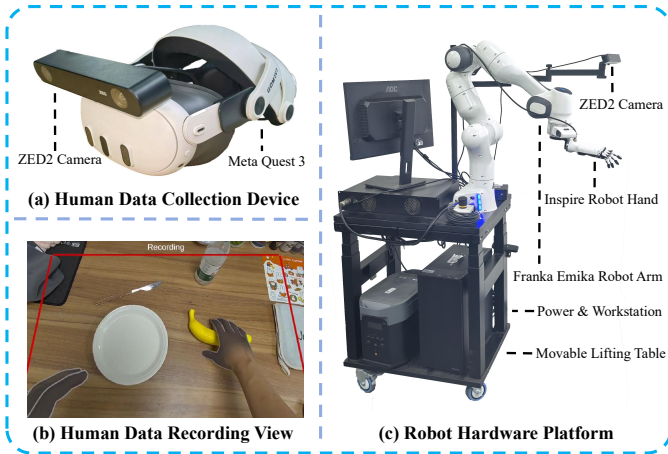


Fig. 3: Illustration of our hardware system, which includes a human VR-based data collection device and a single-arm robot platform. A screenshot of the VR device during human data collection is also provided.

chunk prediction  $A_t \in \mathbb{R}^{T_A \times D}$  [12], where  $T_A$  denotes the action prediction horizon. Next, we describe the details of our human-robot data collection and processing system.

### B. Human-Robot Data Collection System

For human-robot cotraining, we need to collect both robot and human data [25]. For human data collection, we leverage a portable commercial VR device, which allows data to be collected anytime and anywhere. This provides great efficiency in gathering diverse motions and a wider range of tasks [18]. For robot data collection, we use teleoperation to record demonstrations. The top-left side of Figure 2 illustrates the two types of data collection systems.

**Human Data Collection with Portable VR Device.** We extend ARCap [10] to build our human data collection system (Figure 3(a)), incorporating a portable VR headset for recording hand keypoint positions  $K_t$ , wrist poses  $W_t$  and camera poses, and an RGB camera for the image stream  $I_t$ .

For hand pose recording, our goal is to capture the positions of hand keypoints  $K_t$  and human wrist poses  $W_t$  in the coordinate frame of the RGB camera ( $I_t$ ). However, these information is recorded by VR device, placing it in the VR coordinate space. Therefore, we use a self-designed calibration method to transform all hand information from the VR coordinate space to the RGB camera, detailed in Appendix VI-C. For data collection, collectors are instructed to minimize head motion to approximate the static camera setting of real robot hardware, although slight movements are tolerated [54]. To ensure data quality, we provide real-time feedback in the user’s VR view to guide collectors during data acquisition (Figure 3(b)). The feedback includes the RGB camera’s capture range and the positioning of the hands:

- The range of images captured by the RGB camera is used to guide users to ensure their hands are always visible to the RGB camera [10].

- The hand positioning tells collector whether the hand poses recorded by VR is strictly aligned with their hands in real time, thus provide information about the recording delay and accuracy.

We also provide gesture interface to allow collector to abandon current recorded data anytime, if they think the data quality is not good enough considering the principles and feedback mentioned above.

**Robot Data Collection with Teleoperation.** Since our goal is to achieve direct human-to-robot motion transfer, the robot hardware platform need to match the functionality of the human arm and hand. To this end, we choose the combination of a single robot arm and a dexterous robot hand as our hardware platform (Figure 3(c)). We develop our teleoperation system on Open-Television [11], which captures human wrist and hand poses in real time via a VR device and drives the robot to replicate these motions. Based on the collection system above, we collect our *MotionTrans* human-robot datasets (Section IV-A and Figure 4) for multi-task cotraining.

### C. Human Data Transformation to Robot Format

As shown in the previous section, the raw human data collected from the VR device differs in format from robot data, which prevents it from being directly used for cotraining with robot policies [70, 44]. To address this, we propose directly transforming human data into the robot’s observation-action space [11, 43]. After transformation, the human data can serve as a form of “supplementary robot data” for training any mainstream end-to-end **robot** policy.

**Transforming Observation-Action Space.** The observation-action space of the robot includes three components: image observation  $I_t$ , proprioceptive state  $P_t$ , and action  $A_t$  (refer to Section III-A). Both  $P_t$  and  $A_t$  are generated by stacking wrist poses  $W_t$  and hand joint states  $H_t$ . Next, we describe the design for these components:

- **Image observation  $I_t$ :** We use **egocentric** view for both human and robot data, as shown in Figure 4. The use of the similar image view makes the spatial relationships of objects in the scenes similar for accomplishing similar tasks, thus enabling similar motions to achieve those tasks.
- **Wrist poses  $W_t$ :** We use the **egocentric** camera coordinate system (camera captures  $I_t$ ) for both human and robot data. This allows for the measurement of wrist poses in a unified coordinate system, ensuring that the spatial definitions of human and robot data are consistent.
- **Hand joints state  $H_t$ :** we employ the dex-retargeting library [53], an optimization-based inverse kinematics solver, to map human hand keypoints  $K_t$  to robot hand joint state  $H_t$ .

The design above converts human data into the same format as robot data, enabling us to directly replay human data on real-world robots. The replayable property of transformed human data proves how aligned our processed data is with robot data. By replaying human trajectories on a real-robot platform, we derive the following key observations: (O1) the speed of human



manipulation is much faster than that of the robot, which affects safety and motion planning stability; (O2) there is a discrepancy between the distributions of human hand positions and the robot’s comfortable workspace (all defined in egocentric camera coordinate space). To alleviate these problem, we:

- (O1) We **slow down human data by a factor of 2.25** via poses and hand joints state interpolation. More advanced techniques, such as the adaptive speed downsampling strategy [58], are left for future exploration.
- (O2.1) We utilize **action-chunk-based relative poses** [12, 76] as wrist action representation to reduce distribution mismatches between human and robot data. For instance, even if the robot’s and human’s hand positions differ in world space, their relative poses remain the same if they move forward at the same speed.
- (O2.2) We encourage collectors to **change viewpoints between trajectory recordings**. This enhances the diversity of positional relationships between the camera view and the targeted manipulation objects, thereby encouraging policies to adapt to a larger distribution of hand poses and, consequently, a larger workspace for the robot.

The methods and principles described above help reduce the gap between human and robot data, thereby improving the effectiveness of human-to-robot motion transfer. Prior works [29, 30, 31] have proposed rendering robots into human videos to further narrow the visual gap between the two domains. We replicate this rendering approach, as shown in Figure 9, but did not observe significant improvements over directly training on human videos (Section IV-D). Therefore, we do not employ this rendering technique by default in our framework.

#### D. Weighted Multi-Task Human-Robot Cotraining

By unifying the observation and action spaces, we enable joint training of human and robot data under a shared end-to-end robot policy. This section introduce the multi-task policy architectures we use and how we train these policies.

**End-to-End Multi-Task Policy Architectures.** We explore two popular end-to-end policy architectures: (1) **Diffusion Policy (DP)** [12]: unlike the original single-task setup, we extend DP for multi-task training. Each task is associated with a learnable embedding, serving as a unique task condition. The visual encoder is replaced with DINOv2 [48] to enhance visual perception ability [32]. (2) **Vision-Language-Action model ( $\pi_0$ -VLA)**: we adopt network structure from [7], a policy architecture integrating large-scale pretrained Vision-Language Models [59] for multimodal perception and instruction following. Since  $\pi_0$ -VLA supports language input, we directly use instructions to assign tasks. For  $\pi_0$ -VLA, we load  $\pi_0$ -droid pretrained checkpoints [52] before training.

**Unified Action Normalization.** To improve training stability, we apply Z-score normalization to both proprioceptive states and actions before training [12, 13]. Previous human-robot cotraining works [25, 61] typically adopt independent normalization for human and robot data, arguing that it reduces the action gap between the two sources. However, in our motion-level

evaluation setting, where the goal is to directly deploy human tasks on a real robot, this approach introduces a mismatch between training (human normalization) and inference (robot normalization), ultimately causing a performance drop (Section IV-D). Therefore, we adopt unified action normalization across human and robot data within our framework.

**Weighted Human-Robot Cotraining.** Our final step is to design a strategy to train multi-task policies with the processed human-robot dataset. Given the potential imbalance between human and robot data [61, 54], we adopt a weighted cotraining strategy similar to [65]. The training objective over the combined dataset  $D = D_{\text{robot}} \cup D_{\text{human}}$  is defined as:  $\mathcal{L}_D = \alpha \mathcal{L}_{D_{\text{robot}}} + (1 - \alpha) \mathcal{L}_{D_{\text{human}}}$ , where  $\mathcal{L}$  denotes the loss function of imitation learning [12, 7]. In this paper, we set:

$$\alpha = \frac{|D_{\text{human}}|}{|D_{\text{human}}| + |D_{\text{robot}}|}$$

where  $|D_{\text{robot}}|$  and  $|D_{\text{human}}|$  representing the dataset sizes. This weight ensures that the sum of the weights for human and robot data is equal, leading to the balance of these two data sources. We also try domain adaptation training techniques like domain confusion [63, 62] to promote knowledge transfer from human domain to robot domain in our earlier exploration, but do not find it beneficial for motion transfer and it always leads to training instability. Thus, we choose the simplest weighted cotraining strategy in our framework. More details could be found in Appendix VI-E.

## IV. EXPERIMENTS

In this section, we conduct experiments to verify the effectiveness of *MotionTrans* for human-to-robot motion transfer. We first introduce our detailed experiment setup in Section IV-A, including human-robot hardware platform, training datasets and evaluation tasks and metric. We then conduct experiments for both zero-shot (Section IV-B) and few-shot (Section IV-C) settings, as demonstrated in Section III-A. Additionally, an ablation study on the key designs of *MotionTrans* is performed (Section IV-D). We also carry out experiments to explain the mechanism of human-to-robot motion transfer in Section IV-E. The evaluation results of all robot tasks are shown in Appendix VI-G. Finally, we verify the robustness of our results concerning visual backgrounds in Appendix VI-H.

### A. Experiment Setup

**Hardware Platform.** For the robot hardware (Figure 3(c)), we use a Franka Emika robot arm [17] in combination with a 6DoF Inspired Dexterous (Right) Hand [11]. This combination mimics the functionality of a human right hand and arm. The robot is mounted on a movable lift table to facilitate data collection in various locations. A ZED2 camera is fixed to the table in an egocentric view to provide an image observation stream. The recorded images are first cropped to 640×480 resolution and then resize to 224×224. The VR device used for teleoperation is the Meta Quest 3 [11]. Calibration between the robot base and the robot perception camera is achieved through the DROID platform codebase [26].



in the human dataset but not in the robot dataset, such as unplugging, closing, lifting, etc. Overall, the dataset covers a wide range of motions and skills, including pick-and-place, pouring, wiping, pushing, pressing, opening, etc. This wide coverage is proven crucial for successful motion transfer, as demonstrated in subsequent ablation studies (Section IV-E). For simplicity, we name pick-place task with “pick object-place target” format, and name other task with “verb noun” format in the main paper. For tasks with multiple steps, we name it as “step1+step2” format.

To enhance the visual robustness of the policies [75] (Section VI-H), such as robustness to different backgrounds and lighting conditions, we collect these data across various scenes [32]. Each human task is collected in at least 4 different scenes. For robot tasks, about half of the data is collected in the “*green table scenes*” (the scenes for the examples of the “Bread-Pad” and “Unplug Charger” task in Figure 4), with random disturbance objects placed on the table for approximately 80% of the data. *This scene is also designated as the default scene for our evaluation.* The other half of the robot tasks is collected in at least 4 scenes. To enrich language instructions for VLA training, we leverage GPT-4o [20] to paraphrase and expand task descriptions in the dataset.

**Evaluation Tasks and Metrics.** Since our goal is to understand the effectiveness of human-to-robot motion transfer, we focus on **evaluating the performance of robot policies on the human tasks**. Among all 15 tasks in human dataset, there are two tasks (“Fold Towel” and “Pour Milk Bottle”) not been able to deploy to robot due to the hardware design limitation of robot hand (cannot be accomplished even if we use teleoperation). Therefore, we focus on discussing other 13 tasks in this research. The list of all 13 evaluation tasks could be found in Figure 6.

We use the *Success Rate (SR)* to evaluate the policy performance in accomplishing specific tasks. However, this metric alone is insufficient to reflect the effectiveness of motion transfer, as it ignores meaningful motion during task execution. For instance, a policy that demonstrates reaching for the target object should be rated higher than one that does not move at all. To address this limitation, we define a *Motion Progress Score (Score)* to quantify the quality of policy motion for task completion. Detailed scoring rubrics for all tasks are provided in Appendix VI-B. For clarity, we normalize the Score to a [0,1] range in the main paper. For each task, we conduct 10 rollouts and calculate the average results for both metrics. We change **the object arrangement** for each rollout to cover a wide range of configurations of the task across the 10 rollouts.

### B. Zero-shot Experiment

The goal of the zero-shot experiment is to verify the effectiveness of direct human-to-robot motion transfer. We train policies using our *MotionTrans Dataset*. Subsequently, we directly deploy policies to real robot hardware and evaluate the performance of tasks in human data. We refer to this as zero-shot setting because the policies learn motions from humans

without any robot data collected for these human tasks. We seek to answer the following questions:

- (Q1.1) Can the policy directly learn to accomplish tasks in human data by human-robot cotraining, even without collecting any robot data for these tasks?
- (Q1.2) For tasks that cannot be accomplished, can the policy learn meaningful motion for task completion?
- (Q1.3) Is cotraining with robot data the key factor for achieving explicit motion transfer?
- (Q1.4) What is the difference in motion transfer effectiveness between different policy architectures?

**Experiment Details.** We train two end-to-end policies, Diffusion Policy (DP) and  $\pi_0$ -VLA (as mentioned in Section III-D). For DP, we train it for 300 epochs with a learning rate of  $5 \times 10^{-4}$  and 1024 batch size. For  $\pi_0$ -VLA, we train it for 160,000 steps with a learning rate of  $2.5 \times 10^{-5}$  and 192 batch size. Both models are trained with the AdamW optimizer [42]. The training takes approximately 1.5 days and 2.5 days respectively. In this paper, we focus on enabling human-to-robot transfer for mainstream end-to-end policies. Therefore, we do not compare against zero-shot intermediate representation-based methods such as Vid2Robot [22], General-Flow [74], EgoZero [39], ZeroMimic [57] etc., which are not compatible with such policies. Instead, our analysis centers on differences among end-to-end policy architectures (DP vs.  $\pi_0$ -VLA).

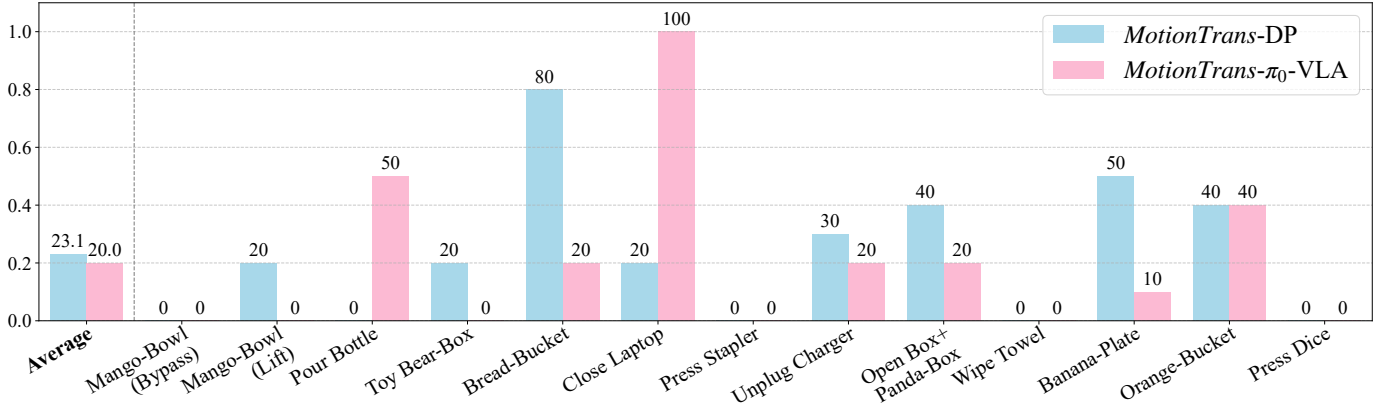
**(Q1.1) *MotionTrans* enables policies to achieve non-trivial success rate across 9 tasks in the human dataset.** The results of the zero-shot experiment are shown in Figure 6. As shown in the results, 9 tasks achieve a non-trivial success rate. The average success rate on all 13 tasks is approximately 20%. The visualization of two examples could be found in the Figure 7(a) (“Orange-Bucket” and “Unplug Charger”). Among these tasks, pick-and-place tasks account for the vast majority. This can be attributed to (1) the simplicity of pick-and-place motion, (2) the similarity of motions between different pick-and-place tasks, and (3) the large number of such tasks in our dataset. Notably, for the cases where even if both the pick objects and place targets are not seen in robot tasks (e.g., the “Orange-Bucket” task, visualized on the left side of Figure 7(a)), this type of task-level transfer is still possible.

Other accomplished tasks include motions such as pouring, unplugging, lifting, opening and closing (pressing). While some tasks (e.g., “Unplug Charger”, 20%) attain only limited success rate, the model consistently exhibit meaningful motion tendencies in unsuccessful rollouts, as will be discussed below. Overall, reaching the target emerged as the most reliable step across tasks, whereas precision-demanding actions such as grasping and infrequent motions in the dataset, such as unplugging, achieved limited success rates.

**(Q1.2) For unsuccessful tasks, *MotionTrans* enables policies to learn meaningful motions toward task completion.** Figure 6 shows that both DP and  $\pi_0$ -VLA achieve positive Motion Progress Scores across all tasks, with an overall average of about 0.5. This indicates that the policies are able to complete



Zero Shot Success Rate Results (SR, %)



Zero Shot Motion Progress Score Results (Score)

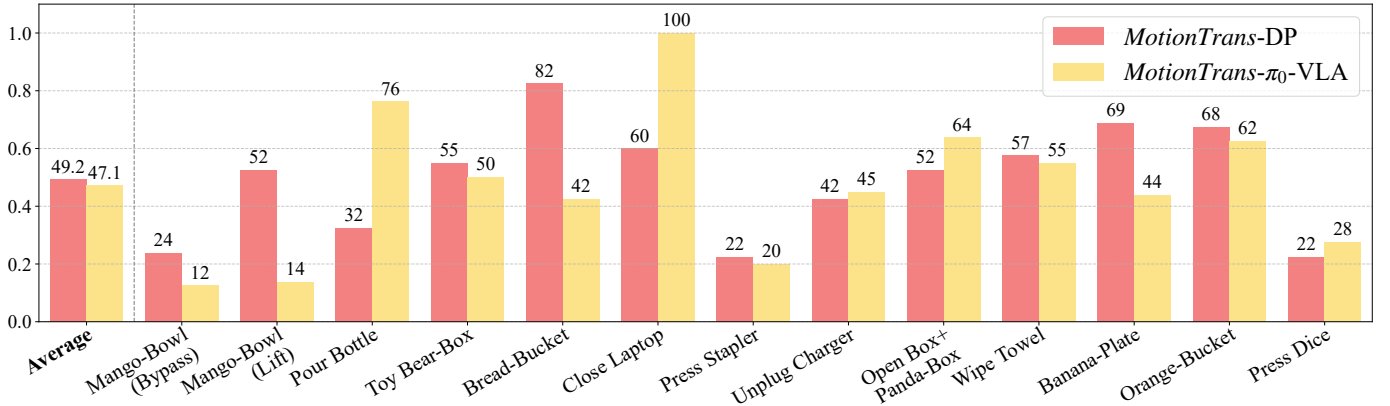


Fig. 6: Results of *MotionTrans* in the zero-shot experiment setting. The results show that both Diffusion Policy (DP)[?] and  $\pi_0$ -VLA[7] achieve successful human-to-robot motion transfer. Even without any robot data for these human tasks, 9 tasks attain a non-zero success rate. For the remaining tasks, *MotionTrans* still generates meaningful motion for task accomplishment, as indicated by a non-trivial Motion Progress Score.

certain sub-processes for all evaluation tasks. For instance, in the “Wipe Towel” task, both DP and  $\pi_0$ -VLA learn the motion of “push towel forward” to some extent (left side of Figure 7(b)). Moreover, we observe that human data enables the policy to identify spatial locations for almost all evaluated human tasks, which is represented as reaching the target manipulated objects (may only appear in human data) to some extent. An example of this is the “Press Stapler” task in Figure 7(b): although the stapler is not seen in the robot data, the policy still performs approaching behavior.

**(Q1.3) Cotraining with robot data is the key factor for successful motion transfer.** We find that when robot data is not included for cotraining, the success rate across all tasks is 0% for zero-shot setting. Generally, the policy trained solely on human data exhibits random motion when deployed on the robot. This demonstrates that cotraining with robot data is essential for explicit human-to-robot motion transfer, which could bridge the gap between humans and robots, allowing human motions to adapt to robot embodiment. A detailed analysis of the mechanism by which robot data support motion transfer can be found in Section IV-E.

	<i>MotionTrans</i> -DP	<i>MotionTrans</i> - $\pi_0$ -VLA
Toy Bear-Box	40	0
Bread-Bucket	100	20
Banana-Plate	50	10
Orange-Bucket	70	50
<i>Average</i>	<b>65</b>	<b>20</b>

TABLE II: The Success Rate of DP and  $\pi_0$ -VLA on all evaluation pick-and-place tasks for zero-shot setting. Generally, DP outperforms  $\pi_0$ -VLA during the grasping stage.

**(Q1.4) DP and  $\pi_0$ -VLA each have their own advantages (manipulation precision and task adherence).** As shown in Figure 6, no single model excels across all tasks. On average, the performance of the two models is nearly identical. However, we observe that different models demonstrate their strengths on different tasks. Generally, DP performs better than  $\pi_0$ -VLA in precise manipulation stage, such as grasping, and exhibits stronger spatial location capabilities. An example of evidence for this is that, for all pick-and-place tasks, the average grasping success rate of  $\pi_0$ -VLA is 20%, while DP achieves 65% (Table II). In contrast,  $\pi_0$ -VLA shows stronger instruction

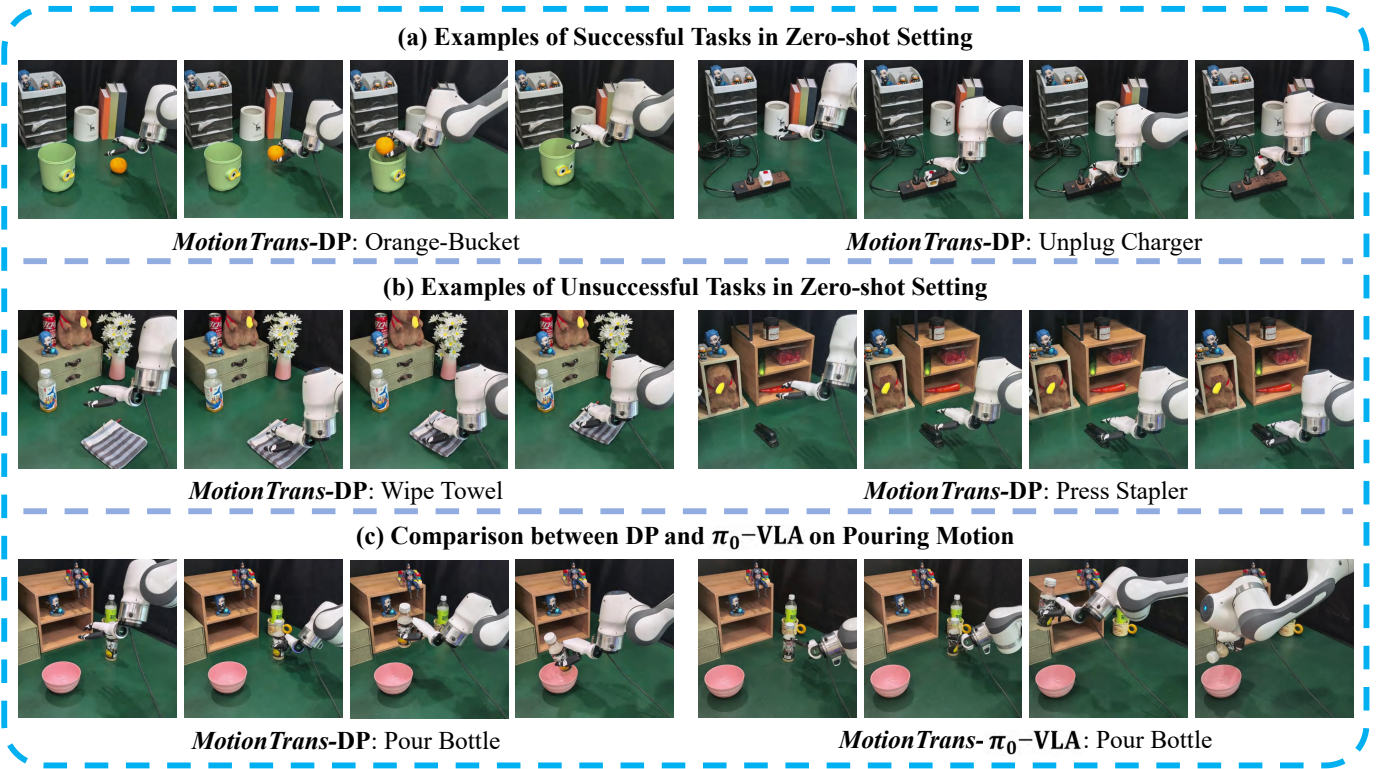


Fig. 7: The visualizations for **zero-shot** human-to-robot motion transfer from our *MotionTrans* framework. All tasks shown here do not involve any robot data collection and are learned from human data. These results demonstrate that the *MotionTrans* enables explicit human-to-robot motion transfer for task completion through human-robot cotraining.

following for motion generation in more cases. For example, in the “Pour Bottle” task, we observed limited wrist rotation with DP, while  $\pi_0$ -VLA successfully performs the complete pouring action (Figure 7(c)). We hypothesize that this difference arises from a balance between visual perception and task semantic following. The model that focuses more on visual perception (DP) tends to achieve greater manipulation precision, whereas the model that emphasizes task semantics and instruction following ( $\pi_0$ -VLA) can adhere to task requirements more stringently.

### C. Few-shot Experiment

In this section, we investigate whether motion transfer from human-robot cotraining can also enhance performance in a few-shot finetuning setting, where a limited number of robot demonstrations of human tasks are available for policy finetuning. We aim to answer the following questions:

- (Q2.1) Will pretraining on *MotionTrans Dataset* help improve policy finetuning performance?
- (Q2.2) What is the contribution of human data versus robot data for policy pretraining?
- (Q2.3) How does pretraining improvement vary with increasing finetuning data?

**Experiment Details.** Considering DP and  $\pi_0$ -VLA exhibit similar average performance in zero-shot experiments, we focus on DP architecture for computational resource efficiency in this part. We additionally collect 20 demonstrations for all

human tasks in the default “green table” evaluation scene, as mentioned in the dataset part in Section IV-A. Subsequently, we perform 5-shot and 20-shot **multi-task finetuning** [6] based on checkpoints previously trained on the *MotionTrans Dataset*. We finetune DP with a learning rate of  $1 \times 10^{-4}$  and a batch size of 256 for 200 epochs, employing the AdamW optimizer [42]. The finetuning process requires 1 hour for the 5-shot setting and 4 hours for the 20-shot setting.

We compared our method with three baselines to investigate the impact of different data components: (1) “**from-scratch**”, which means training policies without any pretraining; (2) “**robot-only**”, which entails pretraining solely on robot data from the *MotionTrans Dataset* before finetuning; and (3) “**human-only**”, which is pretrained exclusively on human data.

**(Q2.1) Pretraining on *MotionTrans Dataset* enable significant improvement for finetuning performance.** The success rate results of the few-shot experiments are presented in Figure 8. The results of Motion Progress Score can be found in Appendix VI-F. We can see that policy pretrained on *MotionTrans Dataset* gains around 40% average success rate improvement compared to “from-scratch” baseline. This is established for both 5-shot and 20-shot setting. These results prove that pretraining on human-robot data could provide useful motion prior [72, 70] for downstream finetuning.

**(Q2.2) Both robot and human data during pretraining are crucial for enhancing performance.** From Figure 8,

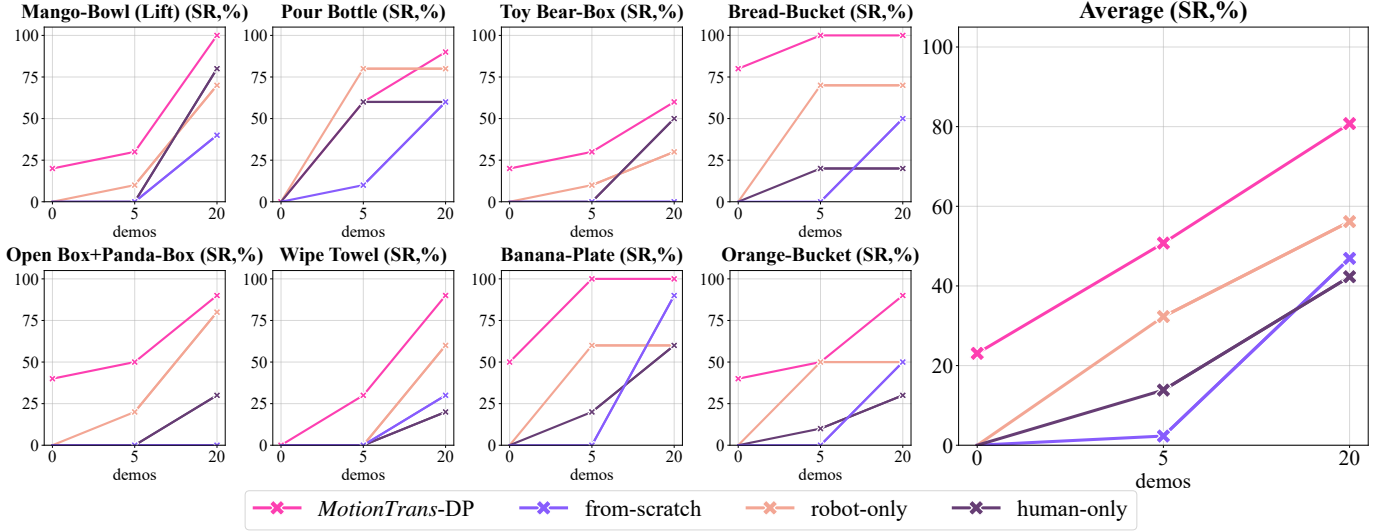


Fig. 8: Results of the success rate for few-shot finetuning experiments. For readability, only the results of 8 example tasks are presented here. The Motion Progress Score results can be found in Appendix VI-F. From these results, we can conclude that both human and robot data during pretraining are important for improving finetuning performance.

we can see that policy pretrained on both human and robot data (*MotionTrans*) shows a significant advantage compared to human-only or robot-only pretraining. Besides, robot-only pretraining outperforms human-only pretraining on average. In our setting, robot pretraining uses data from the same embodiment but different tasks, whereas human pretraining uses data from the opposite case. We therefore conclude that maintaining the same embodiment in pretraining data is more important than exactly matching tasks. This is because the distribution of robot data is generally closer to the downstream robot finetuning distribution than human data, even when the tasks differ. Moreover, motions across different tasks often share similarities, so different robot tasks can still benefit downstream finetuning performance [7].

**(Q2.3) Human-robot pretraining is more effective in low finetuning data region.** Finally, we analyze the impact of pretraining with varying amounts of finetuning data. As shown in Figure 8, the average performance of the policies improves consistently with an increase in finetuning data for all methods. However, the improvements are much larger in the 5-shot setting compared to the 20-shot setting. Moreover, when 20 finetuned demonstrations are available, the advantage of robot-only pretraining becomes minimal, and the benefit of human-only pretraining disappears. However, in the 5-shot setting, all pretraining methods show a significant advantage over the from-scratch baseline.

#### D. Design Ablation

We conduct an ablation study on the key designs of *MotionTrans*. We compare three variants of *MotionTrans* in zero-shot setting experiments, including common techniques used in prior human-to-robot imitation learning:

- **w/ Abs Pose:** We replace the action-chunk-based relative pose [13] with the original absolute egocentric pose for

	Score	SR (%)
w/ Abs Pose	0.370	10.0
w/ ED-Norm [25, 61]	0.341	8.4
w/ Visual Rendering [29, 31]	<u>0.475</u>	<b>23.1</b>
<i>MotionTrans</i> -DP	<b>0.492</b>	<b>23.1</b>

TABLE III: Ablation results of design choices for *MotionTrans*. The results are averages across all 13 evaluation human tasks.

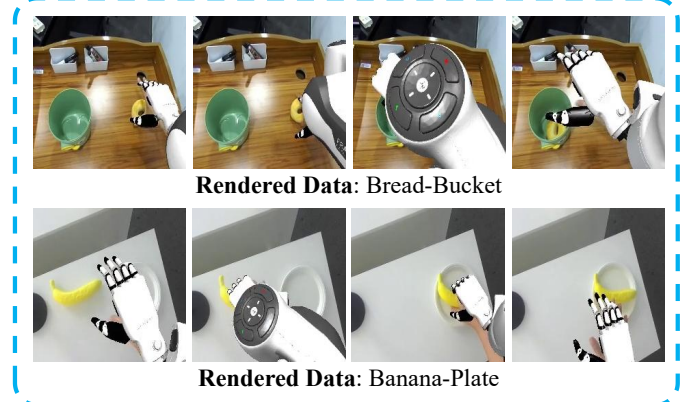


Fig. 9: The visualizations of the rendered RGB observations for the **w/ Visual Rendering** variant in design ablation (Section IV-D).

wrist action representation.

- **w/ ED-Norm** [25, 61]: We use independent action and proprioception normalization for human and robot data before policy training (Embodiment Dependent Normalization).
- **w/ Visual Rendering** [30, 29, 31]: We first replay robot data in simulation, then crop the rendered robot and paste it to the original RGB image observation. Visualizations of the rendered results are shown in Figure 9.

The policy architecture chosen for all variants is Diffusion



	Score	SR (%)
H-bucket	0.0	0
H-bucket + R-pad	0.275	0
H-bucket + R-platform	0.5	30
H-bucket + R-pad + R-platform	0.625	40
H-bucket + R-pad + R-platform + PP-set	0.75	70
all data ( <i>MotionTrans</i> )	0.825	80

TABLE IV: The results of the case study for the “Bread-Bucket” task in zero-shot setting, including outcomes from training on different subsets of *MotionTrans* Datasets. Detailed analysis could be found in Section IV-E.

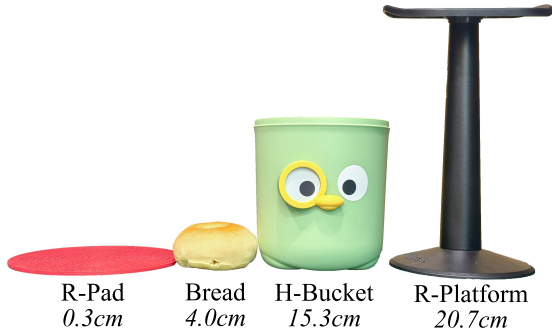


Fig. 10: The visualizations of key objects used in the “Bread-Bucket” case study are presented here. The height of each object is labeled beneath it.

Policy (DP) [12]. The results are averaged across all 13 evaluation tasks and shown in Table III. We observe that **w/ Abs Pose** and **w/ ED-Norm** dramatically decrease the performance of human-to-robot motion transfer. For **w/ Abs Pose**, the usage of absolute pose increases the distribution difference between human and robot action label, prohibiting motion transfer, as discussed in Section III-C. For **w/ ED-Norm**, performance drops because the embodiment-dependent normalization creates a discrepancy in normalization between policy training and deployment. This contrasts with the phenomenon observed in visual robustness evaluations as demonstrated by previous works [25, 61]. When directly learning new motions and skills from human data, it is preferable to keep action normalization consistent between training and inference.

For **w/ Visual Rendering**, we find that performance is nearly the same as the non-rendered version. This may be explained by the fact that, despite appearing realistic to humans, the rendered results still include cues enabling neural networks to identify the embodiment domain. From this perspective, they offer little distinction from the original human videos. One potential solution is to also conduct inpainting during policy inference [30], but may lead to additional computational overhead and policy delay. All these results demonstrate that, when considering *motion-level* transfer and evaluation, the effectiveness of certain designs may differ from their effectiveness when using human data to improve visual robustness [61] or the training efficiency [25] for in-domain robot tasks.

### E. Analysis of Motion Transfer Mechanisms

We have verified the feasibility of explicit human-to-robot motion transfer in our zero-shot experiments (Section IV-B). However, the underlying mechanisms of this transfer remain underexplored. In this section, we design experiments to investigate these mechanisms from three perspectives: (1) how actions transfer, (2) how visual perception transfers, and (3) the scaling trends of motion transfer. We then describe the experimental setup and present the corresponding conclusions.

**(Q3.1) How Do Actions Transfer?** To answer the question, we conduct a case study by down-sampling the number of tasks. We train policies on different subsets of *MotionTrans* Dataset and compare their performance, gaining insights into how varying training tasks and motions influence the actions generated during real-robot inference. We select the task “Bread-Bucket” as the evaluation task for our case study, as it already demonstrates a high success rate (80%) in zero-shot settings, indicating effective motion transfer. Since “action” is an abstract concept, we focus on a concrete dimension: **the height of object placement**. Three tasks with varying placement heights for the “bread” object are selected to create training subsets:

- **(Human Data) Bread-Bucket:** evaluation task, denoted as “H-bucket”.
- **(Robot Data) Bread-Pad:** placing bread on a thin red pad, “R-pad”.
- **(Robot Data) Bread-Platform:** placing bread on a tall black platform, “R-platform”.

Visualizations and placement heights of objects in these tasks are shown in Figure 10, and evaluation results in Table IV (rows 1–4). Trajectory visualizations are provided in Figure 11. Results show that training only on human data tends to cause ambiguity during deployment on robots, consistent with our zero-shot findings (Section IV-B). Cotraining with a single robot task appears to bias the policy toward that specific placement height. When both R-pad and R-platform are included, the policy shows evidence of interpolating across placement heights, leading to bucket-height-aware motions.

Based on the results, we **hypothesize that: the actions for human task completion during real-robot inference are generated by interpolating actions from robot data** (e.g., from R-Pad at 0.3cm and R-Platform at 20.7cm to H-Bucket at 15.3cm). This interpolation ability is learned by training on task-aware motions in human data. When the policy encounters the robot embodiment during inference, it still generates actions within the safe manipulation range as defined by the robot data. However, **task-specific elements**, such as task identifiers or task-related objects in image observation, trigger the policy to activate the interpolation process to generate task-aware motion.

**(Q3.2) How Do Visual Perception Transfers?** We next examine visual perception. To understand the impact on policy visual perception when the embodiment changes from human during training to robot during inference, we visualize attention maps of trajectories from zero-shot *MotionTrans*-DP policy using the DINOv2 encoder [48] and Grad-CAM [56].

### The Performance of “Bread-Bucket” with Different Training Set

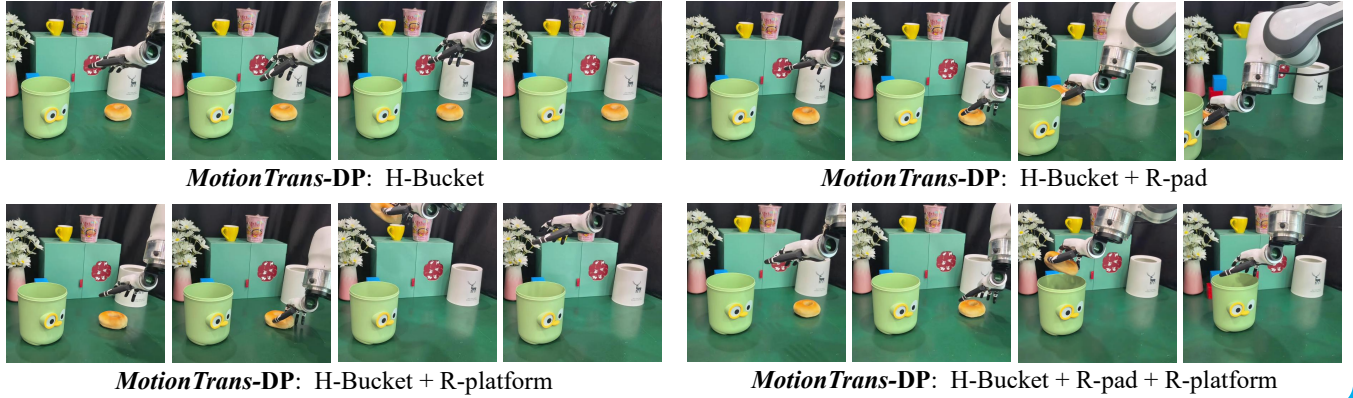


Fig. 11: The visualizations of the *MotionTrans*-DP results for the “Bread-Bucket” task, trained con various combinations of human and robot tasks. By analyzing these results (Section IV-E), we suggest that motion transfer occurs through the use of motion in human data to support robot motion interpolation for generating motions for these human tasks.

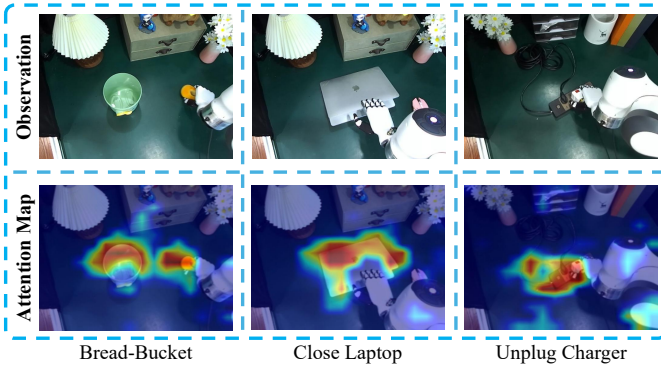


Fig. 12: The visualization of the attention map from the DINO encoder [48] for *MotionTrans*-DP, based on the Grad-Cam toolkit [56]. This shows that the vision encoder learns to focus on the target manipulation objects for tasks in human datasets, even when the embodiment changes to a robot during inference.

Example results are shown in Figure 12. The findings indicate that **the visual encoder attends to target objects in tasks of human data**, despite embodiment shifts at deployment (from human to robot). This embodiment-invariant, task-aware representation allows the policy to generate task-relevant motions during robot deployment and explains its ability to locate target objects even if they appeared only in human data.

#### (Q3.3) Is There a Scaling Trend in Motion Transfer?

Finally, we study the effect of task diversity and motion coverage. we hypothesize that a wider range of motion and task coverage may enhance the policy’s ability for motion interpolation and visual attention as mentioned before, thus leading to improved transfer performance. We verify this through subset training comparison experiment, similar to Q3.1. We introduce a new task subset “PP-set” compared to Q3.1, which includes data from two robot tasks “Mango-Bowl” and “Capybara-PPad” and two human tasks “Banana-Plate” and “Toy Bear-Box”.

The performance of *MotionTrans*-DP trained on different

subsets for the “Bread-Bucket” task is shown in the last 3 rows of Table IV. The results indicate a steady improvement with increased task coverage, suggesting that **motion transfer may benefit from broader task-related motion coverage**. While our study is based on a limited subset, these findings provide preliminary evidence of a potential scaling trend in human-to-robot motion transfer.

#### F. Supplementary Experiment Results

The evaluation results of all tasks in robot data are shown in Appendix VI-G. We also verify the robustness of our results concerning visual backgrounds in Appendix VI-H.

### V. CONCLUSION

In this paper, we propose *MotionTrans*, a framework that achieves motion-level learning from human data for end-to-end robot policies. The experiments show that our method achieves explicit human-to-robot motion transfer in a zero-shot setting and significantly improves finetuning performance in a few-shot setting. We identify two key factors for successful motion transfer: (1) cotraining with robot data, and (2) broader coverage of motions and tasks, which leads to better transfer performance. We hope that the new motion-centric insights that we propose could enhance the utilization of human data in robot policy learning in more effective ways.

**Limitations and Future Directions.** Our largest limitation is that the height perception ability of the policies is still limited, which causes them to sometimes fail to reach the correct height when considering in-the-wild scenes. This limitation arises from our monocular egocentric perception setting, which may be addressed by adding wrist camera for both human and robot hardware platforms [69, 61]. Another limitation is that our study is still limited to self-collected human dataset. Extending motion-level learning to larger, internet-scale datasets, as in [44], is left for future work.

# ACKNOWLEDGMENTS

This work is supported by the National Key R&D Program of China (2022ZD0161700), National Natural Science Foundation of China (62176135, 62476011), Shanghai Qi Zhi Institute Innovation Program SQZ202306, the Tsinghua University Dushi Program, the grant of National Natural Science Foundation of China (NSFC) 12201341.

We would like to express our sincere gratitude to Shuo Wang, Gu Zhang, Enshen Zhou, Haoxu Huang, Jialei Huang, Ruiqian Nai, Zhengrong Xue, Junmin Zhao, and Weirui Ye for their valuable discussions. We are especially grateful to Ruiqian Nai and Fanqi Lin for their assistance with the implementation of  $\pi_0$ -VLA, and to Yankai Fu for his support with the hardware implementation. Our thanks also extend to the SpiritAI and InspireRobot team for their assistance.

# REFERENCES

- [1] Gwon Hwan An, Siyeong Lee, Min-Woo Seo, Kugjin Yun, Won-Sik Cheong, and Suk-Ju Kang. Charuco board-based omnidirectional camera calibration method. *Electronics*, 7(12):421, 2018.
- [2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- [3] Shikhar Bahl, Russell Mendonca, Lili Chen, Unnat Jain, and Deepak Pathak. Affordances from human videos as a versatile representation for robotics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13778–13790, 2023.
- [4] Jose Barreiros, Andrew Beaulieu, Aditya Bhat, Rick Cory, Eric Cousineau, Hongkai Dai, Ching-Hsin Fang, Kunimatsu Hashimoto, Muhammad Zubair Irshad, Masha Itkina, et al. A careful examination of large behavior models for multitask dexterous manipulation. *arXiv preprint arXiv:2507.05331*, 2025.
- [5] Homanga Bharadhwaj, Debidatta Dwibedi, Abhinav Gupta, Shubham Tulsiani, Carl Doersch, Ted Xiao, Dhruv Shah, Fei Xia, Dorsa Sadigh, and Sean Kirmani. Gen2act: Human video generation in novel scenarios enables generalizable robot manipulation. *arXiv preprint arXiv:2409.16283*, 2024.
- [6] Hongzhe Bi, Lingxuan Wu, Tianwei Lin, Hengkai Tan, Zhizhong Su, Hang Su, and Jun Zhu. H-rdt: Human manipulation enhanced bimanual robotic manipulation. *arXiv preprint arXiv:2507.23523*, 2025.
- [7] 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.
- [8] Gary Bradski, Adrian Kaehler, et al. Opencv. *Dr. Dobb’s journal of software tools*, 3(2), 2000.
- [9] Qingwen Bu, Yanting Yang, Jisong Cai, Shenyuan Gao, Guanghui Ren, Maoqing Yao, Ping Luo, and Hongyang Li. Univla: Learning to act anywhere with task-centric latent actions. *arXiv preprint arXiv:2505.06111*, 2025.
- [10] Sirui Chen, Chen Wang, Kaden Nguyen, Li Fei-Fei, and C Karen Liu. Arcap: Collecting high-quality human demonstrations for robot learning with augmented reality feedback. *arXiv preprint arXiv:2410.08464*, 2024.
- [11] Xuxin Cheng, Jialong Li, Shiqi Yang, Ge Yang, and Xiaolong Wang. Open-television: Teleoperation with immersive active visual feedback. *arXiv preprint arXiv:2407.01512*, 2024.
- [12] 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*, page 02783649241273668, 2023.
- [13] Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. *arXiv preprint arXiv:2402.10329*, 2024.
- [14] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Evangelos Kazakos, Jian Ma, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision*, 130(1):33–55, 2022.
- [15] Jakob Engel, Kiran Somasundaram, Michael Goesele, Albert Sun, Alexander Gamino, Andrew Turner, Arjang Talattof, Arnie Yuan, Bilal Souti, Brighid Meredith, et al. Project aria: A new tool for egocentric multi-modal ai research. *arXiv preprint arXiv:2308.13561*, 2023.
- [16] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. In *Conference on robot learning*, pages 357–368. PMLR, 2017.
- [17] Sami Haddadin, Sven Parusel, Lars Johansmeier, Saskia Golz, Simon Gabl, Florian Walch, Mohamadreza Sabaghian, Christoph Jähne, Lukas Hausperger, and Simon Haddadin. The franka emika robot: A reference platform for robotics research and education. *IEEE Robotics & Automation Magazine*, 29(2):46–64, 2022.
- [18] Ryan Hoque, Peide Huang, David J Yoon, Mouli Siva-purapu, and Jian Zhang. Egodex: Learning dexterous manipulation from large-scale egocentric video. *arXiv preprint arXiv:2505.11709*, 2025.
- [19] Cheng-Chun Hsu, Bowen Wen, Jie Xu, Yashraj Narang, Xiaolong Wang, Yuke Zhu, Joydeep Biswas, and Stan Birchfield. Spot: Se (3) pose trajectory diffusion for object-centric manipulation. *arXiv preprint arXiv:2411.00965*, 2024.
- [20] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- [21] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail,



- Michael Equi, Chelsea Finn, Niccolo Fusai, et al. pi\_0.5: a vision-language-action model with open-world generalization. *arXiv preprint arXiv:2504.16054*, 2025.
- [22] Vidhi Jain, Maria Attarian, Nikhil J Joshi, Ayzaan Wahid, Danny Driess, Quan Vuong, Pannag R Sanketi, Pierre Sermanet, Stefan Welker, Christine Chan, et al. Vid2robot: End-to-end video-conditioned policy learning with cross-attention transformers. *arXiv preprint arXiv:2403.12943*, 2024.
- [23] Yuheng Ji, Huajie Tan, Jiayu Shi, Xiaoshuai Hao, Yuan Zhang, Hengyuan Zhang, Pengwei Wang, Mengdi Zhao, Yao Mu, Pengju An, et al. Robobrain: A unified brain model for robotic manipulation from abstract to concrete. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1724–1734, 2025.
- [24] Siddharth Karamcheti, Suraj Nair, Annie S Chen, Thomas Kollar, Chelsea Finn, Dorsa Sadigh, and Percy Liang. Language-driven representation learning for robotics. *arXiv preprint arXiv:2302.12766*, 2023.
- [25] Simar Kareer, Dhruv Patel, Ryan Punamiya, Pranay Mathur, Shuo Cheng, Chen Wang, Judy Hoffman, and Danfei Xu. Egomimic: Scaling imitation learning via egocentric video. *arXiv preprint arXiv:2410.24221*, 2024.
- [26] 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.
- [27] Hanjung Kim, Jaehyun Kang, Hyolim Kang, Meedeum Cho, Seon Joo Kim, and Youngwoon Lee. Uniskill: Imitating human videos via cross-embodiment skill representations. *arXiv preprint arXiv:2505.08787*, 2025.
- [28] Yuxuan Kuang, Junjie Ye, Haoran Geng, Jiageng Mao, Congyue Deng, Leonidas Guibas, He Wang, and Yue Wang. Ram: Retrieval-based affordance transfer for generalizable zero-shot robotic manipulation. *arXiv preprint arXiv:2407.04689*, 2024.
- [29] Marion Lepert, Jiaying Fang, and Jeannette Bohg. Masquerade: Learning from in-the-wild human videos using data-editing. *arXiv preprint arXiv:2508.09976*, 2025.
- [30] Marion Lepert, Jiaying Fang, and Jeannette Bohg. Phantom: Training robots without robots using only human videos. *arXiv preprint arXiv:2503.00779*, 2025.
- [31] Guangrun Li, Yaoxu Lyu, Zhuoyang Liu, Chengkai Hou, Jieyu Zhang, and Shanghang Zhang. H2r: A human-to-robot data augmentation for robot pre-training from videos. *arXiv preprint arXiv:2505.11920*, 2025.
- [32] Fanqi Lin, Yingdong Hu, Pingyue Sheng, Chuan Wen, Jiacheng You, and Yang Gao. Data scaling laws in imitation learning for robotic manipulation. *arXiv preprint arXiv:2410.18647*, 2024.
- [33] Fanqi Lin, Ruiqian Nai, Yingdong Hu, Jiacheng You, Junming Zhao, and Yang Gao. Onetwovla: A unified vision-language-action model with adaptive reasoning. *arXiv preprint arXiv:2505.11917*, 2025.
- [34] Jiaming Liu, Mengzhen Liu, Zhenyu Wang, Lily Lee, Kaichen Zhou, Pengju An, Senqiao Yang, Renrui Zhang, Yandong Guo, and Shanghang Zhang. Robomamba: Multimodal state space model for efficient robot reasoning and manipulation. *arXiv e-prints*, pages arXiv–2406, 2024.
- [35] Jiaming Liu, Hao Chen, Pengju An, Zhuoyang Liu, Renrui Zhang, Chenyang Gu, Xiaoqi Li, Ziyu Guo, Sixiang Chen, Mengzhen Liu, et al. Hybridvla: Collaborative diffusion and autoregression in a unified vision-language-action model. *arXiv preprint arXiv:2503.10631*, 2025.
- [36] Mengzhen Liu, Mengyu Wang, Henghui Ding, Yilong Xu, Yao Zhao, and Yunchao Wei. Segment anything with precise interaction. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 3790–3799, 2024.
- [37] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European conference on computer vision*, pages 38–55. Springer, 2024.
- [38] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024.
- [39] Vincent Liu, Ademi Adeniji, Haotian Zhan, Siddhant Haldar, Raunaq Bhirangi, Pieter Abbeel, and Lerrel Pinto. Egozero: Robot learning from smart glasses. *arXiv preprint arXiv:2505.20290*, 2025.
- [40] Yangcen Liu, Woo Chul Shin, Yunhai Han, Zhenyang Chen, Harish Ravichandar, and Danfei Xu. Immimic: Cross-domain imitation from human videos via mapping and interpolation. In *3rd RSS Workshop on Dexterous Manipulation: Learning and Control with Diverse Data*.
- [41] Yunze Liu, Yun Liu, Che Jiang, Kangbo Lyu, Weikang Wan, Hao Shen, Boqiang Liang, Zhoujie Fu, He Wang, and Li Yi. Hoi4d: A 4d egocentric dataset for category-level human-object interaction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21013–21022, 2022.
- [42] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [43] Chenhao Lu, Xuxin Cheng, Jialong Li, Shiqi Yang, Mazeyu Ji, Chengjing Yuan, Ge Yang, Sha Yi, and Xiaolong Wang. Mobile-television: Predictive motion priors for humanoid whole-body control. *arXiv preprint arXiv:2412.07773*, 2024.
- [44] Hao Luo, Yicheng Feng, Wanpeng Zhang, Sipeng Zheng, Ye Wang, Haoqi Yuan, Jiazheng Liu, Chaoyi Xu, Qin Jin, and Zongqing Lu. Being-h0: Vision-language-action pretraining from large-scale human videos. *arXiv preprint arXiv:2507.15597*, 2025.
- [45] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, et al. Where are we

in the search for an artificial visual cortex for embodied intelligence? *Advances in Neural Information Processing Systems*, 36:655–677, 2023.

- [46] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022.
- [47] Yaru Niu, Yunzhe Zhang, Mingyang Yu, Changyi Lin, Chenhao Li, Yikai Wang, Yuxiang Yang, Wenhao Yu, Tingnan Zhang, Zhenzhen Li, et al. Human2locoman: Learning versatile quadrupedal manipulation with human pretraining. *arXiv preprint arXiv:2506.16475*, 2025.
- [48] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafranec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- [49] 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.
- [50] Sungjae Park, Homanga Bharadhwaj, and Shubham Tulsiani. Demodiffusion: One-shot human imitation using pre-trained diffusion policy. *arXiv preprint arXiv:2506.20668*, 2025.
- [51] Shivansh Patel, Shraddha Mohan, Hanlin Mai, Unnat Jain, Svetlana Lazebnik, and Yunzhu Li. Robotic manipulation by imitating generated videos without physical demonstrations. *arXiv preprint arXiv:2507.00990*, 2025.
- [52] Karl Pertsch, Kyle Stachowicz, Brian Ichter, Danny Driess, Suraj Nair, Quan Vuong, Oier Mees, Chelsea Finn, and Sergey Levine. Fast: Efficient action tokenization for vision-language-action models. *arXiv preprint arXiv:2501.09747*, 2025.
- [53] Yuzhe Qin, Wei Yang, Binghao Huang, Karl Van Wyk, Hao Su, Xiaolong Wang, Yu-Wei Chao, and Dieter Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. *arXiv preprint arXiv:2307.04577*, 2023.
- [54] Ri-Zhao Qiu, Shiqi Yang, Xuxin Cheng, Chaitanya Chawla, Jialong Li, Tairan He, Ge Yan, David J Yoon, Ryan Hoque, Lars Paulsen, et al. Humanoid policy~ human policy. *arXiv preprint arXiv:2503.13441*, 2025.
- [55] Juntao Ren, Priya Sundareshan, Dorsa Sadigh, Sanjiban Choudhury, and Jeannette Bohg. Motion tracks: A unified representation for human-robot transfer in few-shot imitation learning. *arXiv preprint arXiv:2501.06994*, 2025.
- [56] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017.
- [57] Junyao Shi, Zhuolun Zhao, Tianyou Wang, Ian Pedroza, Amy Luo, Jie Wang, Jason Ma, and Dinesh Jayaraman. Zeromimic: Distilling robotic manipulation skills from web videos. *arXiv preprint arXiv:2503.23877*, 2025.
- [58] Modi Shi, Li Chen, Jin Chen, Yuxiang Lu, Chiming Liu, Guanghui Ren, Ping Luo, Di Huang, Maoqing Yao, and Hongyang Li. Is diversity all you need for scalable robotic manipulation? *arXiv preprint arXiv:2507.06219*, 2025.
- [59] Andreas Steiner, André Susano Pinto, Michael Tschannen, Daniel Keysers, Xiao Wang, Yonatan Bitton, Alexey Gritsenko, Matthias Minderer, Anthony Sherbondy, Shangbang Long, et al. Paligemma 2: A family of versatile vlms for transfer. *arXiv preprint arXiv:2412.03555*, 2024.
- [60] Chao Tang, Anxing Xiao, Yuhong Deng, Tianrun Hu, Wenlong Dong, Hanbo Zhang, David Hsu, and Hong Zhang. Functo: Function-centric one-shot imitation learning for tool manipulation. *arXiv preprint arXiv:2502.11744*, 2025.
- [61] Tony Tao, Mohan Kumar Srirama, Jason Jingzhou Liu, Kenneth Shaw, and Deepak Pathak. Dexwild: Dexterous human interactions for in-the-wild robot policies. *arXiv preprint arXiv:2505.07813*, 2025.
- [62] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*, 2014.
- [63] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017.
- [64] Shengjie Wang, Jiacheng You, Yihang Hu, Jiongye Li, and Yang Gao. Skil: Semantic keypoint imitation learning for generalizable data-efficient manipulation. *arXiv preprint arXiv:2501.14400*, 2025.
- [65] Adam Wei, Abhinav Agarwal, Boyuan Chen, Rohan Bosworth, Nicholas Pfaff, and Russ Tedrake. Empirical analysis of sim-and-real cotraining of diffusion policies for planar pushing from pixels. *arXiv preprint arXiv:2503.22634*, 2025.
- [66] Chuan Wen, Xingyu Lin, John So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point trajectory modeling for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- [67] Kun Wu, Chengkai Hou, Jiaming Liu, Zhengping Che, Xiaozhu Ju, Zhuqin Yang, Meng Li, Yinuo Zhao, Zhiyuan Xu, Guang Yang, et al. Robomind: Benchmark on multi-embodiment intelligence normative data for robot manipulation. *arXiv preprint arXiv:2412.13877*, 2024.
- [68] Mengda Xu, Zhenjia Xu, Yinghao Xu, Cheng Chi, Gordon Wetzstein, Manuela Veloso, and Shuran Song. Flow as the cross-domain manipulation interface. *arXiv preprint arXiv:2407.15208*, 2024.
- [69] Mengda Xu, Han Zhang, Yifan Hou, Zhenjia Xu, Linxi Fan, Manuela Veloso, and Shuran Song. Dexumi: Using human hand as the universal manipulation interface for dexterous manipulation. *arXiv preprint*

*arXiv:2505.21864*, 2025.

- [70] Ruihan Yang, Qinxu Yu, Yecheng Wu, Rui Yan, Borui Li, An-Chieh Cheng, Xueyan Zou, Yunhao Fang, Hongxu Yin, Sifei Liu, et al. Egovla: Learning vision-language-action models from egocentric human videos. *arXiv preprint arXiv:2507.12440*, 2025.
- [71] Seonghyeon Ye, Joel Jang, Byeongguk Jeon, Sejeun Joo, Jianwei Yang, Baolin Peng, Ajay Mandlekar, Reuben Tan, Yu-Wei Chao, Bill Yuchen Lin, et al. Latent action pre-training from videos. *arXiv preprint arXiv:2410.11758*, 2024.
- [72] Weirui Ye, Yunsheng Zhang, Mengchen Wang, Shengjie Wang, Xianfan Gu, Pieter Abbeel, and Yang Gao. Foundation reinforcement learning: towards embodied generalist agents with foundation prior assistance. 2023.
- [73] Chengbo Yuan, Geng Chen, Li Yi, and Yang Gao. Self-supervised monocular 4d scene reconstruction for egocentric videos. *arXiv preprint arXiv:2411.09145*, 2024.
- [74] Chengbo Yuan, Chuan Wen, Tong Zhang, and Yang Gao. General flow as foundation affordance for scalable robot learning. *arXiv preprint arXiv:2401.11439*, 2024.
- [75] Chengbo Yuan, Suraj Joshi, Shaoting Zhu, Hang Su, Hang Zhao, and Yang Gao. Roboengine: Plug-and-play robot data augmentation with semantic robot segmentation and background generation. *arXiv preprint arXiv:2503.18738*, 2025.
- [76] Tony Z Zhao, Jonathan Tompson, Danny Driess, Pete Florence, Kamyar Ghasemipour, Chelsea Finn, and Ayzaan Wahid. Aloha unleashed: A simple recipe for robot dexterity. *arXiv preprint arXiv:2410.13126*, 2024.
- [77] Enshen Zhou, Jingkun An, Cheng Chi, Yi Han, Shanyu Rong, Chi Zhang, Pengwei Wang, Zhongyuan Wang, Tiejun Huang, Lu Sheng, et al. Roborefer: Towards spatial referring with reasoning in vision-language models for robotics. *arXiv preprint arXiv:2506.04308*, 2025.
- [78] Huayi Zhou, Ruixiang Wang, Yunxin Tai, Yueci Deng, Guiliang Liu, and Kui Jia. You only teach once: Learn one-shot bimanual robotic manipulation from video demonstrations. *arXiv preprint arXiv:2501.14208*, 2025.
- [79] Xiang Zhu, Yichen Liu, Hezhong Li, and Jianyu Chen. Learning generalizable robot policy with human demonstration video as a prompt. *arXiv preprint arXiv:2505.20795*, 2025.
- [80] 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.



## VI. APPENDIX

### A. Details of MotionTrans Dataset and All Tasks

Here we present the details of all tasks in *MotionTrans Dataset*. The visualization and descriptions / VLA-prompt of all 15 human tasks could be found in Figure 17 and Table V. All 15 robot tasks could be found in could be found in Figure 18 and Table VI.

### B. Rubrics of Motion Progress Score

Table VII provides the detailed rubrics for our Motion Progress Score metric. The scores are allocated to the different motions / stages required to complete the task, with a maximum score of 8 points.

### C. Calibration between VR Headset and RGB Camera

In human data collection (Section III-B), our goal is to record hand pose information captured by a VR device in the RGB camera’s coordinate system. To transform hand poses from the VR coordinate space to the RGB camera, we need to solve the transformation between the two cameras. We achieve this by applying a chain-style calibration. We place an ArUco calibration chessboard [1] on the table and ask users to sit facing it without moving their heads, as illustrated in Figure 14. We then perform two calibrations:

- **Camera-Chessboard Calibration** (Figure 14(a)). Solve  $T_{cam}$ , the pose of the RGB camera based on the chessboard coordinate, using the vision-based calibration method [1] (OpenCV library [8]).
- **VR-Chessboard Calibration** (Figure 14(b)). Solve  $T_{vr}$ , the pose of the VR camera based on the chessboard coordinate by asking the user to place an anchor block on the origin of the chessboard coordinate. We then directly read the coordinate of the anchor block (i.e., the origin) in the VR camera’s coordinate space using the VR app, thus obtaining  $T_{vr}$  by inverting the reading result. To improve placement accuracy, we use the depth sensing built into the VR headset to fit the desktop height (the blue plane in Figure 14(c)), allowing users to only adjust the anchor block’s position forward, backward, left, and right.

By using the chessboard as the bridge, the transformation used to convert hand poses from the VR to the RGB camera coordinate can be expressed as  $T_{cam}^{-1}T_{vr}$ .

### D. Policies Implementation Details.

For the robot observation-action space (Section III-A), we set the proprioception history  $T_p = 2$  and the action horizon  $T_A = 16$ . The representation of the rotation component of wrist poses is chosen as the first two rows of the rotation matrix, as demonstrated in [12]. For policy control, we use 10 fps for both data collection and policy inference. For Diffusion Policy (DP) backbone, the task-embedding dimension is set as 16. The proprioception state is encoded via a 4-layer MLP. The DINOv2 vision encoder utilizes DINOv2-base pretrained checkpoints [48], and during training, we unfreeze the weights of the DINOv2-base encoder. We first concatenate the task

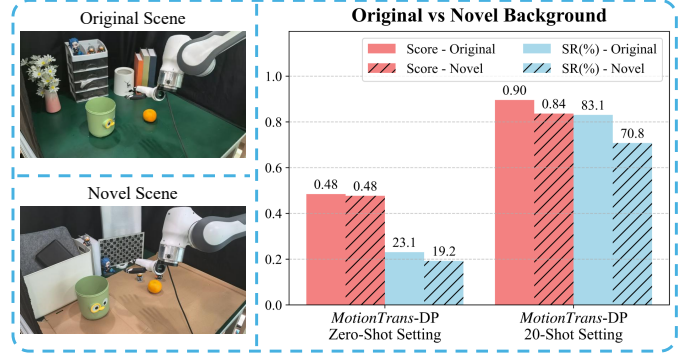


Fig. 13: Illustration of the visual background robustness experiment and results. All results are averaged across all 13 evaluation human tasks. For the novel background, the performance drops slightly but remains at a persuasive level. This prove the robustness of our motion transfer results.

embedding with the features from the vision and proprioception encoder, and then input the concatenated features into the U-Net-based Diffusion head for action generation [12].

### E. Domain Confusion Training Framework

We also tried the domain confusion framework [63, 62] in our earlier exploration. The key idea is to: (1) train a classifier to identify the embodiment from the features generated by the policy encoder. (2) The policy’s target is to generate embodiment-invariant features that mislead the classifier. These embodiment-invariant features thus facilitate better embodiment-agnostic knowledge transfer. (3) Train these two models in an adversarial manner, allowing them to improve themselves by competing against each other.

Following the domain adaptation framework, we additionally train a binary classifier  $C$  to classify whether a data point is from the human or robot domain. The input of  $C$  is the concatenation of image features and proprioception features generated by policy  $P$ , as demonstrated in Appendix VI-D. Since we only want  $C$  to classify based on embodiment/domain, rather than depend on shortcuts like task-specific content in  $D_{human} / D_{robot}$  (e.g., task-related objects in image observations), we trained  $C$  on an augmented version of  $D_{human} / D_{robot}$ : we first used GroundingDINO [37] / RoboSAM [75] to segment the robot out and then pasted it onto a novel MIL-texture background [16, 75]. We represent the augmented version as  $D_{human}^{aug} / D_{robot}^{aug}$ . The final training loss  $\mathcal{L}$  of policy  $P$  can be represented as:

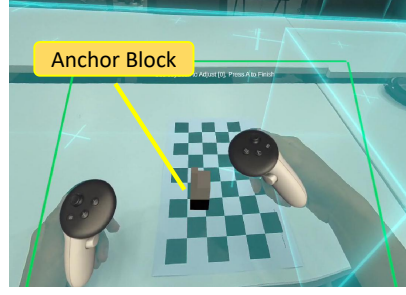
$$\mathcal{L}_{dc} = D_{KL}(C_{frozen}(P_{encoder}(s)) || U) \quad (1)$$

$$\mathcal{L} = \mathcal{L}_D + \alpha \mathcal{L}_{dc} \quad (2)$$

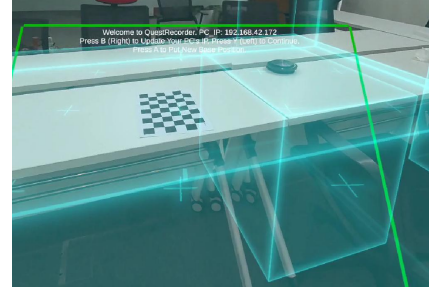
where  $\mathcal{L}_D$  is the imitation learning loss described in Section III-D,  $P_{encoder}$  is the image and proprioception encoder of policy  $P$ ,  $\mathcal{L}_{dc}$  is the domain confusion loss ([63]),  $U$  is a binary uniform distribution, and  $s$  is the data points randomly sampled from  $D_{human}^{aug} / D_{robot}^{aug}$ . The  $\mathcal{L}_{dc}$  encourages policy  $P$  to generate similar features when only considering embodiment differences, thereby leading to embodiment-invariant features that enhance



(a) Camera-Chessboard Calibration



(b) VR-Chessboard Calibration



(c) Height Sensing of VR

Fig. 14: The illustration of the calibration process used to transform data from the VR coordinate space to the RGB camera space. Detailed demonstration can be found in Appendix VI-C.

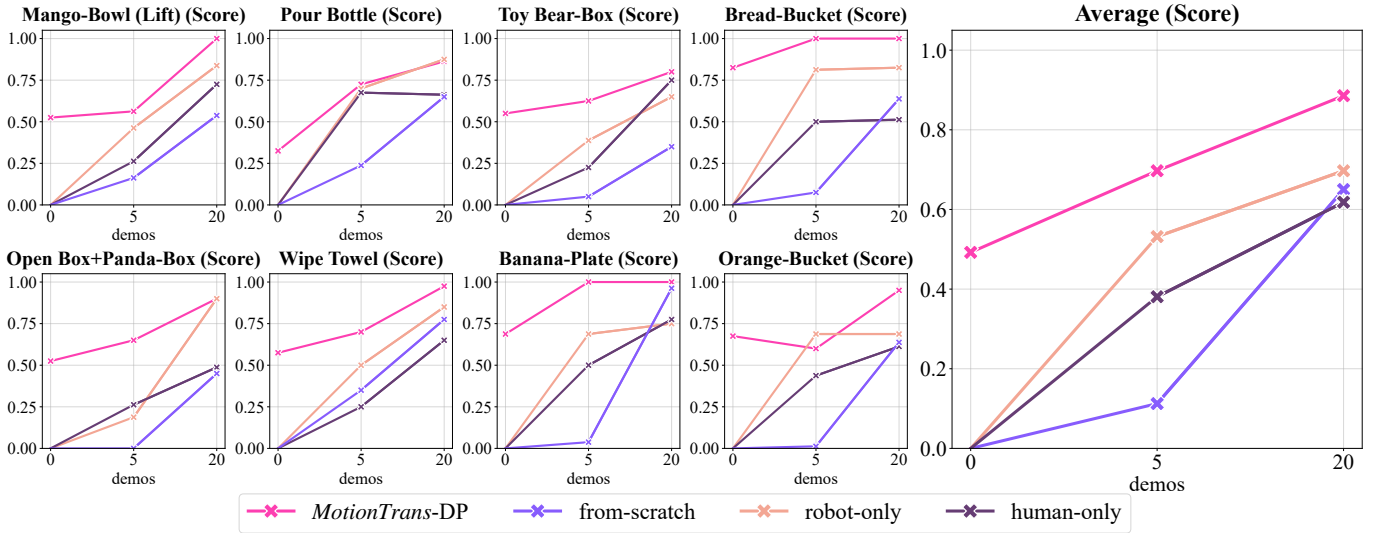


Fig. 15: Results of the Motion Progress Score for few-shot finetuning experiments.

human-to-robot knowledge transfer. We train policy  $P$  (based on  $\mathcal{L}$ ) and classifier  $C$  (based on binary cross entropy loss, BCELoss) in an adversarial manner, which means that we iteratively train these two models. When training one model, we freeze the weight of another model. More details can be found in [63]. The architecture of  $C$  is the same as that of  $P$ , except for changing the final MLP to a classifier head.

However, although we have tried our best to tune different settings, we still find that this kind of adversarial training tends to lead to mode collapse and training instability [2]. This is reflected in the sudden jumps in the value of a certain loss during the training process, as well as the robot’s inability to perform meaningful actions during downstream policy deployment. The key insight we gained from this experience is that rather than only relying on updating the algorithm or model, improving the scale and quality of data may be a more straightforward way to enhance transfer effectiveness. With enough task-related motion coverage, the simplest weighted cotraining framework shows the strongest transfer performance in our setting.

#### F. Few-shot Results of Motion Progress Score

The results of Motion Progress Score for few-shot experiment (Section IV-C) are shown in Figure 15. The conclusion drawn from the Motion Progress Score aligns with that from the Success Rate (Section IV-C).

#### G. Robot Tasks Experiment

In this section, we assess whether cotraining with human data can enhance task performance on robot data. To do this, we compare the Diffusion Policy (DP) trained on the complete *MotionTrans Dataset* (*MotionTrans-DP*) with a version trained solely on the robot data (robot-only). It is important to note that the human tasks does not overlap with robot tasks, which distinguishes our approach from previous cotraining works [25, 54, 47]. The results for all 15 robot tasks are displayed in Figure 16. Our findings indicate that, in our setting, cotraining with non-overlapping human tasks does not significantly impact the policy’s performance on robot tasks, with average success rates of 58.0% for the robot-only model and 58.7% for the *MotionTrans-DP*. We believe this is due to the following reasons:

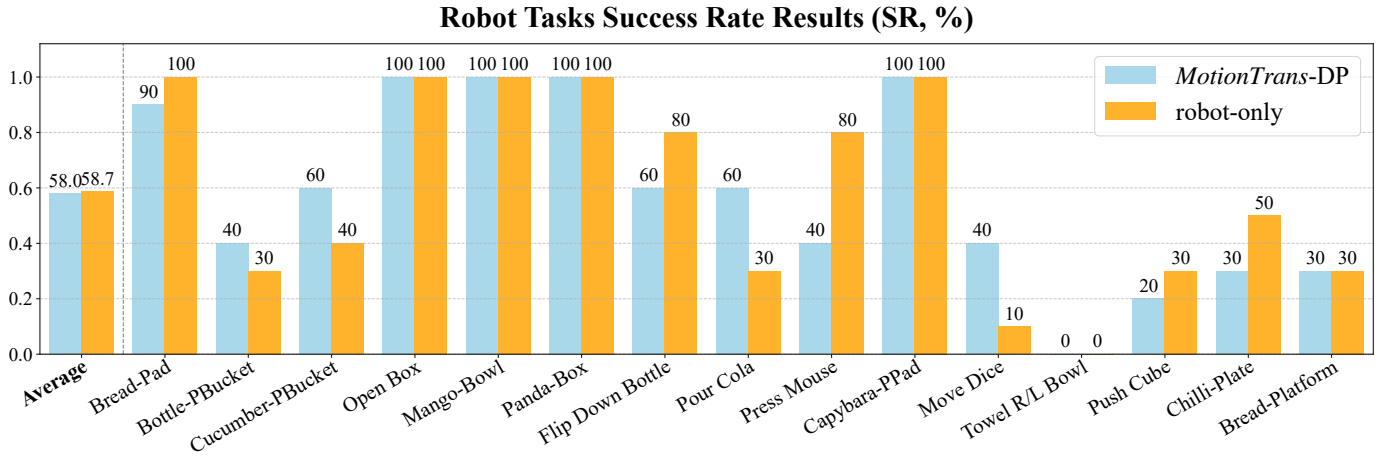


Fig. 16: Results of the success rate for zero-shot experiments on robot tasks. We conclude that when considering motion learning, cotraining with non-overlapping human tasks does not ensure the improvement of policy’s performance on robot tasks, with average success rates of 58.0% for the robot-only model and 58.7% for the *MotionTrans-DP*.

(1) We have already collected sufficient robot demonstrations for all 15 robot tasks in our datasets, which reduces the impact of cotraining (similar to the results in Section IV-C). (2) Although we have conducted data alignment (Section III-C), the motions from non-overlapping tasks still differ too much from those in the robot data, thereby providing insufficient auxiliary guidance for motion learning.

However, performance differences are observed in specific tasks. For example, in the “Pour Cola” task, cotraining with human data improved the policy’s grasping position, resulting in more stable pouring and better performance. Conversely, for the “Press Mouse” task, cotraining negatively affected the final “press” action, leading to instances where the policy only made contact without pressing. Tasks like “Towel R/L Bowl” exhibit low success rates due to insufficient height generalization capabilities of our policies, a limitation we discuss in the limitations section (Section V).

#### H. Visual Background Robustness

Finally, we verify the visual robustness of our experiment results against scene background [70]. We change the background from our default “green table” scenes (mentioned in the dataset part in Section IV-A) to a new scene, as shown in Figure 13, and evaluate Diffusion Policy (DP) performance for both zero-shot and 20-shot settings. The results are averaged across all 13 evaluation human tasks and shown on the right side of Figure 13. We observe that although the performance drops slightly, it still maintains a non-trivial Motion Progress Score and success rate. This proves the robustness of our results on motion-level human data learning. Note that this does not mean we achieve in-the-wild manipulation ability [61], which is not the main focus of this paper and will be discussed in the limitations section.

Human Tasks	Description / VLA-prompt
Unplug Charger	unplug the white charger.
Bread-Bucket	drop bread to the green bucket.
Press Stapler	press the stapler.
Orange-Bucket	put orange to the green bucket.
Wipe Towel	wipe blue towel on the table and push it to the bulky bottle.
Close Laptop	close silver laptop.
Mango-Bowl (Bypass)	put mango to pink bowl while avoiding obstacle by bypassing.
Mango-Bowl (Lift)	put mango to the pink bowl while avoiding obstacle by lifting.
Press Dice	press red dice to make it rotation.
Banana-Plate	put banana to the white plate.
Pour Bottle	pour bottle to the pink bowl.
Toy Bear-Box	put toy bear to the black box.
Open Box + Pand-Box	first open the white cap style box then put toy panda to the box.
Fold Towel	fold the blue towel.
Pour Milk Bottle	pour milk bottle to the yellow pan.

TABLE V: All 15 human tasks with detailed descriptions ( $\pi_0$ -VLA-prompt).

Robot Tasks	Description / VLA-prompt
Push Cube	push orange cube to the bulky bottle.
Panda-Box	put toy panda to the box.
Bread-Pad	put bread to the red pad.
Open Box	open the white cap style box.
Bottle-PBucket	drop black bottle to purple bucket.
Pour Cola	pour cola to the red cup.
Move Dice	move red dice to the bulky bottle.
Flip Down Bottle	flip down the black bottle.
Press Mouse	press the pink mouse.
Bread-Platform	put bread to the high black platform.
Capybara-PPad	put Capybara to the purple pad.
Chilli-Plate	put chilli to the white plate.
Towel R/L Bowl	wipe blue towel on the table and push it left or right to the pink bowl.
Mango-Bowl	put mango to the pink bowl.
Cucumber-PBucket	put cucumber to purple bucket.

TABLE VI: All 15 robot tasks with detailed descriptions ( $\pi_0$ -VLA-prompt).

Human Tasks	Rubrics of Motion Progress Score
Mango-Bowl (Bypass)	(1) show reach-grasp; (1) successful grasp; (2) show bypassing; (2) successful bypassing; (1) show reach-put; (1) successful put;
Mango-Bowl (Lifting)	(1) show reach-grasp; (1) successful grasp; (1) show lifting; (2) successful lifting; (2) show down-putting; (1) successful put;
Pour Bottle	(1) show reach-grasp; (1) successful grasp; (2) show rotation; (2) successful pouring; (2) good pour position;
Toy Bear-Box	(2) show reach-grasp; (2) successful grasp; (2) show reach-put; (1) successful put; (1) good put position;
Bread-Bucket	(1) show reach-grasp; (1) successful grasp; (2) show reach-put; (2) successful put; (2) good put height;
Close Laptop	(2) show reach-press; (2) press finish < 30 degrees; (2) press finish < 15 degrees; (2) press finish = 0 degrees;
Press Stapler	(2) show reach-press; (2) success contact; (2) good contact position; (2) press down;
Unplug Charger	(2) show reach-grasp; (1) successful grasp; (1) show lifting; (2) successful unplug; (2) still holding after unplugging;
Open Box + Panda-Box	(2) open the white box; (1) continue to move; (1) no stop after open the box; (1) reach the panda; (1) successful grasp the panda; (2) successful put;
Wipe Towel	(2) show reach-press; (2) successful press; (2) show pushing (including retry); (2) successful pushing;
Banana-Plate	(1) show reach-grasp; (2) successful grasp; (2) show reach-put; (2) successful put; (1) good put height;
Orange-Bucket	(1) show reach-grasp; (2) successful grasp; (2) show reach-put; (2) successful put; (1) good put height;
Press Dice	(1) show reach-press; (1) successful contact; (2) show press; (2) press > 5 cm; (2) successful press to make it rotate;

TABLE VII: The rubrics of Motion Progress Score for all 13 evaluation human tasks. The scores are allocated to the different motions / stages required to complete the task, with a maximum score of 8 points. The number in () is the score of that stage. The “show reach-{action}” rubric means policy shows approaching motion to achieve {action}.





Fig. 17: The visualizations of all 15 human tasks in the egocentric view.

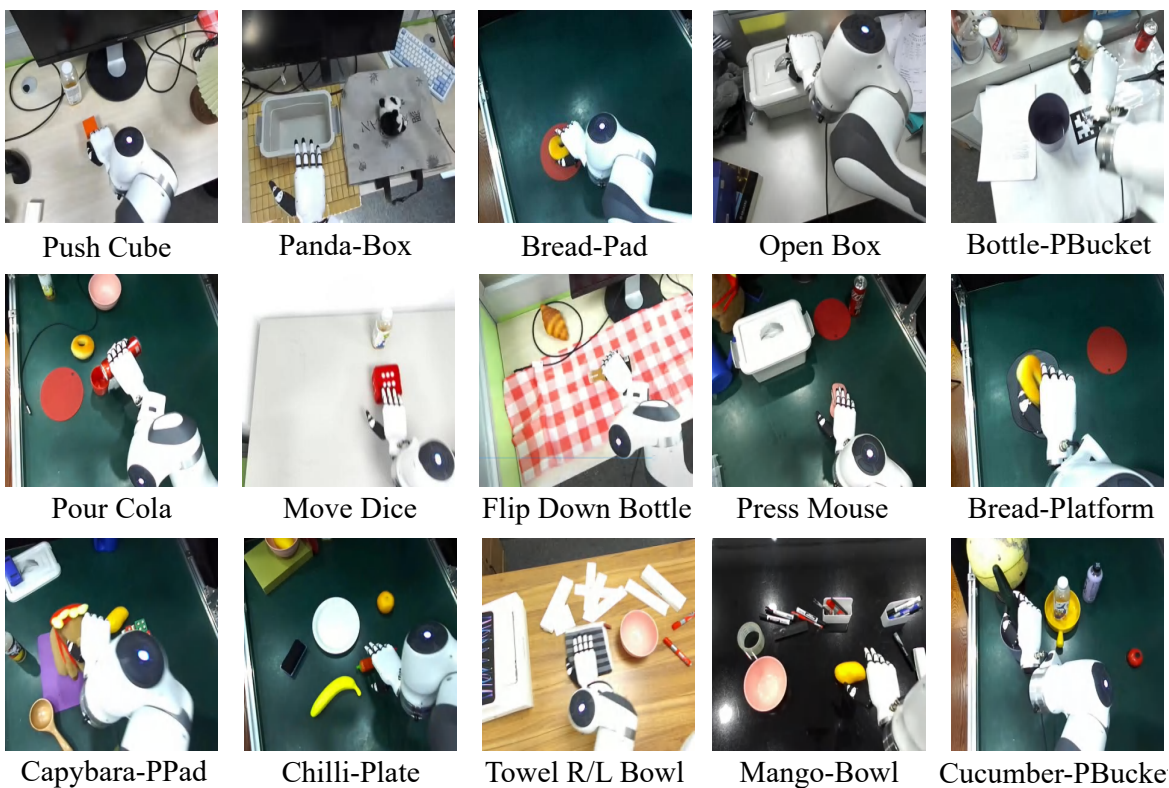


Fig. 18: The visualizations of all 15 robot tasks in the egocentric view.