

ManipLLM: Embodied Multimodal Large Language Model for Object-Centric Robotic Manipulation

Xiaoqi Li¹, Mingxu Zhang², Yiran Geng¹, Haoran Geng¹, Yuxing Long¹,
 Yan Shen¹, Renrui Zhang³, Jiaming Liu¹, Hao Dong[†]

¹School of Computer Science, Peking University

² Beijing University of Posts and Telecommunications ³ MMLab, CUHK

Abstract

Robot manipulation relies on accurately predicting contact points and end-effector directions to ensure successful operation. However, learning-based robot manipulation, trained on a limited category within a simulator, often struggles to achieve generalizability, especially when confronted with extensive categories. Therefore, we introduce an innovative approach for robot manipulation that leverages the robust reasoning capabilities of Multimodal Large Language Models (MLLMs) to enhance the stability and generalization of manipulation. By fine-tuning the injected adapters, we preserve the inherent common sense and reasoning ability of the MLLMs while equipping them with the ability for manipulation. The fundamental insight lies in the introduced fine-tuning paradigm, encompassing object category understanding, affordance prior reasoning, and object-centric pose prediction to stimulate the reasoning ability of MLLM in manipulation. During inference, our approach utilizes an RGB image and text prompt to predict the end effector’s pose in chain of thoughts. After the initial contact is established, an active impedance adaptation policy is introduced to plan the upcoming waypoints in a closed-loop manner. Moreover, in real world, we design a test-time adaptation (TTA) strategy for manipulation to enable the model better adapt to the current real-world scene configuration. Experiments in simulator and real-world show the promising performance of ManipLLM. More details and demonstrations can be found at <https://sites.google.com/view/manipllm>.

1. Introduction

As robot manipulation requires robots to interact with diverse objects, the robustness and explainability of low-level action prediction become essential for manipulating reliability. While certain approaches [9, 11–13, 22, 23]

[†]Corresponding author: hao.dong@pku.edu.cn

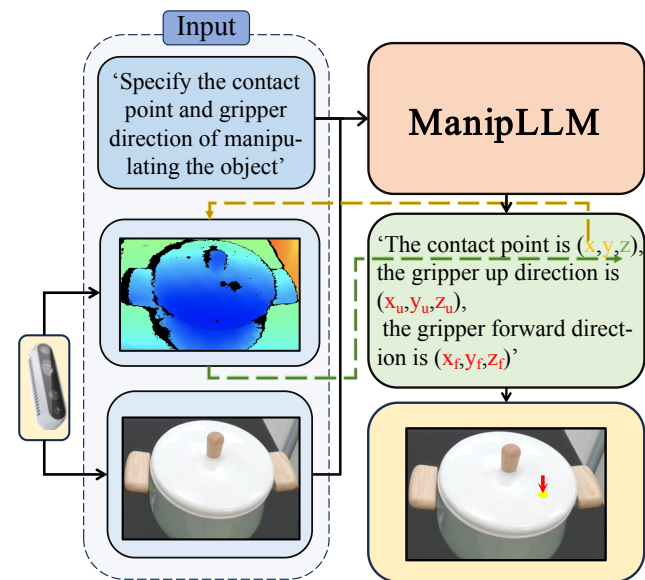


Figure 1. The prediction of ManipLLM. Given the text prompt, RGB image, and depth map inputs, we obtain 3D contact point (x, y, z) . Here, x and y represent the pixel coordinates in the image predicted by ManipLLM, while z corresponds to the depth obtained from the depth camera. Additionally, ManipLLM predicts the gripper’s up direction (x_u, y_u, z_u) and forward direction (x_f, y_f, z_f) , forming the end-effector $SO(3)$ rotation.

demonstrate impressive performance, they often sacrifice interpretability by treating low-level manipulation prediction as a black-box prediction problem and lack the inherent common-sense reasoning abilities inherent in humans, limiting their capacity to manipulate wide-spread categories of objects.

Existing advancements in Multimodal Large Language Models (MLLMs)[1, 19, 22, 38] highlight their proficiency in common sense reasoning and remarkable generalization in vision tasks [2, 8]. However, training Multi-Label Language Models (MLLMs) directly to learn robotic low-level

action trajectories (*i.e.* end-effector trajectories) [4, 40] poses challenges in generalization due to minimal low-level action samples in their pretraining data. Consequently, MLLMs lack prior knowledge in this field while successful training for these tasks necessitates extensive data to achieve desired generalization ability. In contrast, MLLMs exhibit robust capabilities in comprehending objects and demonstrate significant generalization abilities. Given these considerations, transforming MLLMs to object-centric manipulation proves more efficient. This then raises an important question: How can we harness MLLMs to facilitate object-centric robot manipulation? The major challenge is how to enable MLLMs to understand the geometric structure of objects (such as their axis) to predict the movable contact positions for object-centric manipulation. Furthermore, it remains unexplored whether these models, which take 2D inputs, can also predict 3D end-effector directions.

In this study, we aim to exploit the common sense and reasoning ability embedded within MLLMs [26, 38] to realize promising robot manipulation performance. To accomplish this, during training, in order to preserve the powerful ability of MLLMs and empower them with manipulation ability, we only finetune the injected learnable adapters [15] on MLLMs. Furthermore, we design an intricate training paradigm and formulate fine-tuning tasks, including object category identification, affordance prior reasoning, and manipulation-aware pose prediction. The affordance prior considers the geometric intrinsics of the object and reflects the probability of generating movement when acting on a particular pixel. Through this training paradigm, we enable MLLMs to recognize the object at the category level, understand which regions can be manipulated and which cannot at the region level, and ultimately generate precise coordinates and directions for manipulation at the pose level.

During inference, we employ the chain-of-thought [29] flow, consistent with the training flow, to make the model’s predictions more interpretable. This allows us to understand the thought process of the model in obtaining the final pose prediction. The final prediction is depicted in Fig. 1. Given an RGB image featuring an object and a text prompt, our method generates the contact pixel coordinate on the 2D image and an end-effector direction. Additionally, depth information projects the pixel coordinate into 3D space. After the initial contact is established, we design an active impedance adaptation policy to determine the movement by forecasting the upcoming waypoints in a close-loop manner. Specifically, this module applies small forces in the surrounding directions based on the current pose. It aims to identify the direction that yields the maximum movement, which is then chosen as the next pose. This method relies on force feedback generated along the axes and the object to adaptively adjust the direction and predict the trajectory.

In real-world testing, we observe challenges that may

diverge from the simulated learning environment. For instance, in the real world, manipulating a door with a handle using a short suction gripper might require placing the end-effector at a distance from the handle to prevent collisions, which is different from the simulator. To address these variations, we draw inspiration from test-time adaptation (TTA) [21, 35]. TTA involves adjusting partial model parameters during inference based on the current test sample, enhancing the model’s performance for specific real-world scenarios. Subsequently, We design a TTA strategy tailored for robot manipulation, aiming to refine the model’s understanding of real-world configurations. Specifically, with the current test sample, we utilize the outcome of manipulation success or failure to supervise the model’s assessment of whether the predicted pose can result in a successful manipulation and only update partial parameters. This allows the model to retain its original capabilities and adapt to the target domain by distinguishing between effective and ineffective poses in the target domain. Since the model has learned to predict poses that are more likely to result in successful manipulations, when facing upcoming samples, the model tends to predict effective poses, thus enhancing the performance under specific real-world configurations.

Benefiting from the MLLMs and the designed paradigm, our approach exhibits generalization and common-sense reasoning ability in manipulation. Experiments show that in the simulator, our method achieves a promising manipulation success rate across 30 categories. Meanwhile, in real-world experiments, our method shows strong generalization ability, with or without TTA strategy. More real-world videos are shown in the supplement.

In summary, our contributions are as follows:

- We innovatively present a simple yet effective approach that transforms the ability of MLLMs into object-centric robot manipulation.
- We design a chain-of-thought fine-tuning and inference strategy that exploits MLLMs’ reasoning ability to enable robust and explainable end-effector’s pose predictions.
- Experiments across extensive categories demonstrate the generalization ability of ManipLLM.

2. Related Works

2.1. Robotic Manipulation

Robotic manipulation has emerged as a pivotal research domain due to its extensive applicability. One widely used approach is state-based reinforcement learning (RL) [3, 12, 17, 36]. Some works have identified the possibility of using the pure state as the policy input [3]. However, when it comes to more complex settings, vision-based observation [6, 10, 14, 16, 23, 27, 28, 30, 32–34, 40] becomes necessary to perceive the environment and understand the complex scene and objects [5, 18]. Where2Act [23] proposes

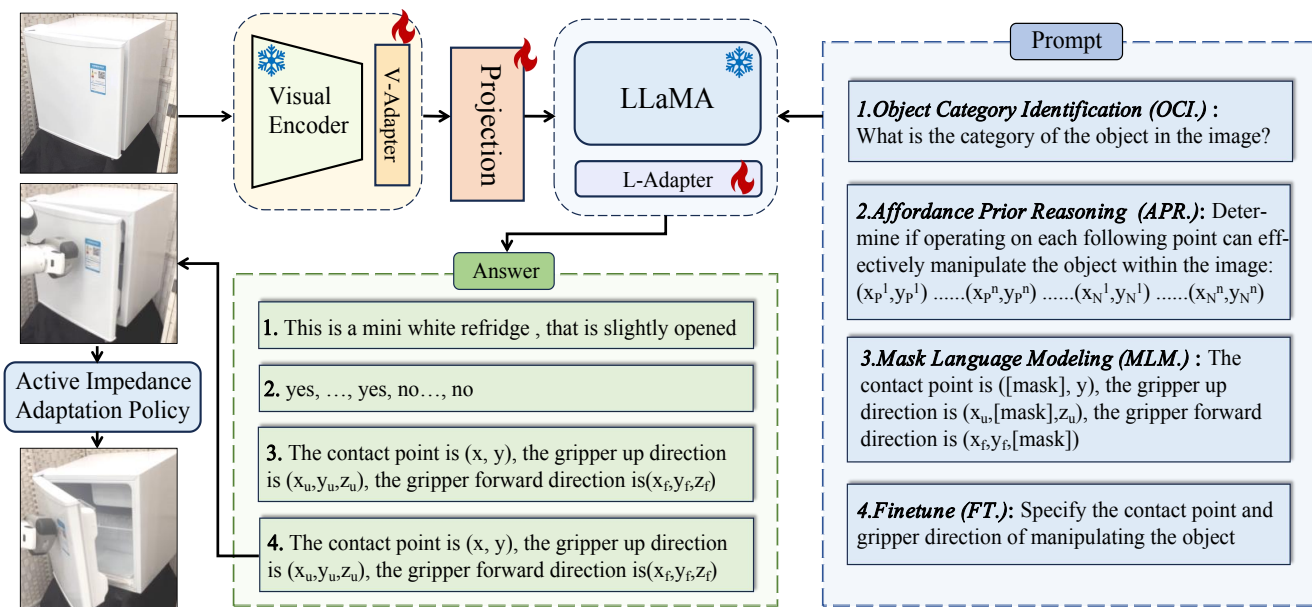


Figure 2. Training details of ManipLLM. This paradigm contains four training tasks, enabling the model to recognize the current object (category-level), understand which regions can be manipulated (region-level), and finally generate a precise end-effector pose (pose-level).

novel networks to predict actionable pixels and movable regions in objects, enabling meaningful interaction in various environments. Flowbot3d [6] also explores a vision-based method for perceiving and manipulating 3D articulated objects through predicting point-wise motion flow. Furthermore, VoxPoser [16] synthesizes adaptable robot trajectories through 3D value maps derived from large language models, based on natural language instructions. RT2 [40], which transfers information to actions, holds promise for adapting more rapidly to novel situations.

However, although these methods achieve noteworthy accomplishments, they formulate the task as a black-box prediction, decreasing its interpretability. This becomes extremely severe when confronted with extensive categories of objects. To reconcile such, ManipLLM harnesses the common sense knowledge and reasoning ability embedded within MLLMs to strengthen robot manipulation performance. We design intricate finetuning and inference strategies to enable interpretable object-centric pose prediction.

2.2. Multimodal Large Language Models

Extensive language models, *i.e.*, LLaMa [26], GPT3 [7] exhibit proficiency in a variety of language tasks given their powerful reasoning ability. Building upon these, Multimodal Large Language Models [1, 19, 20, 37, 38] are introduced to bridge RGB visual images and text. The representative LLaMa-Adapter [38] generalizes to image conditions for multi-modal reasoning, achieving competitive results in both vision and multi-modal tasks.

However, despite the considerable achievements of

MLLMs, their object-centric manipulation ability is still under-explored. Aiming to bridge this gap, our work pioneers in injecting manipulation capabilities into existing MLLMs while preserving their original reasoning ability. By doing so, the finetuned model not only possesses precise manipulation ability but is also capable of dealing with diverse category objects under interpretable thinking.

3. Method

3.1. Fine-tuning Strategy

In this section, we demonstrate how we empower MLLMs with manipulation capabilities. As shown in Fig. 2, we design fine-tuning tasks at the category level, region level, and pose level, allowing the model to progressively and reasonably predict poses for object-centric robot manipulation.

3.1.1 Model Architecture

We adopt the MLLM, LLaMa-Adapter [38], as our backbone and follow its training strategy. Given an RGB image $I \in \mathbb{R}^{H \times W \times 3}$, we adopt the visual encoder of CLIP [25] to extract its visual feature. While text prompts T are encoded into a text feature using the tokenizer of the pre-trained LLaMa [26]. After aligning visual and text feature representation with the multi-modal projection module, LLaMa is required to conduct multi-modal understanding and give correct answers. During training, we only fine-tune the injected adapters [15] in visual CLIP and LLaMa [26], along with the multi-modal projection module, while freezing the

major parameters. This aims to reserve the powerful abilities of existing MLLMs and further empower the model with capabilities in manipulation.

3.1.2 Fine-tuning Tasks Formulation

We design a training paradigm to fine-tune the MLLM and stimulate the model to generate interpretable pose predictions for object-centric manipulation.

Object Category Identification (OCI): To successfully manipulate the object, the model needs to understand the category of the object it is facing, as objects of the same category share common geometric properties. As illustrated in the first prompt in Fig. 2, we formulate the prompt as “What is the category of the object in the image?”. It’s worth mentioning that the MLLMs have been trained on a diverse set of objects in the real world, making them highly capable of category identification and generalization. In contrast, the object categories in the simulator are very limited, with a maximum of 30 to 50 [24]. Updating the learning process in the simulator might lead to a loss of MLLMs’s powerful object category identification ability and robust generalization capability. Therefore, we do not update the model in this stage, and the goal instead is to provide a prior of category cognition for subsequent tasks, helping them extract category-specific manipulation features.

Affordance Prior Reasoning (APR): This stage aims to enable the model aware where of the object region can be manipulated. Affordance map considers the object geometric and indicates the probability of getting a moving distance if operating on certain pixels, reflecting where can act to manipulate the object. It can serve as a region-level affordance prior to enabling the model to have manipulation-aware localization ability. Inspired by Flowbot3D [6], we divide the action type of the object part into “REVOLUTE” and “PRISMATIC”, and collect the affordance map in the simulator accordingly. For the revolute part, we first find the axis of the movable object part and then enable a movement of this part along the axis. We obtain affordance map $\mathcal{A} \in \mathbb{R}^{H \times W}$ following Eq. 1:

$$\mathcal{A} = \frac{\mathcal{D}}{|\max(\mathcal{D}) - \min(\mathcal{D})|} \quad (1)$$

The distance map, denoted as $\mathcal{D} \in \mathbb{R}^{H \times W}$, calculates the Euclidean distance of 3D positions (corresponding to each pixels) before and after the movement. Through applying a normalization operation based on the maximum and minimum value in distance map \mathcal{D} , we obtain the affordance map $\mathcal{A} \in [0, 1]$, indicating the probability of actionability on pixel-level. For the prismatic part, *i.e.*, drawer, operating all points on the surface of the movable part can promote a movement. Therefore, the probability on the affordance map of the prismatic movable part are all equal to 1. We

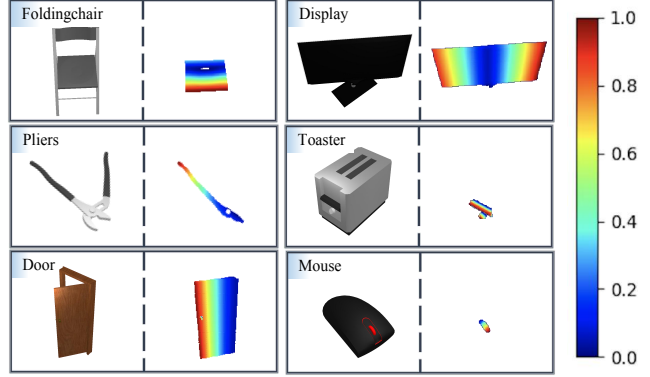


Figure 3. Affordance map for movable parts on objects. It indicates the probability of actionability on the pixel level.

visualize affordance maps in Fig. 3. For revolute parts, the affordance map reflects the regions where manipulation is possible, *i.e.*, regions away from the axis.

After we obtain the affordance map, our goal is to enable the model to learn from such manipulation prior. Since we only have a language decoder (LLaMa) instead of a visual decoder, the model is not able to generate an affordance map directly. Therefore, we aim to translate the visually represented affordance map to linguistic affordance prior. Specifically, we randomly select n positive pixels with an affordance score higher than 0.8 and select n negative pixels with an affordance score lower than 0.2 as training samples. The negative samples cover both the pixels on parts that cannot be moved and the pixels on parts that can be moved but have a low affordance score, *i.e.*, pixels close to the revolute axis. As shown in the second prompt in Fig. 2, we formulate the text prompt with the coordinates of the selected pixels as “Determine if operating on each following point can effectively manipulate the object within the image: $(x_P^1, y_P^1)..(x_P^n, y_P^n)..(x_N^1, y_N^1)..(x_N^n, y_N^n)$ ”, where P and N denote positive and negative samples. The corresponding ground truth answer is formulated as “yes, yes no, no...” with n “yes” and n “no” based on affordance scores. This is supervised under cross-entropy loss \mathcal{L}_A , enabling the model aware where of the object region can be manipulated and facilitating the model latter predict contact position that can promote a movement.

Finetuning (FT.) and Mask Language Modeling (MLM.): These tasks aim to enable the model to generate the precise end-effector pose. In the simulator, when pre-collecting training data, if the manipulation is successful, we record the RGB image and the corresponding end-effector pose, which are used as model input and answer ground truth. For task finetuning (FT.), as shown in the last prompt in Fig. 2, we design the input text prompt for pose prediction as “Specify the contact point and gripper direction of manipulating the object.” The answer is formulated as “The

contact point is (x, y) , the gripper up direction is (x_u, y_u, z_u) , and the gripper forward direction is (x_f, y_f, z_f) ". To decrease the difficulty of direction regression prediction, we transform it to classification prediction by discretizing the continuous numbers in the normalized direction vectors into 100 discrete bins $[-50, 50]$, with each bin spanning 0.02. The output is supervised under cross-entropy loss \mathcal{L}_F .

However, we found that directly fine-tuning the model for pose prediction leads to inaccuracies. Therefore, to facilitate the prediction of pose, in task Masked Language Modeling (MLM), we mask out the value of coordinate or direction vectors in the input text prompt and promote the model to infill the missing characters, as shown in the third prompt in Fig. 2. This is supervised by the unmasked answer under cross-entropy loss \mathcal{L}_M to stimulate the model's ability in pose prediction. The model learns to predict reasonable contact positions benefit from affordance prior learning. As for predicting the appropriate direction, we observe that MLLMs inherently possess direction awareness, such as being able to reason out "pull the door toward you". The training maps such direction cognitive descriptions and direction vectors to a consistent representation, enabling the prediction of the end-effector direction.

Training and Inference. During training, the aforementioned tasks are trained simultaneously under the total objective function: $\mathcal{L} = \mathcal{L}_A + \mathcal{L}_M + \mathcal{L}_F$. During inference, we adopt chain-of-thought reasoning to simulate the model to generate a precise initial contact end-effector pose interpretively. As shown in Fig. 4, the reasoning process follows the three steps that are consistent with the training tasks. The model finally outputs pixel coordinate (x, y) , gripper up direction (x_u, y_u, z_u) , and gripper forward direction (x_f, y_f, z_f) . We utilize the depth map $\mathbb{D}^{H \times W}$ to project the contact point into the 3D manipulation space (x, y, z) . The gripper up direction and gripper forward direction jointly formulate the end effector's rotation. Together with the predicted direction, they jointly determine the pose of the end-effector to establish the initial interaction with the object.

3.2. Active Impedance Adaptation Policy

After the initial interaction with the object is established, we apply a close-loop heuristic policy to adaptively generate upcoming waypoints under impedance control, *i.e.* the trajectory for opening the door. In the task of manipulating articulated objects, where we have limited freedom to move things around, it can be quite tough to figure out the best way to do it. For example, when trying to open a door, the best way to do it often involves moving it in a very specific direction along the axis of the door frame. To deal with these difficulties, the proposed policy aims to adjust how we interact with things based on impedance force feedback, which can handle different scenarios effectively. In con-

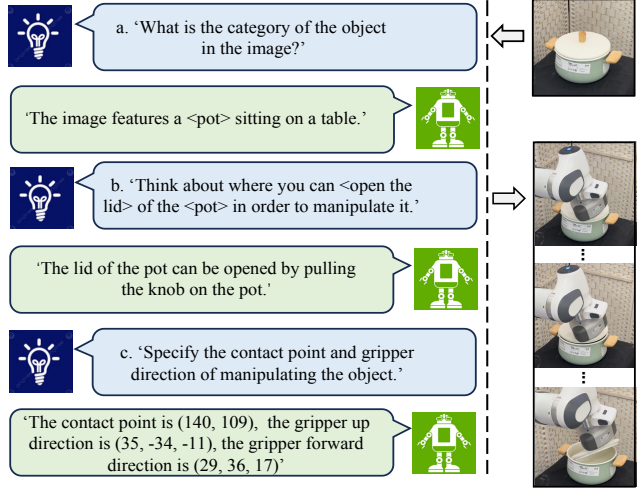


Figure 4. The chain-of-thought inference process of ManipLLM.

trast with leveraging a model to predict each following pose, such a heuristic policy is much more efficient.

This policy is employed within a loop to incrementally accomplish a manipulation task. At each iteration of the loop, it adaptively predicts the current direction that is best suited for the current state of the object, based on the previous step's forward direction. We use the initial iteration as an example. Given the predicted forward direction $d_i = (x_f, y_f, z_f), i = 1$, a random perturbation ζ is introduced such that $\|\zeta\| < \epsilon_1$, where ϵ_1 denotes a small positive constant. This procedure is reiterated N times for d_i , resulting in a set of directions represented as $D = \{d_{ij} = d_i + \zeta_j\}_{j \in \{0, 1, 2, \dots, N\}}$, where d_{i0} represents d_i , $\zeta_0 = \mathbf{0}$. Utilizing impedance control, a force f_j (f_j has a direction defined by d_{ij} and $\|f_j\| = \epsilon_2$, where ϵ_2 denotes another small positive constant) is applied to each direction d_{ij} within D . The optimal direction d_{opt} is then determined based on the observed end-effector movements δ_j . We assume that in constrained object manipulation tasks, greater movements represent the efficacy of the applied force direction. Thus, the best forward direction is generated as the following to determine the current end-effector's pose:

$$d_{opt, opt} = \arg \max_{j \in \{0, 1, \dots, N\}} \|\delta_j\|$$

By doing so, we determine the optimal movement pose given the current object state by considering the force feedback along the axis and ensuring a smooth trajectory.

3.3. Sim-to-real Transfer

Though the model is trained to perform well in the simulator, in the real world, robots often encounter more unique situations, *i.e.*, environment, or hard device configuration, which may differ significantly from what simulators simulate, leading to a sim-to-real gap. To bridge this gap, we

design a Test-Time Adaptation (TTA) strategy tailored for manipulation. TTA, as described by [21, 35], involves updating partial model parameters during inference based on the current test sample, enhancing the model’s performance for specific real-world scenarios. To determine which parameters to update for pose prediction during TTA, we analyze the outcomes of the reasoning steps during inference in Fig. 4. We observe that the reasoning abilities of ManipLLM, benefit from LLaMa [26], continue to exhibit strong performance in real-world scenarios. It can accurately recognize objects depicted in images and comprehend how to manipulate them. Its orientation awareness is also robust, ensuring the robustness of ManipLLM’s direction predictions. Even though there might be imprecise directions, with the active impedance adaptation policy introduced in Sec. 3.2, we can adjust the direction to a more optimal state. In contrast, position predictions are susceptible to domain gaps caused by factors like lighting and texture. Consequently, we adjust the visual perception for the target domain during TTA by only updating the V-Adapter in Fig. 2.

Specifically, given the current test sample, we introduce an additional reasoning step to prompt the model to assess whether the predicted position can lead to a successful manipulation. The text prompt used in this step is consistent with the training phase of “Affordance Prior Reasoning”, which is “Determine if operating on the following point can effectively manipulate the object within the image: (x, y).” The contact position predicted by the model is the region that they believe can lead to a successful manipulation. Therefore, the responses to this question are consistently “yes.” We obtain the ground-truth result based on whether, in the real world, the object was successfully manipulated, forming either a “yes” or “no” as the supervision signal to supervise the previous answer. By implementing this process, we enable the model to distinguish between effective and ineffective predictions in the target domain. This adjustment allows the model to predict valid poses when facing subsequent test samples, thereby adapting to the specific real-world configuration.

4. Experiment Results

4.1. Training Details

Data Collection. We adopt SAPIEN [31] and the PartNet-Mobility dataset to set up an interactive environment for our task, with VulkanRenderer of high-efficiency rasterization-based renderer. We use a Franka Panda Robot with flying suction gripper as the robot actuator. We sample the training data offline with approximately 10,000 manipulation success samples across 20 categories. We randomly select a contact point p on the movable part and use the opposite direction of its normal vector as the end-effector orientation for interacting with the object. If successful manipulation

is achieved, we record it as a successful sample. The tasks involve several pulling action primitives, where the movement direction of the object part and the opposite of end effector direction are within the same hemisphere, i.e., open the drawer, open the door, rotate the pliers, lift the lid, etc.

Training Details. We finetuned LLaMA-Adapter [38] on a 40G A100 GPU for 10 epochs, with an epoch costs around an hour. It includes pre-trained CLIP [25] as the visual encoder, 7B LLaMA [26] model as the decoder, and multi-modal projection module of 32 transformer layers. The n in Affordance Prior Reasoning is set to 20.

Evaluation Metric. We adopt the manipulation success rate to reflect the outcome of the manipulation which is the ratio of the number of successfully manipulated samples divided by the total number of all test samples. As for the definition of success sample, we adopt the binary success definition which is measured by thresholding the move distance of the object part at δ : $\text{success} = 1(\delta_{dis} > \delta)$. We set $\delta = 0.01$ or $\delta = 0.1$ for initial movement or long-distance movement, respectively, meaning that the gap between the start and end part 1-DoF pose is greater than 0.01 or 0.1 unit length. Initial movement is a prerequisite for long-distance movement and can effectively reflect the model’s ability to predict end-effector pose. Both initial and long-distance movements apply active impedance adaptation policy to adjust movement direction.

4.2. Quantitative Comparison

We compare ManipLLM against four representative baselines, including Where2Act [23], UMPNet [33], Flowbot3D [6], and Implicit3D [39]. For simplicity, we conduct the comparisons with other approaches only on initial movement setting since this can reflect how well the model can perform given the initial state of the object, which is the preliminary condition in realizing the whole long-distance movement. For a fair comparison, all methods are under the same train/test split and end-effector setting.

Where2Act [23]: It takes point-cloud as input and estimates per-point scores, selecting the point with the highest score as contact point. It further predicts 100 end-effector orientations and selects the orientation with the highest score to formulate the contact pose. For fair comparison, we alter the used parallel gripper to suction gripper.

UMPNet [33]: Following UMPNet, we execute manipulation on the contact point that UMPNet predicts with the orientation perpendicular to the object surface.

Flowbot3D [6]: It predicts motion direction on the point cloud, denoting it as ‘flow’. The point with the largest flow magnitude serves as the interaction point, while the direction of the flow represents end-effector’s orientation.

Implicit3D [39]: It develops a manipulation policy for downstream tasks that utilizes the Transporter to detect keypoints for 3D articulated objects. The keypoints are then

Method	Train Categories															
Where2Act [23]	0.26	0.36	0.19	0.27	0.23	0.11	0.15	0.47	0.14	0.24	0.13	0.12	0.56	0.68	0.07	0.40
UMPNet [33]	0.46	0.43	0.15	0.28	0.54	0.32	0.28	0.56	0.44	0.40	0.10	0.23	0.18	0.54	0.20	0.42
FlowBot3D [6]	0.67	0.55	0.20	0.32	0.27	0.31	0.61	0.68	0.15	0.28	0.36	0.18	0.21	0.70	0.18	0.26
Implicit3D [39]	0.53	0.58	0.35	0.55	0.28	0.66	0.58	0.51	0.52	0.57	0.45	0.34	0.41	0.54	0.39	0.43
Ours	0.68	0.64	0.36	0.77	0.43	0.62	0.65	0.61	0.65	0.52	0.53	0.40	0.64	0.71	0.60	0.64
Ours (long)	0.68	0.62	0.28	0.76	0.43	0.62	0.65	0.61	0.61	0.45	0.43	0.38	0.62	0.71	0.60	0.63

Method	Train Categories					Test Categories										
					AVG											
Where2Act [23]	0.13	0.18	0.13	0.40	0.26	0.18	0.35	0.38	0.28	0.05	0.21	0.17	0.20	0.15	0.15	0.21
UMPNet [33]	0.22	0.33	0.26	0.64	0.35	0.42	0.20	0.35	0.42	0.29	0.20	0.26	0.28	0.25	0.15	0.28
FlowBot3D [6]	0.17	0.53	0.29	0.42	0.37	0.23	0.10	0.60	0.39	0.27	0.42	0.28	0.51	0.13	0.23	0.32
Implicit3D [39]	0.27	0.65	0.20	0.33	0.46	0.45	0.17	0.80	0.53	0.15	0.69	0.41	0.31	0.30	0.31	0.41
Ours	0.41	0.75	0.44	0.67	0.56	0.38	0.22	0.81	0.86	0.38	0.85	0.42	0.83	0.26	0.38	0.51
Ours (long)	0.37	0.75	0.44	0.67	0.54	0.34	0.22	0.81	0.86	0.30	0.85	0.42	0.80	0.26	0.38	0.50

Table 1. Comparisons of our method against baseline methods.

used to determine end-effector pose.

Our current experimental settings involve training on a wider range of object categories. Consequently, this poses challenges in extracting features and learning characteristics from these wide categories, which may lead to a decrease in manipulation success rate compared to the original papers. However, our method, by retaining the common-sense reasoning capabilities of MLLMs and injecting manipulation abilities, ensures strong generalization across diverse categories. It’s worth noting that other methods also require a movable mask during testing, whereas our method achieves this **without the need for movable part mask**.

Voxposer [16]: We also compare our method with VoxPoser for the “Opening Drawer” task. The success rate of VoxPoser is 14.0% while ours is 69.0%. We found that VoxPoser struggles to find the correct grasping pose and a suitable moving trajectory for challenging cases.

4.3. Ablation and Analysis

To elucidate the contribution and effectiveness of individual modules within our approach, we conduct extensive ablation studies. The results are measured on novel instances in train categories with initial movement.

Effectiveness of tasks in the training paradigm. In Table 2, we gradually add each task in the training paradigm to show the effectiveness of each. *Finetuning(FT.):* In the first row of Table 2 with only fine-tuning, the last prompt in Fig. 2, we found that introducing this single task allows the model to have some manipulation capability, showcasing the strong learning ability of the large model. *Object Category Identification(OCI.):* Subsequently, in the second row of Table 2, we introduce the task of object category identi-

Train				Test		AVG
FT.	OCI.	MLM.	APR.	COT.	AIA.	
✓	-	-	-	✓	✓	0.41
✓	✓	-	-	✓	✓	0.44
✓	✓	✓	-	✓	✓	0.49
✓	✓	✓	✓	✓	✓	0.56
✓	✓	✓	✓	-	✓	0.54
✓	✓	✓	✓	✓	-	0.53

Table 2. Ablation analysis of each training task in the training paradigm and strategies in inference.

fication, the first prompt in Fig. 2. It enables the model to discover commonalities in manipulating objects of the same or similar categories, thus enhancing the model’s manipulation abilities by +3%. *Mask Language Modeling(MLM.):* Next, in the third row of Table 2, also the third prompt in Fig. 2, we randomly mask values in the coordinates or direction vectors to force the model to predict precise pose. This task stimulates the model’s ability of localization that it previously lacked and enables the model to map the direction common sense reasoning with $SO(3)$ direction representation, thus improving the performance by +5%. *Affordance Prior Reasoning(APR.):* Finally, in the fourth row of Table 2, we introduce the affordance prior reasoning task, the second prompt in Fig. 2, allowing the model to learn manipulation-aware localization and predict accurate contact positions. It thus significantly improves the manipulation success rate by +7%.

Effectiveness of strategies in inference. *Chain-Of-Thought Reasoning (COT.):* The COT reasoning strategy

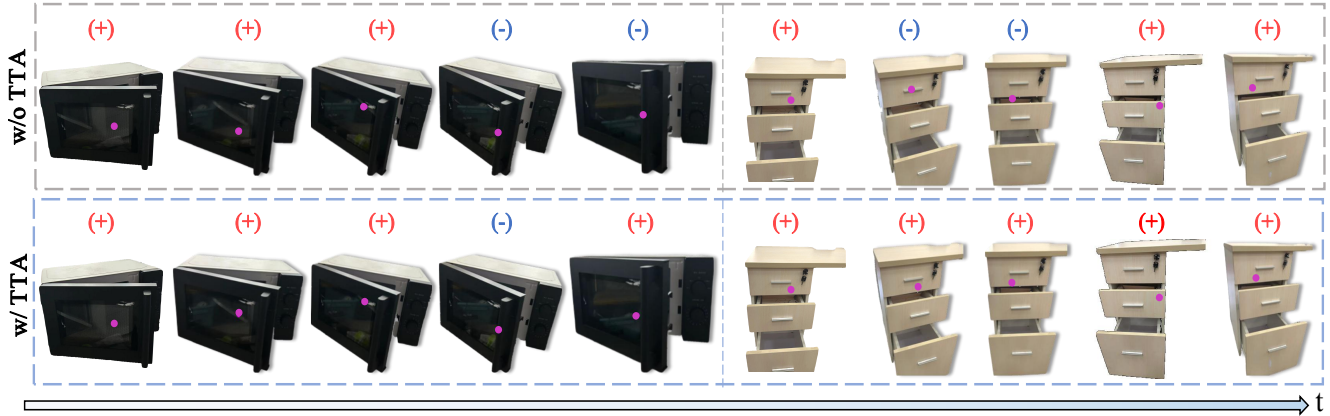


Figure 5. Visualizations of TTA process in real-world scenarios. The center of pink dot represents the predicted contact position.

aims to enable the model to generate the final pose prediction under a transparent and reasonable thinking process. For comparison, we ask the model to generate the final pose prediction directly without the thinking process in Fig. 4. As shown in the second-to-last row of Table 2 w/o COT, we find that this decreases the performance by -2% compared with applying COT during inference. This emphasizes the essential of enabling the model to predict under a transparent and interpretable process.

Active Impedance Adaptation (AIA.): The active impedance adaptation policy adaptively adjusts the pose to adapt the current object state under impedance control. In the last row of Table 2 w/o AIA., we employ a straightforward control policy, which operates by moving directly to the desired position under invariant direction. In contrast, Active Impedance Adaptation policy applies impedance control to effectively adjust the direction based on force feedback, thus enabling a smooth manipulation trajectory in the long-term movement and improving the long-distance movement by -3% . This policy is particularly crucial for long-distance movement.

4.4. Real-world Evaluation

We conduct experiments that involve interacting with various real-world household objects. We employ a Franka Emika robotic arm with cobot pump suction gripper and utilize a RealSense 415 camera to capture RGB image and depth map. To address the sim-to-real problem: 1) during training, we utilize the LLAMA-Adapter pretrained in the real-world and employ a method that combines injection and finetuning of the adapter to make it learn new downstream tasks. This training strategy allows the model to retain its strong perception abilities in the real-world while equipping it with the capability to perform manipulation. 2) When collecting data in the simulator, we employ domain randomization to increase scenario diversity by varying elements such as object part poses, camera view angles,

and lighting, among others, in order to mitigate the potential sim-to-real gap. 3) During testing, we design test-time adaptation strategy to help the model better adapt to the current scene’s configuration. By finetuning only the adapters, we shift its attention, enabling MLLMs to generate predictions that are better suited to the current scene.

The results of real-world experiments are shown in Table 3. As illustrated in Fig. 5, the devised TTA strategy addresses discrepancies arising from real-world hardware configurations. In our specific hardware configuration, the suction gripper is unable to grasp the handle due to the non-smooth surface. Additionally, its head is relatively short, which presents a collision risk when interacting with the protruding handle. The TTA process learns from both successful and unsuccessful scenarios, gradually adapting its predictions to align with the current configuration.





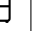
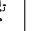

Object Category							
Success/Total	4/5	5/5	4/5	3/5	4/5	4/5	4/5
Distance(m)	0.17	0.28	0.10	0.08	0.14	0.15	0.18

Table 3. Real world experiments.

5. Conclusion

We transform MLLMs to robotic manipulation through the chain-of-thought training paradigm and equip the model with the ability to predict poses. We introduce an active impedance adaptation policy that adjusts direction based on force feedback of the correct state to ensure a smooth moving trajectory. ManipLLM shows strong generalization ability across extensive categories and in real-world.

Acknowledgements. The project was supported by the National Youth Talent Support Program (8200800081) and National Natural Science Foundation of China (No. 62136001).

References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. 1, 3
- [2] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6077–6086, 2018. 1
- [3] OpenAI: Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020. 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. 2
- [5] Congyue Deng, Jiahui Lei, Bokui Shen, Kostas Daniilidis, and Leonidas Guibas. Banana: Banach fixed-point network for pointcloud segmentation with inter-part equivariance. *arXiv preprint arXiv:2305.16314*, 2023. 2
- [6] Ben Eisner, Harry Zhang, and David Held. Flowbot3d: Learning 3d articulation flow to manipulate articulated objects. *arXiv preprint arXiv:2205.04382*, 2022. 2, 3, 4, 6, 7
- [7] Luciano Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30: 681–694, 2020. 3
- [8] Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. *arXiv preprint arXiv:1606.01847*, 2016. 1
- [9] Haoran Geng, Ziming Li, Yiran Geng, Jiayi Chen, Hao Dong, and He Wang. Partmanip: Learning cross-category generalizable part manipulation policy from point cloud observations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2978–2988, 2023. 1
- [10] Haoran Geng, Songlin Wei, Congyue Deng, Bokui Shen, He Wang, and Leonidas Guibas. Sage: Bridging semantic and actionable parts for generalizable articulated-object manipulation under language instructions. *arXiv preprint arXiv:2312.01307*, 2023. 2
- [11] Haoran Geng, Helin Xu, Chengyang Zhao, Chao Xu, Li Yi, Siyuan Huang, and He Wang. Gapartnet: Cross-category domain-generalizable object perception and manipulation via generalizable and actionable parts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7081–7091, 2023. 1
- [12] Yiran Geng, Boshi An, Haoran Geng, Yuanpei Chen, Yaodong Yang, and Hao Dong. End-to-end affordance learning for robotic manipulation. In *International Conference on Robotics and Automation (ICRA)*, 2023. 2
- [13] Yiran Geng, Boshi An, Haoran Geng, Yuanpei Chen, Yaodong Yang, and Hao Dong. Rlafford: End-to-end affordance learning for robotic manipulation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5880–5886. IEEE, 2023. 1
- [14] Ran Gong, Jianguo Huang, Yizhou Zhao, Haoran Geng, Xiaofeng Gao, Qingyang Wu, Wensi Ai, Ziheng Zhou, Demetri Terzopoulos, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. Arnold: A benchmark for language-grounded task learning with continuous states in realistic 3d scenes. *arXiv preprint arXiv:2304.04321*, 2023. 2
- [15] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2, 3
- [16] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. *arXiv preprint arXiv:2307.05973*, 2023. 2, 3, 7
- [17] Shirin Joshi, Sulabh Kumra, and Ferat Sahin. Robotic grasping using deep reinforcement learning. In *2020 IEEE 16th International Conference on Automation Science and Engineering (CASE)*, pages 1461–1466. IEEE, 2020. 2
- [18] Jiahui Lei, Congyue Deng, Bokui Shen, Leonidas Guibas, and Kostas Daniilidis. Nap: Neural 3d articulation prior. *arXiv preprint arXiv:2305.16315*, 2023. 2
- [19] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. 1, 3
- [20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 3
- [21] Jiaming Liu, Senqiao Yang, Peidong Jia, Ming Lu, Yandong Guo, Wei Xue, and Shanghang Zhang. Vida: Homeostatic visual domain adapter for continual test time adaptation. *arXiv preprint arXiv:2306.04344*, 2023. 2, 6
- [22] Jiageng Mao, Yuxi Qian, Hang Zhao, and Yue Wang. Gpt-driver: Learning to drive with gpt. *arXiv preprint arXiv:2310.01415*, 2023. 1
- [23] Kaichun Mo, Shilin Zhu, Angel X. Chang, Li Yi, Subarna Tripathi, Leonidas J. Guibas, and Hao Su. PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 1, 2, 6, 7
- [24] Kaichun Mo, Shilin Zhu, Angel X Chang, Li Yi, Subarna Tripathi, Leonidas J Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 909–918, 2019. 4
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,

- Amanda Aspell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 3, 6
- [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, 3, 6
- [27] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. Unidexgrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. *arXiv preprint arXiv:2304.00464*, 2023. 2
- [28] Qianxu Wang, Haotong Zhang, Congyue Deng, Yang You, Hao Dong, Yixin Zhu, and Leonidas Guibas. Sparsedff: Sparse-view feature distillation for one-shot dexterous manipulation. *arXiv preprint arXiv:2310.16838*, 2023. 2
- [29] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022. 2
- [30] Ruihai Wu, Yan Zhao, Kaichun Mo, Zizheng Guo, Yian Wang, Tianhao Wu, Qingnan Fan, Xuelin Chen, Leonidas Guibas, and Hao Dong. Vat-mart: Learning visual action trajectory proposals for manipulating 3d articulated objects. *arXiv preprint arXiv:2106.14440*, 2021. 2
- [31] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su. SAPIEN: A simulated part-based interactive environment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 6
- [32] Yinzhen Xu, Weikang Wan, Jialiang Zhang, Haoran Liu, Zikang Shan, Hao Shen, Ruicheng Wang, Haoran Geng, Yijia Weng, Jiayi Chen, et al. Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation and goal-conditioned policy. *arXiv preprint arXiv:2303.00938*, 2023. 2
- [33] Zhenjia Xu, Zhanpeng He, and Shuran Song. Universal manipulation policy network for articulated objects. *IEEE Robotics and Automation Letters*, 7(2):2447–2454, 2022. 6, 7
- [34] Jingyun Yang, Congyue Deng, Jimmy Wu, Rika Antonova, Leonidas Guibas, and Jeannette Bohg. Equivact: Sim(3)-equivariant visuomotor policies beyond rigid object manipulation. *arXiv preprint arXiv:2310.16050*, 2023. 2
- [35] Senqiao Yang, Jiarui Wu, Jiaming Liu, Xiaoqi Li, Qizhe Zhang, Mingjie Pan, and Shanghang Zhang. Exploring sparse visual prompt for cross-domain semantic segmentation. *arXiv preprint arXiv:2303.09792*, 2023. 2, 6
- [36] Denis Yarats, Rob Fergus, Alessandro Lazaric, and Lerrel Pinto. Mastering visual continuous control: Improved data-augmented reinforcement learning. *arXiv preprint arXiv:2107.09645*, 2021. 2
- [37] Yang You, Bokui Shen, Congyue Deng, Haoran Geng, He Wang, and Leonidas Guibas. Make a donut: Language-guided hierarchical emd-space planning for zero-shot deformable object manipulation, 2023. 3
- [38] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023. 1, 2, 3, 6
- [39] Chengliang Zhong, Yuhang Zheng, Yupeng Zheng, Hao Zhao, Li Yi, Xiaodong Mu, Ling Wang, Pengfei Li, Guyue Zhou, Chao Yang, et al. 3d implicit transporter for temporally consistent keypoint discovery. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3869–3880, 2023. 6, 7
- [40] 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 *7th Annual Conference on Robot Learning*, 2023. 2, 3