

Transferring Foundation Models for Generalizable Robotic Manipulation

Jiange Yang¹, Wenhui Tan², Chuhao Jin², Keling Yao³, Bei Liu⁴,
Jianlong Fu⁴, Ruihua Song², Gangshan Wu¹, Limin Wang^{1,5*}

¹State Key Laboratory for Novel Software Technology, Nanjing University, China

²Renmin University of China

³The Chinese University of Hong Kong, Shenzhen

⁴Microsoft Research ⁵Shanghai AI Lab

jianggeyang.jgy@gmail.com, {tanwenhui404, jinchuhao, rsong}@ruc.edu.cn,
120090220@link.cuhk.edu.cn, {jianf, Bei.Liu}@microsoft.com, {gsu, lmwang}@nju.edu.cn

Abstract

Improving the generalization capabilities of general-purpose robotic manipulation in real world has long been a significant challenge. Existing approaches often rely on collecting large-scale robotic data which is costly and time-consuming. However, due to insufficient diversity of data, they typically suffer from limiting their capability in open-domain scenarios with new objects and diverse environments. In this paper, we propose a novel paradigm that effectively leverages language-reasoning segmentation mask generated by internet-scale foundation models, to condition robot manipulation tasks. By integrating the mask modality, which incorporates semantic, geometric, and temporal correlation priors derived from vision foundation models, into the end-to-end policy model, our approach can effectively and robustly perceive object pose and enable sample-efficient generalization learning, including new object instances, semantic categories, and unseen backgrounds. We first introduce a series of foundation models to ground natural language demands across multiple tasks. Secondly, we develop a two-stream 2D policy model based on imitation learning, which processes raw images and object masks to predict robot actions with a local-global perception manner. Extensive real-world experiments conducted on a Franka Emika robot and a low-cost dual-arm robot demonstrate the effectiveness of our proposed paradigm and policy. Demos can be found in [link1](#) or [link2](#) and our code will be released at <https://github.com/MCG-NJU/TPM>.

1. Introduction

Creating a general-purpose robotic agent that is capable of performing diverse tasks in real-world remains a long-

standing and challenging research. In this paper, we endeavor to develop a model that enables generalizable robotic manipulation. The first challenge in creating generalizable agents is effectively converting abstract task condition instructions into specific robot inputs. Current approaches utilize various forms, including task identifiers [59], goal images [41], video showcasing of human demonstrations [26], and natural language [4, 18, 26, 42, 62]. Language, in particular, provides a natural and scalable manner for human-robot interaction, but may be under-specified and ambiguous. The second challenge involves enhancing the generalization capabilities of a single robot model to handle multiple tasks, encompassing both unseen objects and environments. To address the above challenges, recent advancements [4, 25, 26] have predominantly embraced data-driven learning-based methods. Notably, one of the groundbreaking work is RT-1 [4], which has introduced a comprehensive model capable of executing diverse instructions using an extensive dataset of approximately 130,000 demonstrations spanning over 17 months and involving 13 robots. However, the collection of real-world data poses significant resource requirements, and the approach exhibits limitation to compositional generalization, struggling with unseen objects and environments, due to insufficient diversity of data.

In this paper, we propose to achieve sample-efficient generalization for robotic manipulation by introducing language-reasoning mask modality containing semantics, geometry, and temporal correlation priors inherent from internet-scale vision foundation models into an end-to-end policy model. Specifically, segmentation has been proven significant as grasping priors in manipulation tasks [32, 38, 77], and we introduce language-reasoning mask as a new condition modality for policy model and conduct end-to-end training using imitation learning. An intuitive pipeline for robot manipulation is to first achieve visual perception

*Corresponding author.



Figure 1. A demonstration of our task. Receiving human instruction “I want to take a shower”, our model can reason out the desired object (i.e., the towel), and then precisely pick and place it near the target object (i.e., the user represented by a Lego toy).

and then conduct motion planning like [9, 33, 64]. However, such pipeline necessitates an efficient motion planner, while also requiring completely accurate object mask and depth for constructing point clouds, which may not perform well when dealing with transparent and disturbed objects, as well as unstructured environments with collision situation. In contrast, the end-to-end 2D policy pipeline we adopt in this work, as well as other works [4, 26], can dynamically receive raw image as input and output continuous action in a close-loop manner, which does not rely on depth calibration and completely accurate object masks. Furthermore, *our paradigm aims to fully unify the generalizability of internet-scale models and the potential of imitation learning to capture multimodal action distribution from human skills at the lowest possible training and inference cost, while also mitigating the ambiguity of language as condition. Therefore, this paradigm is highly scalable.*

In order to further build a holistic robot system with human-robot interaction, we utilize large language model to reason human demands and design a two-stream policy model to predict actions with a local-global perception manner. Specifically, we first use GPT-4 [52] to interpret language instructions and generate desired object prompts. Second, we identify and locate objects by open-vocabulary detection [37] and tracking [7] models. We then adopt the vision foundation model SAM [31] to generate segmentation masks of desired objects. Subsequently, we propose a **Two-stream architecture Policy Model, TPM**, which uses a deeper branch to capture global RGB information and a shallower branch to capture local object-related RGB-M information as well as fuses multi-view features and robot proprioception states through attention mechanism. These designs enable more robust 3D perception and thus leading to accurate action prediction. Finally, to verify our system, we primarily select the widely popular pick-and-place (picking A and placing near B) tasks for quantitative evaluation. Specifically, we collect a dataset consisting of 1000 demonstrations with 40 objects for training. Our experiments results show the effectiveness of our proposed paradigm and policy model architecture, particularly in gen-

eralizing to unseen objects, complex backgrounds, and multiple distractors. A simple demonstration of our task is shown in Figure 1. In addition, we conduct further experiments to show the capabilities of our method to transfer more manipulation skills, including opening drawer, picking A and placing inside of B, placing on top of B, folding and stacking.

In summary, the contributions of our paper are summarized as follows:

- We propose a novel paradigm by transferring internet-scale foundation models for robotic manipulation with language-reasoning segmentation mask, aiming to enhance its generalization capabilities in a sample-efficient way. Additionally, the mask modality also provide a more specified and unambiguous condition representation than unprocessed human language.
- We develop a two-stream policy model for handling images and language-reasoning masks with a local-global perception manner, which achieves better spatial relationship understanding.
- Extensive real-world experiment results demonstrate that our paradigm and policy model architecture can effectively improve the performance and generalize to handle unseen objects, new backgrounds, more distractors, and even expand to more manipulation skills.

2. Related Work

Pre-trained Foundation Models for Robotics. Recently several works [21–25, 33, 34, 36, 64, 69] explored to leverage off-the-shelf large language models to plan feasible tasks for robotics. Some works [3, 12, 27, 50, 53, 56] further use in-domain data to fine-tune LLMs for embodied reasoning. InstructRL [35] employs masked autoencoder M3AE [15] to encode visual observations and language instructions. PAFF [14] utilizes CLIP [57] to provide feedback for relabeling demonstrations. MOO [65] leverages pre-trained vision-language model to extract object-centric representations, which is based on a single pixel

in the center of the bounding box on the first frame. In contrast, our approach distinctively employs more specified object masks using a detection-tracking-segmentation manner, which not only enhances the precision and reliability of object representation for robotic manipulation tasks, but also provides stronger geometric and temporal correlation priors. In addition, different from [22, 33, 64], which use prompt-based segmentation mask and depth to construct object point cloud and then call for additional grasp and motion planner, *we directly utilize the mask conditioning the end-to-end and learnable policy model with closed-loop manner, without requiring depth, camera calibration and completely accurate object masks.*

Vision-based Robot Learning. Vision-based robot learning plays a crucial role in robotics. Recently pre-trained visual models for robotics have been rapidly developing [29, 43, 44, 46, 51, 54, 58, 75, 81]. R3M [51] and VC-1 [46] demonstrate the effectiveness of pre-training on egocentric videos. STP [75] further considers temporal motion cues. [44] proposes a vision model capable of producing dense reward signals and LIV [43] expand it to multi-modal. Li et al. show that pre-training on semantic tasks like classification and segmentation helps in improving efficiency and generalization of grasping [77]. Some works [1, 2, 61, 70, 71, 73, 76] also focus on learning skill priors from large-scale human video data. As for model architecture, several works [20, 28] point out that convolution-based models tend to outperform transformer-based models and [6, 79] demonstrates the superiority of generative modeling of the policy. Our work is orthogonal to the contributions of these works. In addition, some works [13, 45, 66, 67] focus on point cloud-based grasping pose generation. In contrast, our method requires dynamically mapping observations to continuous actions with closed-loop manner. With different problem formulations, direct comparisons are usually not performed.

Language-Conditioned Robotics Control. The goal of building a robotic model that can follow diverse natural language instructions has consistently been an active research field [4, 16–18, 26, 42, 47–49, 62, 63]. Noteworthy advancements have been made by various approaches. Hiveformer [18] proposes a unified approach to encode the full history of observation-action pairs. Perceive-actor [63] encodes RGB-D voxel observations with a Perceiver Transformer to provide a strong structural prior. BC-Z [26] and RT-1 [4] focus on scaling and expanding the collection of real-world data to facilitate generalization of robots, which encompasses unseen tasks, environments, and objects. However, due to the limitations in both quantity and diversity of the collected data compared to large-scale datasets in the vision and language domains [10, 60, 74], the generalization ability of the trained models is still relatively poor, especially in unseen object categories.

3. Approach

In this section, we begin by presenting the problem formulation of our method in Section 3.1. Following this, we detail the pipeline of language-reasoning mask generation with foundation models in Sections 3.2. Finally, we elaborate on our two-stream policy model and its training methodology in Section 3.3. We primarily illustrate our approach using pick-and-place skill, and the same manner can be applied to other skills as well.

3.1. Problem formulation

We aim to develop a generalizable robot capable of interpreting high-level language instructions from users and executing precise actions to fulfill their needs. For instance, when a user says, “I am thirsty”, the robot should pick the milks up from the tabletop, and place it near the user’s hands. This process seamlessly integrates perception, reasoning, and control into a unified pipeline.

In our problem setting, we aim to develop a robotic system, $\Phi_\varepsilon(a|p, o, l)$, parameterized by ε , which maps robot proprioception p , visual observations o , and human-provided language instructions l to continuous actions a on a physical robot. To enhance the generalization capabilities and sample efficiency, we divide our system into two parts: (1) A segmentation mask generator based on foundation models discussed in Section 3.2, which fully leverages the potential of internet-based foundation models to obtain the desired object masks based on human instructions. (2) A two-stream policy model described in Section 3.3, which encodes multi-modal inputs, including raw images, language-reasoning object masks, and robot proprioception, mapping them to specific robot actions.

3.2. Mask Generation with Foundation Models

To enhance the generalization capabilities of robot agents, we propose a novel paradigm that utilizes a series of internet-scale foundation models to accurately interpret abstract natural language instructions, locate objects, and segment their geometry mask for subsequent specified task condition. Specifically, this paradigm comprises target object reasoning based on LLMs, a multi-modal prompt generator based on open-vocabulary detector and tracker, and segmentation mask generation based SAM. The framework of our robotic agent system is provided in Figure 2.

Target objects reasoning. We employ an LLM, GPT-4 [52], as the central cognitive component of our model. When a user submits a high-level language instruction (e.g., “I spilled the milk”), it is integrated into a designed prompt template, which is subsequently fed into the GPT-4 to reason out the target objects for task condition (e.g., pick a sponge and place it near the user). This process enables our model to effectively deduce the target objects that fulfills the user’s requirements, ensuring the correct task condition.

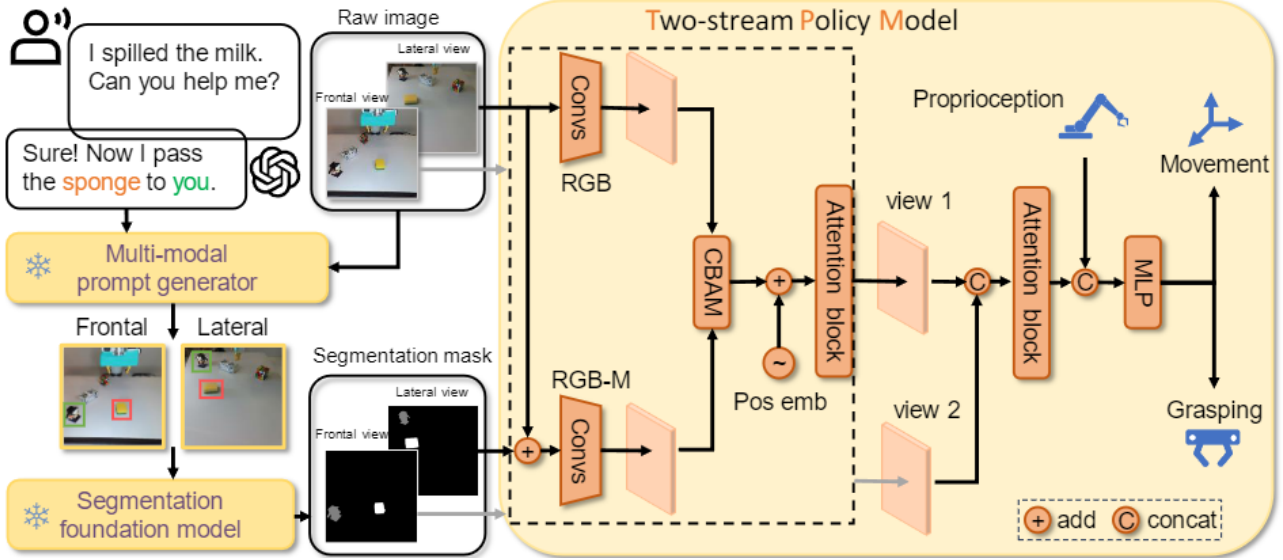


Figure 2. Our model comprises four components: (1) GPT-4 reasons target objects based on human demands. (2) A multi-modal prompt generator, comprising object detection and tracking models, transforming input images and target object prompts into bounding boxes. (3) The Segment Anything model, which uses bounding boxes as prompts to segment target objects and generate task-relevant masks. (4) A two-stream policy model that processes images, language-reasoning segmentation masks, and robot proprioception to predict actions.

Multi-modal prompt generator. After reasoning out which targets should be interacted, we apply a state-of-the-art open-vocabulary object detector Grounding DINO [37] with powerful semantic concept generalization ability to locate the target objects given language expression. However, during the process of a robotic agent performing a task, the target to be manipulated inevitably encounters occlusion (e.g., being blocked by the robotic arm), disturbance (e.g., dynamic objects) and potential distractors. Compared to object detection, object tracking is more robust in handling these challenging scenarios due to its inherent spatio-temporal correlation. Therefore, after the robotic agent completes the first step of the action, we switch to a state-of-the-art and efficient tracker MixFormer [7] to obtain all subsequent bounding boxes. To conclude, we obtain the bounding box by a object detector in the first step, and tracks it after this, which consequently serve as prompts for object mask generation.

Mask generation based on SAM. The image segmentation model SAM [31] has been widely explored in various fields due to its powerful object generalization capability and promptability. The appeal of SAM is enhanced by its ability to flexibly integrate the semantic concept generalization of open-vocabulary detectors with a bounding box prompt. In robot manipulation tasks, object masks are closely related to affordance maps because of their shape and geometry prior. Thus, after identifying and locating the target objects, we provide the bounding box as a prompt to SAM to generate the segmentation masks of them.

In our approach, we transform abstract language instructions into specified target object masks, which incorporates

rich semantic, geometry, and temporal correlation priors derived from vision foundation models into the end-to-end 2D policy model. By collecting a small amount of real robotic data, the policy model is able to efficiently learn how to adapt these capabilities onto specific manipulation actions, thereby achieving sample-efficient learning of generalizable manipulation. In summary, the process of target object masks generation can be formulated as follow:

$$m = \text{SAM}(o, \text{PG}(o, \text{RM}(o, l))), \quad (1)$$

where PG and RM denote the multi-modal prompt generator and the reasoning model, respectively. m , o and l stand for the generated object masks, raw images and language instructions, respectively.

3.3. Two-stream Policy Model Architecture

Policy model architecture. As shown in the right side of Figure 2, the two-stream Policy Model, which we refer to as **TPM**, maps the robot proprioception, object masks, as well as the raw image observations of two different views to continuous actions. The process can be formulated as

$$\pi_{\theta}(p, (o_1, m_1), (o_2, m_2)), \quad (2)$$

where the subscripts index different perspective of views, and TPM is parameterized by θ .

Inspired by the success of mask features for memory bank in video object segmentation [5], we concatenate the image RGB frame and object mask along the channel dimension to form RGB-M. To make it clear, we further explain details as follows: (1) we first adopt a two-stream

convolution-based architecture to separately encode RGB-M and RGB, where RGB-M branch employs a shallower ResNet-18 [19] network to capture local features of task-related objects, and RGB-only branch utilizes a deeper ResNet-50 network to capture the global spatial relationships within the entire scene. We both take stage-4 features with stride 16 as standard feature maps, and utilize a CBAM [72] block to fuse them from space and channel dimension. (2) In order to enhance the spatial perception capability, we add a 2D positional embedding for each feature point and subsequently apply an attention [68] block to promote spatial interaction. As for multi-view feature fusion, we concatenate two view features and also employ a global self-attention block to ensure dynamic spatial alignment. (3) Finally, for better integrating embodied perception to our policy model, we further inject the robot proprioception (pose of the end-effector) state. In our work, the proprioception embeddings and visual embeddings are fused via concatenation to obtain the final state representation, which is then fed into two MLPs, to generate the predicted next action, i.e., the movement along x, y and z-axis and the opening state of the gripper.

Policy model training. We train our TPM using behavior cloning with our collected offline dataset \mathcal{D}_1 . In overall, the loss function can be formulated as follow:

$$\min_{\theta} \sum_{\mathcal{D}} \text{CE}(a_g, \pi_{\theta}(p, o, m)) + \mu \cdot \text{MSE}(a_m, \pi_{\theta}(p, o, m)), \quad (3)$$

where a_g and a_m denote the gripper state and the movement of the robot’s end-effector. We adopt a deterministic policy manner. The model learns the gripper open or closed state as a binary classification task with cross-entropy loss (CE), and learns the continuous movement of end-effector through mean square error (MSE). The μ denotes the loss weight of the MSE loss.

4. Experiment

In this section, we first introduce our pick-and-place dataset. Then we elaborate on the implementation details, experimental setup and evaluation results in the following sections. Finally, we elucidate how our pipeline can be flexibly extended to other skills, and provide both qualitative and quantitative experiment results.

4.1. Dataset

To facilitate the experiments of our work, we collect a dataset using Franka Emika Research 3 robot arm to perform imitation learning, which contains 1000 episodes. For each episode, we annotate the language instruction and use GroundingDINO [37] and SAM [31] to obtain the language-conditioned masks. We select manipulated objects of various types, such as those with different shapes, sizes, textures, and colors, to ensure the diversity of the dataset. Specifically, we capture multi-view images from

both frontal and lateral perspectives. We select 40 common objects categorized into 5 typical shapes from daily life, in a total of 1000 pick-and-place trajectories across 3 different table-top backgrounds, as shown in Figure 3(c). For every pick-and-place demonstration, we randomly place 0 to 2 additional objects as distractors.

4.2. Implementation Details

In our evaluation experiments, all vision foundation models [7, 31, 37] employ the base version [11, 39]. To optimize our proposed TPM, we use Adam [30] with decoupled weight decay [40] of $5 \cdot 10^{-4}$. The peak learning rate is set to $5 \cdot 10^{-5}$, decaying with a cosine learning rate schedule to $5 \cdot 10^{-6}$. We empirically set the weight μ of movement loss to gripper state loss at 1,000 to keep them in comparable magnitude. We train on 224×224 images without data augmentation, with a batch size of 24 for 500k iterations. The model code is implemented in PyTorch [55] and trained on an NVIDIA RTX A6000 GPU.

4.3. Experimental Setup

Real-World Environment. In our real-world experiments, we use a Franka Emika Research 3 robot arm in a table-top environment, consistent with data collection. The Intel RealSense camera ($1,280 \times 720$ resolution) and Azure Kinect DK camera ($2,048 \times 1,536$ resolution) are mounted on fixed supports, as shown in Figure 3(a).

Evaluation Metric. Our quantitative experiments primarily focus on pick-and-place tasks described as “pick A and place it near B”. We measure the percentage of successful pick-and-place tasks as success rate. A successful pick-and-place is defined as (1) grasping A and (2) placing it within 2 inches of B.

Evaluation Setup. We define the 40 objects in the collected training data as seen objects and hold out another 20 objects not present in the training data as unseen objects. To comprehensively evaluate the generalization capabilities of our method, we evaluate our model from three aspects: (1) *Seen/Unseen objects*. We have two settings, one is with seen objects for both A and B and the other is with unseen objects for A and B. Unlike prior works [4] focusing on compositional generalization for unseen tasks, our unseen objects are strictly new categories or instances, testing the model’s ability to handle new objects. We present some objects in Figure 3(b). (2) *New background*. We introduce an environment with a complex-textured tablecloth, altering lighting, materials, and backgrounds, at the same time. This evaluates our model’s robustness to out-of-distribution generalization. We show the new background in Figure 3(d). (3) *More distractors*. We add a scenario with more than 2 distractor objects (ranging from 3 to 6 objects), creating a congested tabletop scene with target objects disturbance and occlusion. This evaluates the model’s robustness against



Figure 3. (a): Overview of our workstation, which has a Franka robot arm, a frontal view camera, and a lateral view camera. (b): Seen and unseen objects in the experiments. (c): Three backgrounds in the training data. (d): A challenging background with complex texture for new background evaluation.

| Scenario | Seen | Unseen | Average |
|------------------|------|--------|---------|
| Standard | 82.5 | 80.0 | 81.25 |
| New background | 65.0 | 55.0 | 60.0 |
| More distractors | 75.0 | 70.0 | 72.5 |

Table 1. Experimental results evaluated on different scenarios.
more distractor objects.

4.4. Experiment results

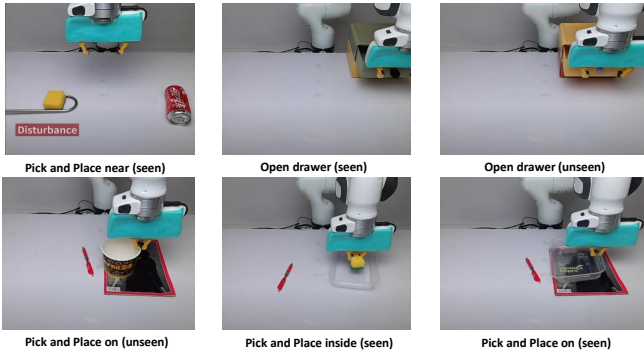


Figure 4. Some demonstration examples of disturbances scene and other manipulation skills.

We initially conduct a comprehensive evaluation of our method’s effectiveness and its ability to generalize to unseen objects, new backgrounds, and more distractors. To ensure a more extensive evaluation, we randomly select 20 tasks under each setting. For each task, the source object initial positions are placed once on the left and once on the right, resulting in a total of **40** trials. The *standard* environment refers to a tabletop background with 0-2 distractors, consistent with the training data. The experiment results are presented in Table 1. These results indicate slight performance decrease when introducing a new background. However, our model demonstrates robustness to a greater number of distractors and unseen objects, which ex-

ceeds the limit of RT-1 and can be attributed to the inclusion of language-reasoning segmentation mask modality derived from foundation models for action prediction.

To further verify the effectiveness of our proposed approach and its individual components, we compare our method with several variants. We evaluate each method in four settings: (1) seen objects in the standard environment; (2) unseen objects in the standard environment; (3) seen objects in a new background; and (4) seen objects with more distractors (randomly 3-6). Each setting includes 20 tasks, resulting in a total of **80** trials for each method. The comparative experiment results are presented in Table 2. Additionally, we also follow the setup in Table 1, replacing the masks obtained from foundation models with initial language instructions, called as RT-1-like. The experiment results show that due to overfitting, this scheme is comprehensively inferior to our paradigm, especially in the setting with unseen objects. Finally, we analyze the results to address the following questions:

- Does the segmentation mask outperform the bounding box for action prediction?
- Is tracking more robust for prompt generation than frame-by-frame detection?
- Is attention based multi-view fusion more advantageous than a single front view?
- Does incorporating a separate RGB-only branch yield better performance?

Segmentation mask is more effective than bounding box for action prediction. A related work MOO [65] uses object-centric pixel and necessitates the model to implicitly perform temporal correlation due to fixing the object-centric pixel after the first frame. To re-implement it, we attempt to extend our model to a temporal version by adding temporal attention and use the center pixel of the bounding box in the first frame to condition tasks. However, we find that this model even struggles with picking correct object instances. We conjecture that point-based prompt learning



Figure 5. Our policy model can be conditioned by assigning different values to object masks for different manipulation skills.

| Method | Seen | Unseen | New background | More distractors | Average |
|-------------------------|------|--------|----------------|------------------|---------------|
| Ours | 82.5 | 80.0 | 65.0 | 75.0 | 75.625 |
| -MOO-like [65] | 50.0 | 42.5 | 27.5 | 35.0 | 38.75 |
| -RT-1-like [4] | 65.0 | 0.0 | 20.0 | 60.0 | 36.25 |
| -replace mask with bbox | 50.0 | 40.0 | 25.0 | 30.0 | 36.25 |
| -w/o tracking | 70.0 | 50.0 | 55.0 | 70.0 | 61.25 |
| -single view | 65.0 | 80.0 | 20.0 | 70.0 | 58.75 |
| -RGB-M only | 85.0 | 70.0 | 50.0 | 70.0 | 68.75 |

Table 2. Comparison of our method and its variants on various settings.

is relatively challenging and prone to over-fit in our limited low-data regime. Additionally, we also replace the object mask with its bounding box. The results (75.625 vs. 36.25) show that the segmentation mask significantly outperforms the bounding box. In addition to provide more rich geometry and shape priors, the segmentation mask also demonstrates greater robustness to complex textures and distractors. In contrast, the bounding box struggles to achieve such precision. Moreover, explicitly incorporating a tracker allows our model to easily handle dynamic objects or those subjected to disturbances, as shown in Figure 4.

Tracking is more robust for prompt generation than frame-by-frame detection. We then replace the paradigm of first-frame detection and subsequent-frame tracking with frame-by-frame detection for prompt generation. The average success rate significantly decreases, particularly when the robot arm severely obstructs the object during grasping, illustrating the robustness of the detection-tracking paradigm. Moreover, the detection-tracking paradigm also substantially improves the inference speed.

Multi-view fusion is more beneficial compared to single view. We further investigate the conversion of the multi-view model to a single-view model, retaining only the front view. The experiment results demonstrate the effectiveness of multi-view fusion. Specifically, there is a significant drop (65.0 vs. 20.0) in the new background. We believe this is because multi-view vision can estimate depth through disparity, making it more robust than single-view vision.

Incorporating a separate RGB branch is beneficial. Finally, we implement a single-branch architecture RGB-M policy model based on ResNet-50 for a fair comparison. The experiment results demonstrate the effectiveness of our

two-stream architecture. This two-stream approach allows the model to better capture both local and global features. Additionally, the disparity (50.0 vs. 65.0) between the two models is most noticeable when altering the background, with the RGB-M model even exhibiting hovering motion in mid-air. We hypothesize that incorporating an RGB-only branch could contribute to better generalization in depth estimation, and a separate RGB-M branch may be more prone to overfitting the distribution of the training data.

4.5. Extension to other skills

Although we mainly demonstrate the effectiveness of our paradigm and policy on quantitative pick-and-place tasks evaluation, our multi-task policy model can be conditioned by assigning different values to object masks for different manipulation skills, as shown in Figure 5. Therefore, our method can be flexibly extended to some common skills by transferring language instruction of skills to mask values of manipulated objects. To verify this assumption, we further collect 100 demonstrations of other tasks to co-fine-tune our model, including “open drawer”, “pick A and place inside of B”, and “pick A and place on top of B”. Some demonstration examples are shown in Figure 4. For drawer-opening task, we mask the handle of the drawer based on the output of GPT-4, and conduct **40** trials, achieving a success rate of approximately **50%**. More qualitative results of other skills can be found in the demo video.

Additionally, we also validate folding and stacking skills on a low-cost dual-arm robot. Specifically, we co-fine-tune the previous TPM model weight using 50 new demonstrations and deploy it using Mixformer-small [8] and SAM-tiny [78]. The demos can be seen Fig 6. It is worth mention-

| | GPT-4 [52] | DetGPT [56] | MiniGPT-4 [80] |
|--------------|-------------|-------------|----------------|
| Success Rate | 0.95 | 0.75 | 0.2 |

Table 3. The reasoning performance comparison of LLMs.

| | GroundingDINO-B | Mixforemr-B | MixformerV2-S | SAM-B | SAM-T | TPM |
|---------------------|-----------------|-------------|---------------|-------|-------|------|
| Inference Time (ms) | 148.6 | 103.4 | 17.0 | 18.2 | 10.1 | 34.8 |

Table 4. The inference time for different modules and model sizes.

ing that our low-cost dual-arm does not require additional motion planning time due to without movelt calls.

The success rates for folding (seen), folding (unseen), stacking (seen), and stacking (unseen) are **85%**, **65%**, **50%**, and **30%**, respectively.

5. Discussion and Limitations

Although our approach presents a promising direction for achieving generalizable robotic manipulation, there remain some clarifications and future works.

Which manipulation skills might benefit from our paradigm? Our paradigm provides generalizable semantic, geometry, and temporal correlation priors in interacting objects. Therefore it can benefit all skills. As we observe that many end-to-end policy models fail to accurately locate objects due to lacking the equivariance of translation and rotation. Of course, we acknowledge that this requires manually designing complex prompt templates. We leave this issue to future works, such as better code generation. In addition, for contact-rich skills, extensive demonstrations are still needed to learn complex behaviors.

How scalable is our paradigm? Our paradigm aims to fully unify the generalizability of internet-scale models and the potential of imitation learning to capture multimodal action distribution from human skills at the lowest possible training and inference cost. We believe it is difficult for robot data to reach the scale of internet data, and even RT-1 cannot handle unseen object categories. Additionally, a large policy model such as RT-2 is difficult to interpret and incurs significant training and inference costs. Therefore, our multi-model paradigm has strong scalability, and improving the respective performance by scaling laws and co-ordination ability of various components in this system is a worthwhile exploration in the future.

How to improve the performance of our paradigm? We observe that the main performance bottleneck of our system lies in the connection between the language reasoning module and the detection module, as current detectors still lack many visual concepts. Therefore, a good solution is to add prompts in the LLMs: *Please add some descriptions of shape and color to the objects.* In addition, our system does not depend on completely precise masks and these situations are included in our training data.

How to improve the execution speed of our system?

The integration of foundation models in our system introduces a more noticeable latency problem. To facilitate more efficient deployment in real-world scenarios, we also recommend offline language models and distilled lightweight vision models, such as MixformerV2 [8] and MobileSAM [78]. We report a comparison of the reasoning success rate of several language models in complex environments in Table 3. Although GPT-4 still performs the best, using a relatively lightweight dedicated model (DetGPT [56]) is still a good choice. In addition, as for RTX A6000 GPU, we report the specific inference times of different modules and model sizes in Table 4, which indicates that distilled lightweight models have great value for use.

6. Conclusion

In this paper, we propose to improve the generalization abilities of robotic manipulation by transferring internet-scale vision foundation models. By utilizing specified language-reasoning mask as our condition representation, which incorporates rich semantic, geometry, and temporal correlation priors derived from vision foundation models, into the policy model, our approach significantly improves the sample efficiency. In addition, our two-stream policy model also performs excellently with its local-global manner. Extensive experiments demonstrate the effectiveness and generalization capabilities of our paradigm and policy model, especially for unseen objects and backgrounds. Finally, we show that our policy model can be conditioned by assigning different values to object masks for different skills, therefore our system is scalable for new skills.

Acknowledgement

This work is supported by the National Key R&D Program of China (No. 2022ZD0160900), the National Natural Science Foundation of China (No. 62076119), the Fundamental Research Funds for the Central Universities (No. 020214380119), Jiangsu Frontier Technology Research and Development Program (No. BF2024076), and the Collaborative Innovation Center of Novel Software Technology and Industrialization.

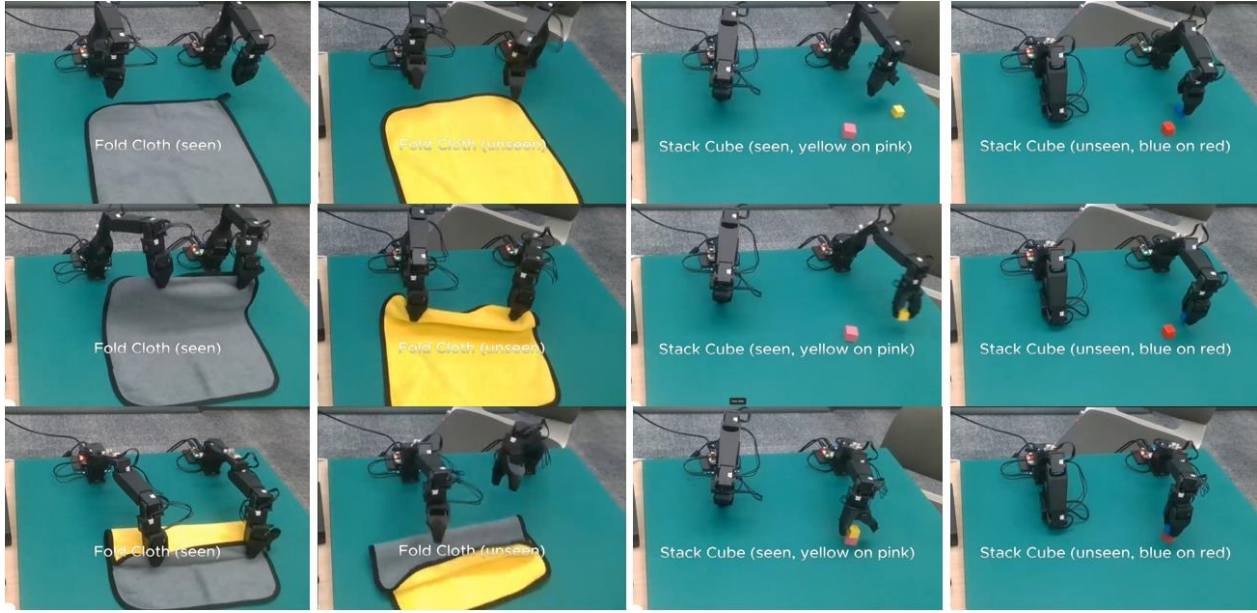


Figure 6. The demos of folding cloth and stacking cube skills.

References

- [1] 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. 3
- [2] Homanga Bharadhwaj, Abhinav Gupta, Vikash Kumar, and Shubham Tulsiani. Towards generalizable zero-shot manipulation via translating human interaction plans. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6904–6911. IEEE, 2024. 3
- [3] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023. 2
- [4] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, et al. RT-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022. 1, 2, 3, 5, 7
- [5] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation. *Advances in Neural Information Processing Systems*, 34:11781–11794, 2021. 4
- [6] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023. 3
- [7] Yutao Cui, Cheng Jiang, Limin Wang, and Gangshan Wu. MixFormer: End-to-end tracking with iterative mixed attention. In *CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 13598–13608. IEEE, 2022. 2, 4, 5
- [8] Yutao Cui, Tianhui Song, Gangshan Wu, and Limin Wang. Mixformerv2: Efficient fully transformer tracking. *arXiv preprint arXiv:2305.15896*, 2023. 7, 8
- [9] Aidan Curtis, Xiaolin Fang, Leslie Pack Kaelbling, Tomás Lozano-Pérez, and Caelan Reed Garrett. Long-horizon manipulation of unknown objects via task and motion planning with estimated affordances. In *ICRA*, pages 1940–1946. IEEE, 2022. 2
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009. 3
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR Virtual Event, Austria, May 3-7, 2021*. 5
- [12] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. PALM-E: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023. 2
- [13] Hao-Shu Fang, Chenxi Wang, Hongjie Fang, Minghao Gou, Jirong Liu, Hengxu Yan, Wenhai Liu, Yichen Xie, and Cewu Lu. Anygrasp: Robust and efficient grasp perception in spatial and temporal domains. *TRO*, 2023. 3

- [14] Yuying Ge, Annabella Macaluso, Li Erran Li, Ping Luo, and Xiaolong Wang. Policy adaptation from foundation model feedback. In *CVPR*, pages 19059–19069, 2023. 2
- [15] Xinyang Geng, Hao Liu, Lisa Lee, Dale Schuurams, Sergey Levine, and Pieter Abbeel. Multimodal masked autoencoders learn transferable representations. *arXiv preprint arXiv:2205.14204*, 2022. 2
- [16] Theophile Gervet, Zhou Xian, Nikolaos Gkanatsios, and Katerina Fragkiadaki. Act3d: 3d feature field transformers for multi-task robotic manipulation. In *CoRL*, 2023. 3
- [17] Ankit Goyal, Jie Xu, Yijie Guo, Valts Blukis, Yu-Wei Chao, and Dieter Fox. Rvt: Robotic view transformer for 3d object manipulation. In *CoRL*, pages 694–710. PMLR, 2023. 3
- [18] Pierre-Louis Guhur, Shizhe Chen, Ricardo Garcia Pinel, Makarand Tapaswi, Ivan Laptev, and Cordelia Schmid. Instruction-driven history-aware policies for robotic manipulations. In *CoRL*, pages 175–187. PMLR, 2023. 1, 3
- [19] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 5
- [20] Yingdong Hu, Renhao Wang, Li Erran Li, and Yang Gao. For pre-trained vision models in motor control, not all policy learning methods are created equal. *arXiv preprint arXiv:2304.04591*, 2023. 3
- [21] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pages 9118–9147. PMLR, 2022. 2
- [22] 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
- [23] Wenlong Huang, Fei Xia, Dhruv Shah, Danny Driess, Andy Zeng, Yao Lu, Pete Florence, Igor Mordatch, Sergey Levine, Karol Hausman, et al. Grounded decoding: Guiding text generation with grounded models for robot control. *arXiv preprint arXiv:2303.00855*, 2023. 2
- [24] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Tomas Jackson, Noah Brown, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning through planning with language models. In Karen Liu, Dana Kulic, and Jeffrey Ichnowski, editors, *CoRL*, volume 205 of *Proceedings of Machine Learning Research*, pages 1769–1782. PMLR, 2022. 2
- [25] Brian Ichter, Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, Dmitry Kalashnikov, Sergey Levine, Yao Lu, Carolina Parada, Kanishk Rao, Pierre Sermanet, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Mengyuan Yan, Noah Brown, Michael Ahn, Omar Cortes, Nicolas Sievers, Clayton Tan, Sichun Xu, Diego Reyes, Jarek Rettinghouse, Jorrell Quiambao, Peter Pastor, Linda Luu, Kuang-Huei Lee, Yuheng Kuang, Sally Jesmonth, Nikhil J. Joshi, Kyle Jeffrey, Rosario Jauregui Ruano, Jasmine Hsu, Keerthana Gopalakrishnan, Byron David, Andy Zeng, and Chuyuan Kelly Fu. Do as I can, not as I say: Grounding language in robotic affordances. In Karen Liu, Dana Kulic, and Jeffrey Ichnowski, editors, *CoRL 2022, 14-18 December 2022, Auckland, New Zealand*, volume 205, pages 287–318. PMLR, 2022. 1, 2
- [26] Eric Jang, Alex Irpan, Mohi Khansari, Daniel Kappler, Fredrik Ebert, Corey Lynch, Sergey Levine, and Chelsea Finn. BC-Z: zero-shot task generalization with robotic imitation learning. In Aleksandra Faust, David Hsu, and Gerhard Neumann, editors, *CoRL, 8-11 November 2021, London, UK*, volume 164 of *PMLR*, pages 991–1002. PMLR, 2021. 1, 2, 3
- [27] Chuhao Jin, Wenhui Tan, Jiange Yang, Bei Liu, Ruihua Song, Limin Wang, and Jianlong Fu. Alphablock: Embodied finetuning for vision-language reasoning in robot manipulation. *arXiv preprint arXiv:2305.18898*, 2023. 2
- [28] Ya Jing, Xuelin Zhu, Qie Sima, Tao zheng Yang, Yun hai Feng, Tao Kong, et al. Exploring visual pre-training for robot manipulation: Datasets, models and methods. In *CoRL 2022 Workshop on Pre-training Robot Learning*. 3
- [29] 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. 3
- [30] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [31] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 2, 4, 5
- [32] Ian Lenz, Honglak Lee, and Ashutosh Saxena. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research*, 34(4-5):705–724, 2015. 1
- [33] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9493–9500. IEEE, 2023. 2, 3
- [34] Kevin Lin, Christopher Agia, Toki Migimatsu, Marco Pavone, and Jeannette Bohg. Text2motion: From natural language instructions to feasible plans. *arXiv preprint arXiv:2303.12153*, 2023. 2
- [35] Hao Liu, Lisa Lee, Kimin Lee, and Pieter Abbeel. Instruction-following agents with jointly pre-trained vision-language models. *arXiv preprint arXiv:2210.13431*, 2022. 2
- [36] Peiqi Liu, Yaswanth Orru, Chris Paxton, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Ok-robot: What really matters in integrating open-knowledge models for robotics. *arXiv preprint arXiv:2401.12202*, 2024. 2
- [37] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 2, 4, 5

- [38] YuXuan Liu, Xi Chen, and Pieter Abbeel. Self-supervised instance segmentation by grasping. *arXiv preprint arXiv:2305.06305*, 2023. 1
- [39] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *CVPR*, pages 10012–10022, 2021. 5
- [40] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5
- [41] Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre Sermanet. Learning latent plans from play. In Leslie Pack Kaelbling, Danica Kragic, and Komei Sugiura, editors, *CoRL 2019, Osaka, Japan, October 30 - November 1, 2019, Proceedings*, volume 100 of *Proceedings of Machine Learning Research*, pages 1113–1132. PMLR, 2019. 1
- [42] Corey Lynch, Ayzaan Wahid, Jonathan Tompson, Tianli Ding, James Betker, Robert Baruch, Travis Armstrong, and Pete Florence. Interactive language: Talking to robots in real time. *arXiv preprint arXiv:2210.06407*, 2022. 1, 3
- [43] Yecheng Jason Ma, William Liang, Vaidehi Som, Vikash Kumar, Amy Zhang, Osbert Bastani, and Dinesh Jayaraman. Liv: Language-image representations and rewards for robotic control. *arXiv preprint arXiv:2306.00958*, 2023. 3
- [44] Yecheng Jason Ma, Shagun Sodhani, Dinesh Jayaraman, Osbert Bastani, Vikash Kumar, and Amy Zhang. Vip: Towards universal visual reward and representation via value-implicit pre-training. *arXiv preprint arXiv:2210.00030*, 2022. 3
- [45] Jeffrey Mahler, Jacky Liang, Sherdil Niyaz, Michael Laskey, Richard Doan, Xinyu Liu, Juan Aparicio Ojea, and Ken Goldberg. Dex-net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics. *arXiv preprint arXiv:1703.09312*, 2017. 3
- [46] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Yecheng Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Pieter Abbeel, Jitendra Malik, et al. Where are we in the search for an artificial visual cortex for embodied intelligence? *arXiv preprint arXiv:2303.18240*, 2023. 3
- [47] Oier Mees, Jessica Borja-Diaz, and Wolfram Burgard. Grounding language with visual affordances over unstructured data. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11576–11582. IEEE, 2023. 3
- [48] Oier Mees, Lukas Hermann, and Wolfram Burgard. What matters in language conditioned robotic imitation learning over unstructured data. *IEEE Robotics and Automation Letters*, 7(4):11205–11212, 2022. 3
- [49] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and Automation Letters*, 7(3):7327–7334, 2022. 3
- [50] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. EmbodiedGPT: Vision-language pre-training via embodied chain of thought. *arXiv preprint arXiv:2305.15021*, 2023. 2
- [51] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3M: A universal visual representation for robot manipulation. In Karen Liu, Dana Kulic, and Jeffrey Ichnowski, editors, *CoRL 2022, 14-18 December 2022, Auckland, New Zealand*, volume 205, pages 892–909. PMLR, 2022. 3
- [52] OpenAI. GPT-4 technical report, 2023. 2, 3, 8
- [53] Abhishek Padalkar, Acorn Pooley, Ajinkya Jain, Alex Bewley, Alex Herzog, Alex Irpan, Alexander Khazatsky, Anant Rai, Anikait Singh, Anthony Brohan, et al. Open x-embodiment: Robotic learning datasets and rt-x models. *arXiv preprint arXiv:2310.08864*, 2023. 2
- [54] Simone Parisi, Aravind Rajeswaran, Senthil Purushwalkam, and Abhinav Gupta. The unsurprising effectiveness of pre-trained vision models for control. In *International Conference on Machine Learning*, pages 17359–17371. PMLR, 2022. 3
- [55] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 5
- [56] Renjie Pi, Jiahui Gao, Shizhe Diao, Rui Pan, Hanze Dong, Jipeng Zhang, Lewei Yao, Jianhua Han, Hang Xu, and Lingpeng Kong Tong Zhang. DetGPT: Detect what you need via reasoning. *arXiv preprint arXiv:2305.14167*, 2023. 2, 8
- [57] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICLR*, pages 8748–8763. PMLR, 2021. 2
- [58] Ilija Radosavovic, Tete Xiao, Stephen James, Pieter Abbeel, Jitendra Malik, and Trevor Darrell. Real-world robot learning with masked visual pre-training. In Karen Liu, Dana Kulic, and Jeffrey Ichnowski, editors, *CoRL 2022, 14-18 December 2022, Auckland, New Zealand*, volume 205 of *Proceedings of Machine Learning Research*, pages 416–426. PMLR, 2022. 3
- [59] Rouhollah Rahmatizadeh, Pooya Abolghasemi, Ladislau Bölöni, and Sergey Levine. Vision-based multi-task manipulation for inexpensive robots using end-to-end learning from demonstration. In *ICRA*, pages 3758–3765. IEEE, 2018. 1
- [60] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*. 3
- [61] Kenneth Shaw, Shikhar Bahl, and Deepak Pathak. Videodex: Learning dexterity from internet videos. In *Conference on Robot Learning*, pages 654–665. PMLR, 2023. 3
- [62] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. CLIPort: What and where pathways for robotic manipulation. In Aleksandra Faust, David Hsu, and Gerhard Neumann, editors, *CoRL, 8-11 November 2021, London, UK*, volume 164 of

- Proceedings of Machine Learning Research*, pages 894–906. PMLR, 2021. 1, 3
- [63] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *CoRL*, pages 785–799. PMLR, 2023. 3
- [64] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11523–11530. IEEE, 2023. 2, 3
- [65] Austin Stone, Ted Xiao, Yao Lu, Keerthana Gopalakrishnan, Kuang-Huei Lee, Quan Vuong, Paul Wohlhart, Brianna Zitkovich, Fei Xia, Chelsea Finn, et al. Open-world object manipulation using pre-trained vision-language models. *arXiv preprint arXiv:2303.00905*, 2023. 2, 6, 7
- [66] Martin Sundermeyer, Arsalan Mousavian, Rudolph Triebel, and Dieter Fox. Contact-grasnet: Efficient 6-dof grasp generation in cluttered scenes. In *ICRA*, pages 13438–13444. IEEE, 2021. 3
- [67] Andreas Ten Pas, Marcus Gualtieri, Kate Saenko, and Robert Platt. Grasp pose detection in point clouds. *IJRR*, 36(13-14):1455–1473, 2017. 3
- [68] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 5
- [69] Naoki Wake, Atsushi Kanehira, Kazuhiro Sasabuchi, Jun Takamatsu, and Katsushi Ikeuchi. Gpt-4v (ision) for robotics: Multimodal task planning from human demonstration. *arXiv preprint arXiv:2311.12015*, 2023. 2
- [70] Chen Wang, Linxi Fan, Jiankai Sun, Ruohan Zhang, Li Fei-Fei, Danfei Xu, Yuke Zhu, and Anima Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play. *arXiv preprint arXiv:2302.12422*, 2023. 3
- [71] 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. 3
- [72] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. CBAM: convolutional block attention module. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part VII*, volume 11211 of *Lecture Notes in Computer Science*, pages 3–19. Springer, 5
- [73] 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. 3
- [74] Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and Baining Guo. Advancing high-resolution video-language representation with large-scale video transcriptions. In *CVPR*, pages 5036–5045, 2022. 3
- [75] Jiange Yang, Bei Liu, Jianlong Fu, Bocheng Pan, Gangshan Wu, and Limin Wang. Spatiotemporal predictive pre-training for robotic motor control. *arXiv preprint arXiv:2403.05304*, 2024. 3
- [76] Jiange Yang, Haoyi Zhu, Yating Wang, Gangshan Wu, Tong He, and Limin Wang. Tra-moe: Learning trajectory prediction model from multiple domains for adaptive policy conditioning. *arXiv preprint arXiv:2411.14519*, 2024. 3
- [77] Lin Yen-Chen, Andy Zeng, Shuran Song, Phillip Isola, and Tsung-Yi Lin. Learning to see before learning to act: Visual pre-training for manipulation. In *ICRA*, pages 7286–7293. IEEE, 2020. 1, 3
- [78] Chaoning Zhang, Dongshen Han, Yu Qiao, Jung Uk Kim, Sung-Ho Bae, Seungkyu Lee, and Choong Seon Hong. Faster segment anything: Towards lightweight sam for mobile applications. *arXiv preprint arXiv:2306.14289*, 2023. 7, 8
- [79] Tony Z Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. *arXiv preprint arXiv:2304.13705*, 2023. 3
- [80] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 8
- [81] Haoyi Zhu, Honghui Yang, Yating Wang, Jiange Yang, Limin Wang, and Tong He. Spa: 3d spatial-awareness enables effective embodied representation. *arXiv preprint arXiv:2410.08208*, 2024. 3