# *RoboMM*: All-in-One Multimodal Large Model for Robotic Manipulation

Feng Yan,* Fanfan Liu,* Liming Zheng, Yufeng Zhong, Yiyang Huang,
Zechao Guan, Chengjian Feng, Lin Ma†

*Meituan Inc.*

https://github.com/RoboUniview/RoboMM

## Abstract

*In recent years, robotics has advanced significantly through the integration of larger models and large-scale datasets. However, challenges remain in applying these models to 3D spatial interactions and managing data collection costs. To address these issues, we propose the multimodal robotic manipulation model, RoboMM, along with the comprehensive dataset, RoboData. RoboMM enhances 3D perception through camera parameters and occupancy supervision. Building on OpenFlamingo, it incorporates Modality-Isolation-Mask and multimodal decoder blocks, improving modality fusion and fine-grained perception. RoboData offers the complete evaluation system by integrating several well-known datasets, achieving the first fusion of multi-view images, camera parameters, depth maps, and actions, and the space alignment facilitates comprehensive learning from diverse robotic datasets. Equipped with RoboData and the unified physical space, RoboMM is the generalist policy that enables simultaneous evaluation across all tasks within multiple datasets, rather than focusing on limited selection of data or tasks. Its design significantly enhances robotic manipulation performance, increasing the average sequence length on the CALVIN from 1.7 to 3.3 and ensuring cross-embodiment capabilities, achieving state-of-the-art results across multiple datasets.*

## 1. Introduction

In recent years, machine learning has experienced profound advancements, from the advent of CLIP [53, 67] to the progression of foundational models like the GPT series [1, 8, 9], Llama [59, 60], LLaVA [37], and Flamingo [2, 3]. These strides are largely due to larger transformer-based architectures and the utilization of "internet-scale" datasets [11, 13, 14, 36]. These innovations have not only extended the frontiers of natural language processing [16]

---

*Equal contribution.
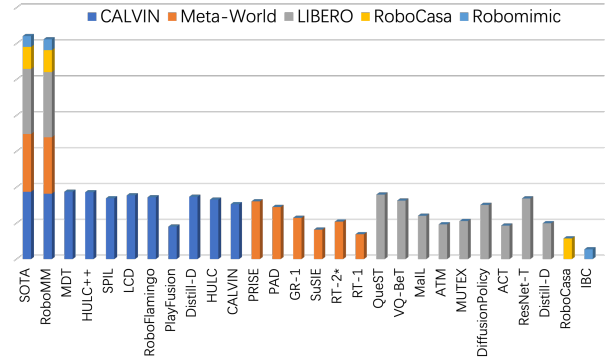†Correspondence to Lin Ma at malin11@meituan.com

Figure 1. Stacked bar chart depicting the performance of various models across five datasets. The "SOTA" label represents the best results achieved by specialized models for each dataset. Notably, *RoboMM* is the only generalist policy evaluated on multiple datasets, demonstrating competitive performance relative to "SOTA". For more details, please refer to Section 5.1.

and computer vision [22, 54] but have also galvanized researchers to integrate these models into Embodied Artificial Intelligence (EAI) [49], thereby enabling more complex and varied tasks in real-world environments.

In terms of modeling, there has been a gradual shift from single-task or single-dataset learning [42, 57, 67] towards transfer learning approaches [7, 25, 30, 35, 58, 61]. These models leverage robust foundational models that are pretrained on extensive "internet-scale" datasets or diverse data sources. Subsequently, they undergo fine-tuning on specific robotic datasets to produce precise control actions. On the data front, researchers has collected data through various means to augment models. For example, the Open X-Embodiment [51] amalgamates diverse datasets containing vision-language-action pairs, whereas the RH20T [18] gathers data via teleoperation. Despite the impressive robustness of these efforts, they still encounter significant challenges in practical applications.

Firstly, **is the direct application of multimodal models to EAI the optimal solution?** Robots must interact with the physical 3D space; however, current multimodal models
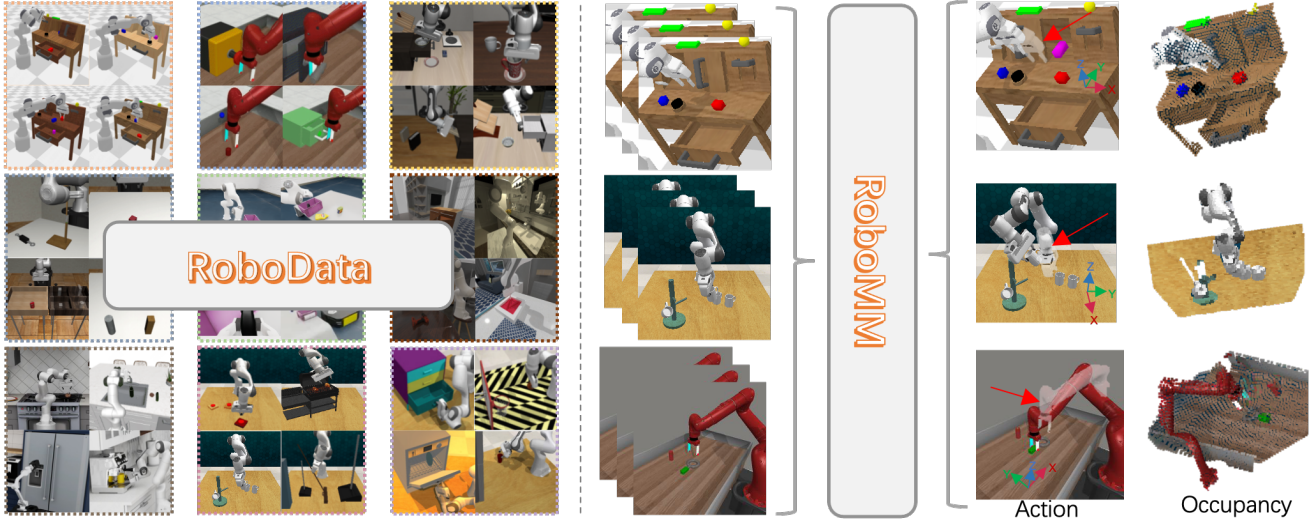
Figure 2. **Left:** *RoboData* integrates nine diverse and complex datasets (CALVIN [41], Meta-World [65], LIBERO [34], Robomimic [39], RoboCAS [68], ManiSkill2 [21], RoboCasa [47], RLBench [26], and Colosseum [52]), covering various robot embodiments, environments, and task types. Standardized input and output spaces achieve a unified dataset. **Right:** *RoboMM* features comprehensive 3D environmental perception capabilities, flexible and diverse multimodal outputs, and significantly enhances the robotic manipulation generalization capabilities.

predominantly focus on 2D image understanding and generation, which may limit their practical applicability. Secondly, **is it essential to address the cost and efficiency of dataset construction?** For instance, collecting around 130,000 episodes from the RT-1 [6] took 17 months. Therefore, it is imperative to integrate as many existing multi-platform, multi-robot datasets from the industry as possible to address this urgent issue.

In response to these challenges, this paper introduces *RoboMM*, the native multimodal large model for robotic manipulation, and *RoboData*, the comprehensive dataset integrating datasets across various platforms and robots for evaluation and training purposes.

In terms of modeling, *RoboMM* combines camera parameters and occupancy supervision to enhance 3D environmental perception. Additionally, leveraging large language models like OpenFlamingo [3], we design the plug-and-play efficient Modality-Isolation-Mask (MIM), which flexibly introduces multimodal supervision. This not only grants the model fine-grained perception capabilities but also enables the effective utilization of vast amounts of internet data. On the data side, while Open X-Embodiment [51] integrates multiple datasets, it lacks critical information such as multi-view images, camera parameters, and depth maps, making it more suitable for 2D multimodal training. Additionally, the lack of data space alignment prevents the robot's 6D pose from being consistent across different datasets. In contrast, *RoboData* addresses these limitations by integrating a wide range of well-known industry datasets, including CALVIN [41],

Meta-World [65], LIBERO [34], Robomimic [39], RoboCAS [68], ManiSkill2 [21], RoboCasa [47], RLBench [26], and Colosseum [52], which are not covered in Open X-Embodiment [51]. We dedicate 200 person-days to collect and organize these data, supplementing missing modalities such as depth maps and camera parameters. Furthermore, thanks to the unified physical space, *RoboData* aligns the input and output spaces of cross robots and platforms, ensuring consistency and facilitating integrated learning from diverse robotic datasets.

***RoboData* aims to provide the industry with a comprehensive and fair evaluation system, while *RoboMM* is the generalist policy to incorporate training and testing across multiple datasets.** Extensive experiments demonstrate that the various components of *RoboMM* significantly improve performance in robotic manipulation tasks, enhancing the average sequence length on the CALVIN [41] benchmark from 1.7 to 3.3. Additionally, *RoboMM* ensures cross-embodiment capabilities and achieves state-of-the-art results across multiple datasets, as shown in Figure 1. This paper underscores the critical role of advanced models and curated datasets in advancing robotics, highlighting their potential to drive significant improvements in robotic performance and functionality.

## 2. Related Work

**Robotic Datasets.** In the early stages of robotics research, it is typically necessary to collect specific datasets for each robot, task, and environment, such as RLBench [26] and
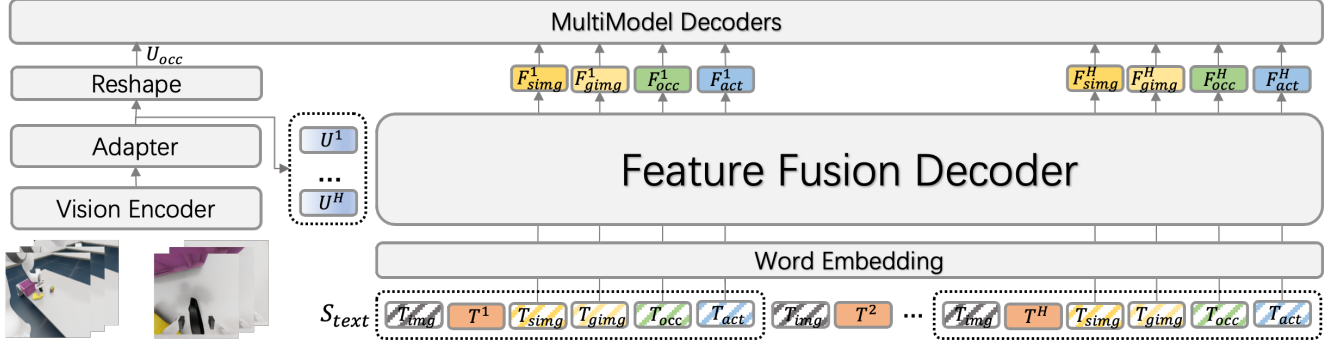
Figure 3. Architecture of *RoboMM*. Vision Encoder Block for extracting multi-view features, Adapter Block that leverages occupancy supervision to unify features and enhance spatial perception, Feature Fusion Block based on LLMs for merging text and visual information and Multimodal Decoder Blocks that enhance fine-grained perception and understanding through multimodal outputs.

CALVIN [41]. Although these datasets are highly customized and of high quality, they are limited in quantity and have poor generalization capabilities. To further enhance model performance and generalization, researchers have collected large amounts of data through teleoperation methods, such as RT-1 [6] and RH20T [18]. These large-scale datasets cover more scenarios and tasks, supporting multi-task learning, but also bring high data annotation costs. As research progresses, methods for integrating multiple datasets, such as Open X-Embodiment [51] and DROID [28], have been proposed to improve model generalization and data utilization efficiency by merging data from different sources. However, these methods also face issues of data inconsistency and potential biases. This paper proposes *RoboData*, which efficiently integrates multiple datasets and unifies the input and output spaces, thereby addressing data heterogeneity. Additionally, it breaks the limitation of training for a single specific task, providing a unified benchmark for robotic manipulation.

**Robotic Policies.** Previous works such as R3M [46], VC-1 [38], ACT [67], and HULC++ [42] typically employ strategies with a small number of parameters. Subsequent models like RoboFlamingo [33], Corki [25], and RoboUniView [35] have built on multimodal large models but have only fine-tuned on limited datasets. Despite advancements in multi-task learning and few-shot learning, recent models such as RT-X [51], Octo [58], HPT [61], CrossFormer [17], GR-2 [10], and OpenVLA [30] have trained vision-language-action robotic policies on various datasets. However, these works often pre-train on data from real robots [18, 28], human videos [20, 46], and simulation domains [41, 68], neglecting the uniformity of physical space, and achieve good performance only after fine-tuning on specific datasets. Given that robots operate in 3D physical environments, their perception and interaction capabilities must integrate 3D sensing, akin to the requirements of autonomous driving systems.

## 3. *RoboMM*

### 3.1. Review

Multimodal Large Language Models (MLLMs) typically consist of three main components: the modality encoder (Enc), the adapter (Adapter), and the large language model (LLM), mathematically expressed as follows:

$$\begin{aligned} O_T &= \text{MLLM}(I, T) \\ &= \text{LLM}\left(\text{Adapter}\left(\text{Enc}(I)\right), \text{WE}(T)\right). \end{aligned} \tag{1}$$

Here, WE denotes the word embedding layer. The modality encoder transforms inputs from single modality into appropriate representations. For instance, the image encoder extracts features $F_I$ from input images $I$. Common visual encoders like CLIP [53] are pre-trained on image-text pairs, aligning visual and textual semantics for easier integration with LLMs. The adapter maps features from visual and other modalities into inputs $U$ that the LLM can understand. For example, Blip2 [32] uses Q-Former for feature interaction; LLaVA [37] employs MLPs to align visual features with text features. The large language model is the core component of our framework, referred to in this paper as the Feature Fusion Decoder. It typically employs auto-regressive models such as LLaMA [59] or GPT [1], as well as cross-attention models like Flamingo [2] or LLaMA3.2 [44]. This model fuses the feature representations $U$ with text features $F_T$ extracted from the word embedding layer to generate the final textual output $O_T$. This integration of features enhances the model's ability to produce contextually relevant responses.

### 3.2. Network Architecture

For robotic manipulation tasks based on language instructions $T$, which typically rely on historical frames $I$ from $N$ perspectives at $H$ time steps, the task can be mathematically expressed as $O_A = \Theta(I, T)$, where $I \in \mathbb{R}^{H \times N \times H \times W \times 3}$. Integrating the principles of MLLMs, this paper proposes the novel native multimodal robotic manipulation model

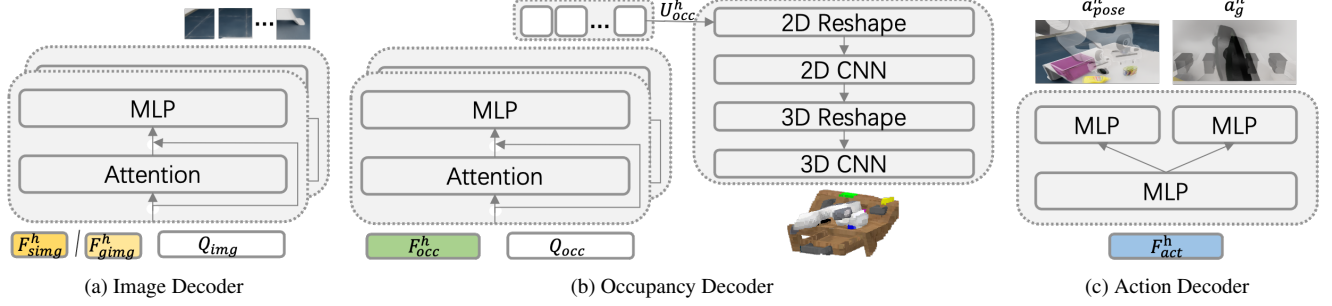| (a) Image Decoder | (b) Occupancy Decoder | (c) Action Decoder |

Figure 4. Overview of the multimodal decoders. (a) Image Decoder, (b) Occupancy Decoder, and (c) Action Decoder. Each decoder processes input features through a series of Multi-Layer Perceptrons (MLPs), attention mechanisms, and convolutional neural networks (CNNs) to generate appropriate output representations.

named *RoboMM*, as show Figure 3, which has the capability of 3D environment perception and handles multimodal inputs (text $T$, vision $I$, camera parameters $Cam$) and outputs (actions $O_A$, images $O_I$, occupancy $O_o$):

$$(O_A, [O_I, O_O]) = \text{RoboMM}(T, I, Cam). \quad (2)$$

*RoboMM* consists of the following key components: (1) Vision Encoder Block: Designed to extract observation features $F_I^{h,n}$ from $H$ time steps and $N$ perspectives. (2) 3D Perception Adapter Block: Enhances physical space perception by integrating camera parameters. (3) Feature Fusion Decoder Based on Large Language Models: Merges text and visual information to output multimodal features, with the use of Modality-Isolation-Mask (MIM) increasing the flexibility of modality fusion. (4) Multimodal Decoder Blocks: Enhances the model's fine-grained perception and understanding through multimodal outputs. Notably, thanks to MIM, $O_I, O_o$ are optional outputs.

**Adapter:** We employ UVFormer from RoboUni-View [35], which is a simple yet powerful 3D environment perception model. UVFormer takes image features $X^h = \{F_I^{h,n}\}_n^N$, camera parameters $Cam^h = \{Cam^{h,n}\}_n^N$, and learnable unified view queries $Q$ as inputs and outputs the unified view representation $U_I^h$:

$$U_I^h = \text{UVFormer}(Q, X^h, Cam^h). \quad (3)$$

Here, $Q = \{Pos, Emb\}$, $Pos \in \mathbb{R}^{L \times B \times 3P}$ and $Emb \in \mathbb{R}^{L \times B \times C}$ represent the positions and learnable features of the queries, respectively. $L$, $B$, and $P$ define the spatial shape of the 3D grid within the operational space of the robot. Specifically, $Emb_{l,b} \in \mathbb{R}^C$ is responsible for the corresponding pillar cell area in unified view space. $U_I^h \in \mathbb{R}^{L \times B \times C}$ is the unified view representation, containing all relevant information in the $L \times B \times P$ 3D grid.

**Feature Fusion Decoder:** Due to the need to support multi-frame or video inputs, we abandon the Auto-Regressive (AR) mechanism used in LLaVA [36] and adopt OpenFlamingo [3] with cross-attention as the Feature Fusion Decoder. It integrates unified visual representations

with language and other modality placeholders through cross-attention layers.

(a) To support multimodal output, we first construct the text sequence $T'$, which includes text and multiple modality read-out tokens:

$$T' = \{[T_{\text{img}}, T^h, T_{\text{simg}}, T_{\text{gimg}}, T_{\text{occ}}, T_{\text{act}}]\}_h^H. \quad (4)$$

Here, $T' \in \mathbb{R}^{\sum_h^H (1 + L^h + 8*3+1)}$, $T_{\text{simg}}, T_{\text{gimg}}, T_{\text{occ}}, T_{\text{act}}$ represent read-out tokens for static images, wrist images, occupancy, and actions, respectively. $L^h$ represents the length of $T^h$. $T_{\text{img}}$ is used to indicate the position of the original image. $T_{\text{simg}}, T_{\text{gimg}}, T_{\text{occ}}$ each use 8 tokens. We then feed the constructed text sequence into the word embedding layer to obtain text features:

$$F_T = \text{WE}(T'). \quad (5)$$

(b) Attention Fusion: Continuing with the use of cross-attention in OpenFlamingo [3], we fuse visual and text features wherein the text features $F_T^h$ serve as the query, and the visual features $U_I^h$ serve as the key and value. It is worth noting that the self-attention layer incorporates MIM(see in Figure 5 Left), which allows training with auxiliary modality supervision and omitting unnecessary modalities during inference, significantly increasing the flexibility of modality fusion.

(c) Multimodal Output Features: According to the modality read-out tokens, their corresponding output features are indexed as $F_{\text{simg}}^h$, $F_{\text{gimg}}^h$, $F_{\text{occ}}^h$, and $F_{\text{act}}^h$.

**Multimodal Decoders:** We design different decoder modules to accommodate various modalities.

(a) Image Decoder: As shown in Figure 4a, we design a simple structure that includes 2 attention decoder layers. This structure outputs image patches, which are then assembled into a complete image(static images $O_{simg}^h$ or wrist image $O_{gimg}^h$) based on their coordinates.

(b) Occupancy Decoder: The initial part of this structure, as shown in Figure 4b, is similar to Image Decoder,

generating the feature $U_{occ}^h$. Then, $U_{occ}^h$ is reshaped, upsampled, and processed through 3D convolutions to reconstruct the entire 3D occupancy $O_o^h = \{o_{pos}^h, o_{rgb}^h\}$. The flexible model architecture allows the visual module to generate the occupancy map $O_o^h = O_{o_v}^h$ from multi-view features using UVFormer, while LLM also outputs $O_o^h = O_{o_m}^h$, corresponding to $T_{occ}$. Experimental validation shows that $O_{o_v}^h$ and $O_{o_m}^h$ provide similar assistance for robotic manipulation. Unless otherwise specified, in this paper, $O_o^h = O_{o_v}^h$.

(c) Action Decoder: In Figure 4c, we use a few MLP layers to output actions $O_A^h$ consisting of delta 6D poses $a_{pose}^h = \{\Delta\text{pos}_x^h, \Delta\text{pos}_y^h, \Delta\text{pos}_z^h, \Delta\text{rot}_x^h, \Delta\text{rot}_y^h, \Delta\text{rot}_z^h\}$ and 1-DoF gripper actions $a_g^h$.

## 3.3. Training Objective

In multimodal learning tasks, we design the comprehensive loss function $l$ to enhance overall model performance by optimizing different modality outputs:

$$l = l_a + \lambda_{\text{image}} (l_{simg} + l_{gimg}) + \lambda_{\text{occ}} l_o. \quad (6)$$

Here, $\lambda_{\text{image}}$ and $\lambda_{\text{occ}}$ are weight coefficients for image and occupancy losses. Notably, $l_{simg}$, $l_{gimg}$, and $l_o$ can be excluded from training if the corresponding modality is unavailable, providing a flexible optimization framework for multimodal learning.

**Action Loss** $l_a$**:** The action loss function optimizes $a_{pose}^h$ and $a_{gripper}^h$ using the combination of Mean Squared Error (MSE) and Binary Cross-Entropy (BCE):

$$l_a = \sum_h \left( \text{MSE}(a_{pose}^h, \hat{a}_{pose}^h) + \lambda_g \text{BCE}(a_g^h, \hat{a}_g^h) \right). \quad (7)$$

Here, $\lambda_g$ is a parameter used to adjust the weight of the gripper state loss, and $\hat{a}_{pose}^h$ and $\hat{a}_g^h$ represent the example action pose and gripper state at time step $h$, respectively.

**Image Loss** $l_{simg}$ **or** $l_{gimg}$**:** The image loss measures pixel-level differences between the predicted images and the next frame images $\hat{I}_{simg}^{h+1}$ or $\hat{I}_{gimg}^{h+1}$ using L2 loss:

$$l_{simg} = \sum_h \sum_{\text{pixels}} \|O_{simg}^h - \hat{I}_{simg}^{h+1}\|_2^2, \quad (8)$$

$$l_{gimg} = \sum_h \sum_{\text{pixels}} \|O_{gimg}^h - \hat{I}_{gimg}^{h+1}\|_2^2. \quad (9)$$

This effectively guides the model in predicting future scenes.

**Occupancy Loss** $l_o$**:** The occupancy loss incorporates spatial position and RGB color information:

$$l_o = \sum_h \sum_{\text{points}} l_{o'}, \quad (10)$$

$$l_{o'} = \|o_{pos}^h - \hat{o}_{pos}^h\|_2^2 + \lambda_{rgb}\|o_{rgb}^h - \hat{o}_{rgb}^h\|_2^2, \quad (11)$$

where $\lambda_{rgb}$ adjusts the contributions of position and RGB color losses.
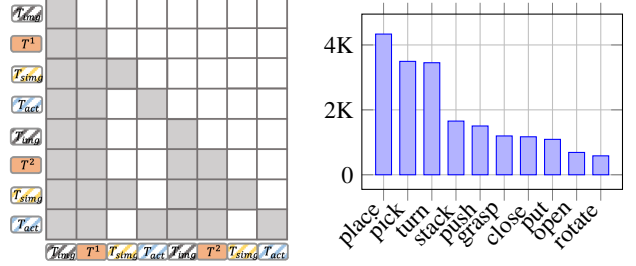


Figure 5. **Left:** Modality Isolation Mask (MIM). The KQ mask structure regulates attention interactions among different modalities (e.g., <text>, <image>, <action>). Dark squares indicate allowed attention connections between keys (K) and queries (Q), while white squares denote prohibited attention, ensuring modality isolation. **Right:** Frequency of Tasks. This section illustrates the distribution of tasks within the dataset, detailing the number of episodes associated with each task. The bars represent the frequency of various tasks, including "place," "pick," and "turn," highlighting the diversity and focus areas of the dataset. The y-axis indicates the number of episodes, emphasizing the relative frequency of each task.

## 4. *RoboData*

The rise of ChatGPT [50] and large AI models [2, 4, 59] signifies a paradigm revolution in artificial intelligence, all built upon the foundation of rich "internet-scale" datasets. However, in the domain of embodied intelligence, research still focuses on single, specific tasks such as grasping, path planning, and pick-and-place, aiming to train agents tailored for particular scenarios. Although projects like Open X-Embodiment [51] and ARIO [62] compile multiple datasets, they still present numerous issues. For example, they lack essential 3D information—such as multiview, camera intrinsic and extrinsic parameters, and depth maps—making these datasets suitable only for 2D multimodal training. Moreover, there is a lack of proper spatial alignment across datasets; specifically, the 6D poses (i.e., position and orientation) of the robotic end-effectors recorded exhibit inconsistencies due to different world coordinate systems.

To address these challenges, we curate well-known datasets from the industry, including CALVIN [41], Meta-World [65], LIBERO [34], Robomimic [39], RoboCasa [47], ManiSkill2 [21], RoboCAS [68], RLBench [26], and Colosseum [52], forming a comprehensive dataset we call *RoboData*. This dataset aims to provide the industry with a complete and fair evaluation system, comprising 70,000 episodes and 7 million samples. It encompasses a diverse range of tasks, including placing, picking, turning, and stacking. Figure 5 Right illustrates these tasks along with their corresponding number of episodes, highlighting the distribution of tasks within the dataset.

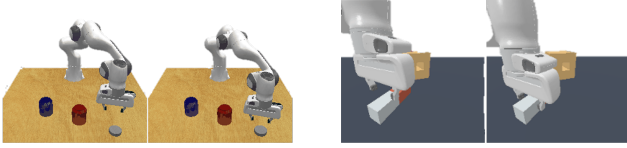As shown in Table 4, each dataset encompasses different simulation platforms and possesses unique world coor-

Figure 6. The downward movement of the robot in the RLBench and ManiSkill2 environments.

dinate systems, workspaces, perspectives, and other characteristics. Therefore, we align the input and output spaces of the models based on multiple influencing factors.

**3D Space Alignment:** We focus on the unification of world coordinate systems, workspaces, and action spaces. Different datasets adopt their own coordinate systems. For example, the RLBench [26] and Colosseum [52] reference the robot's body, setting the X-axis to point forward, the Y-axis to point **left**, and the Z-axis to point **upward**; whereas ManiSkill2 [21] orient the X-axis forward, the Y-axis **right**, and the Z-axis **downward**. Figure 6 illustrate the movements of robotic arms from RLBench [26] and ManiSkill2 [21] datasets, respectively. Although both exhibit similar motion directions (moving from top to bottom), the representation of actions differs significantly due to the variations in coordinate systems. For instance, in RLBench [26], $a_{pose} = [0.0, 0.0, -0.1, 0.0, 0.0, 0.0]$, while in ManiSkill2 [21], $a_{pose} = [0.0, 0.0, 0.1, 0.0, 0.0, 0.0]$.

Unifying all data into the same coordinate system is crucial for conducting cross-platform joint training. If this unification is not achieved, conflicts may arise during the training process, negatively impacting the final learning outcomes. **Benefiting from the supplementation of 3D information in the *RoboData* dataset**, we rotate the coordinate systems to unify all data to the same orientation: the X-axis pointing right, the Y-axis pointing forward, and the Z-axis pointing upward. For example, the original world coordinate system of Robomimic [39], $W_{Robomimic}^{ori}$, is transformed into $W_{Robomimic}$:

$$W_{Robomimic} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0.4 \\ 0 & 0 & 0 & 1 \end{bmatrix} W_{Robomimic}^{ori}.$$

In parallel, we limit the workspace through translation to the range of [-0.5, -0.5, 0] to [0.5, 0.5, 1].

**Action Representation Alignment:** Various datasets employ different methods to obtain robot actions, and this diversity can lead to inconsistencies in the data. For instance, CALVIN [41] uses the Euler Angle Difference Method (EADM) for action representation, while LIBERO [34], Robomimic [39], RoboCasa [47] utilize the Composite Rotation Matrix Method (CRMM), and ManiSkill2 [21] adopts the Pose Composition Method (PCM). To address this issue, we unify the action representation

across all datasets by regenerating them using the CRMM, which has been validated in SRT [29]. This choice not only enhances data consistency but also provides a more reliable foundation for subsequent research. For more detailed technical information, please refer to Chapter 8.

**Missing Data Imputation:** In many datasets, camera intrinsics and extrinsics are often not directly provided, which poses challenges for research. To tackle this problem, we reconstruct the original simulations and render them anew, leveraging the original data to drive the robots in order to obtain these critical parameters. This process ensures data completeness and provides essential support for subsequent experiments.

In conclusion, by effectively aligning the input and output spaces, we propose a new dataset standard aimed at comprehensively optimizing existing datasets. The implementation of this standard will facilitate the development of more versatile and general-purpose embodied AI agents. **We plan to release this dataset along with evaluation code to promote further research and advancement in the field of embodied intelligence. Additionally, we welcome all like-minded researchers to join our efforts in driving the development of this field.**

## 5. Experiments

In the previous sections, we elaborate on the *RoboMM* framework and the characteristics of the *RoboData* dataset. Next, we address the following questions using task success rates: its performance across multiple datasets, the importance of each module within *RoboMM*.

### 5.1. Results and Analysis

**Protocol.** Due to the substantial time and computational resources required for model training, this study chooses to train the *RoboMM* models without Image Loss on five datasets: CALVIN [41], Meta-World [65], LIBERO [34], RoboCasa [47], and Robomimic [39], as shown in the Figure 8. A total of 10 epochs are trained, utilizing 2.1 million samples, employing 32 A100 GPUs with 80GB each, and taking 130 hours.

During the evaluation phase, we adhere to the official configurations provided for each dataset. Specifically, for the CALVIN dataset [41], a sequence length of 1000 is used, while for the other datasets, each task was evaluate over 20 episodes with varying initial conditions. The evaluation metric adopted is success rate(SR), to comprehensively assess the model's performance.

**Experiment.** The results summarized in Table 1 reveal that *RoboMM* exhibits exceptional performance across the evaluated datasets, achieving state-of-the-art results in particular on the LIBERO [34] and RoboCasa [47] datasets. Notably, in the CALVIN dataset [41], *RoboMM* attain a success rate (SR) of 91.0%, which is competitive with the per-

| Dataset | Model | Source | SR↑ |
|---|---|---|---|
| CALVIN | MDT [55] | RSS'24 | 93.7% |
| | HULC++ [43] | ICRA'24 | 93.0% |
| | SPIL [70] | RA-L'24 | 84.6% |
| | LCD [66] | arXiv'23 | 88.7% |
| | RoboFlamingo [33] | arXiv'23 | 86.0% |
| | PlayFusion [12] | CoRL'23 | 45.2% |
| | Distill-D [24] | CoRL'23 | 86.7% |
| | HULC [40] | RA-L'22 | 82.7% |
| | CALVIN [41] | RA-L'22 | 76.4% |
| | RoboMM (ours) | - | 91.0% |
| | RoboMM⁻ (ours) | - | 74.7% |
| Meta-World | PRISE [69] | ICML'24 | 80.4% |
| | PAD [23] | NeurIPS'24 | 72.5% |
| | GR-1 [64] | ICLR'24 | 57.4% |
| | SuSIE [5] | ICLR'24 | 41.0% |
| | RT-2* [7] | arXiv'23 | 52.2% |
| | RT-1 [6] | RSS'23 | 34.6% |
| | RoboMM (ours) | - | 78.6% |
| | RoboMM⁻ (ours) | - | 79.3% |
| LIBERO | QueST [45] | arXiv'24 | 89.8% |
| | VQ-BeT [31] | ICML'24 | 81.4% |
| | MDT [55] | RSS'24 | 67.2% |
| | MaIL [27] | CoRL'24 | 60.3% |
| | PRISE [69] | ICML'24 | 54.4% |
| | ATM [63] | RSS'24 | 48.4% |
| | MUTEX [56] | CoRL'23 | 53.0% |
| | DiffusionPolicy [15] | IJRR'23 | 75.4% |
| | ACT [67] | RSS'23 | 46.6% |
| | ResNet-T [34] | NeurIPS'23 | 84.4% |
| | Distill-D [24] | CoRL'23 | 49.9% |
| | RoboMM (ours) | - | 90.7% |
| | RoboMM⁻ (ours) | - | 64.2% |
| RoboCasa | RoboCasa [48] | RSS'24 | 28.8% |
| | RoboMM (ours) | - | 30.6% |
| | RoboMM⁻ (ours) | - | 27.0% |
| Robomimic | IBC [19] | CoRL'21 | 13.6% |
| | RoboMM (ours) | - | 15.0% |
| | RoboMM⁻ (ours) | - | 8.0% |

Table 1. Performance on Various Datasets.

formance of MDT [55] and HULC++ [43], significantly exceeding the performance of models in the second tier. In the Meta-World dataset [65], *RoboMM* achieve a success rate of 78.6%, positioning it within the top tier of models. Furthermore, *RoboMM* consistently outperformed other models in LIBERO [34], RoboCasa [47], and Robomimic [39], demonstrating its robustness and adaptability across a variety of tasks and environments. These findings not only underscore *RoboMM*'s versatility but also highlight its effectiveness in handling tasks with varying levels of complexity, indicating its potential for real-world applications.

In contrast, *RoboMM⁻* reflects the results obtained from training on unaligned *RoboData*. The findings indicate that the misalignment between input and output spaces significantly hampers model performance when training across multiple datasets. It is important to note that in the Meta-World dataset [65], the impact of this misalignment is relatively minimal. This is because the actions involved are limited to three positional changes without rotation ($a^h_{pose} = \{\Delta\text{pos}^h_x, \Delta\text{pos}^h_y, \Delta\text{pos}^h_z, 0.0, 0.0, 0.0\}$), and the original spatial coordinate system aligns with that of the *RoboData* coordinate system. As a result, there is no need for additional alignment processes, allowing the model to perform adequately despite the lack of training on aligned data.

The successful alignment of input and output spaces within *RoboData*, combined with the compatible architecture of *RoboMM*, enables effective joint training and evaluation across multiple datasets. **It is important to note that, unlike *RoboMM*, a Generalist Policy, all other compared models are specialized and limited to training or fine-tuning on a single dataset.** While many studies in the field may involve training on multiple datasets, they often necessitate fine-tuning on the final evaluation dataset. This domain-specific fine-tuning can enhance performance but also increases training costs and may compromise the generalization capabilities of the models.

### 5.2. Ablation Study

**Protocol.** This experiment aims to systematically evaluate the contributions of each module in the *RoboMM* framework to overall performance, using the CALVIN [41] dataset for analysis. We design six different experimental setups to progressively analyze the impact of each module. First, in the baseline setup, the model receives $H$ frames as input and outputs an action only at the last frame to assess its basic prediction capability. Next, with the addition of the FFA module, the model is able to output actions after each frame, allowing us to examine its immediate prediction capability and performance in continuous prediction scenarios. As the experiment progresses, the introduction of the Image module enables the model to generate the next frame, thereby enhancing its contextual understanding and fine-grained perception. Subsequently, we added the UVFormer module to assess its specific contribution to the model's performance. Finally, by incorporating the occupancy (OCC) output, we further enhance the model's understanding of 3D spatial information, thereby improving its predictive capabilities.

**Experiment.** The results of the experiments are summarized in Table 2, which demonstrates a clear trend of performance improvement as different modules are added. The baseline configuration achieves a success rate of 81.0% for completing the first task, while the addition of the Frame by Frame Action (FFA) module increases this success rate to 85.0%. The introduction of image output capability fur-

| ID | FFA | Image | UVFormer | OCC | Task Completed in a Sequence | | | | | Avg Len |
|----|-----|-------|----------|-----|------|------|------|------|------|---------|
| | | | | | 1 | 2 | 3 | 4 | 5 | |
| 1 | ✗ | ✗ | ✗ | ✗ | 81.0% | 48.1% | 25.7% | 14.5% | 8.6% | 1.77 |
| 2 | ✓ | ✗ | ✗ | ✗ | 85.0% | 63.3% | 42.0% | 28.7% | 18.8% | 2.37 |
| 3 | ✓ | ✓ | ✗ | ✗ | 88.5% | 74.7% | 60.7% | 49.1% | 39.6% | 3.13 |
| 4 | ✓ | ✗ | ✓ | ✗ | 94.2% | 74.7% | 55.1% | 38.3% | 25.8% | 2.88 |
| 5 | ✓ | ✗ | ✓ | ✓ | 94.5% | 78.4% | 61.1% | 46.6% | 35.4% | 3.18 |
| 6 | ✓ | ✓ | ✓ | ✓ | 94.7% | 80.3% | 65.1% | 51.4% | 39.0% | 3.31 |

Table 2. Performance on different modules of *RoboMM*. FFA and OCC refer to frame by frame action module and occupancy supervision.

ther enhances success rates across all tasks. Notably, the UVFormer setup shows a significant improvement in the first task, achieving a success rate of 94.2%, indicating the effectiveness of the UVFormer module in enhancing the model's predictive capabilities. The OCC configuration builds on this by incorporating occupancy loss, resulting in even higher success rates, particularly in later tasks, suggesting that the model benefits from enriched contextual and spatial information.

Finally, row 6 shows the highest performance across all tasks, with a significant increase in the average length of completed tasks. Combining insights from Table 2 and Figure 7, it is evident that although the generated images and occupancy do not exhibit ideal quality, they significantly contribute to enhancing action performance. This indicates that even with suboptimal generation quality, the integration of visual and spatial information still has a positive impact on the execution of robotic tasks.

Overall, this ablation study highlights the contributions of each module within *RoboMM*, illustrating how enhancements can lead to significant performance gains in task completion.

### 5.3. Comparison with OpenVLA

To evaluate the superiority of our model architecture, we compare *RoboMM* with the currently best-performing OpenVLA. To ensure a fair comparison, we set the window size to 1, and train *RoboMM* from scratch, while OpenVLA is fine-tuned on the officially released weights.

As shown in Table 3, *RoboMM* outperforms OpenVLA (LoRA) in multiple metrics, particularly in Task 3-5 and average sequence length. This indicates that *RoboMM* can capture longer dependencies when handling tasks, thereby improving model accuracy. These results not only demonstrate the superior performance of the *RoboMM* architecture but also provide valuable references for future research.

| | Task Completed in a Sequence | | | | | Avg Len |
|---|------|------|------|------|------|---------|
| | 1 | 2 | 3 | 4 | 5 | |
| OpenVLA (LoRA) | 78% | 55% | 29% | 17% | 8% | 1.86 |
| RoboMM (ours) | 81% | 54% | 37% | 25% | 16% | 2.15 |

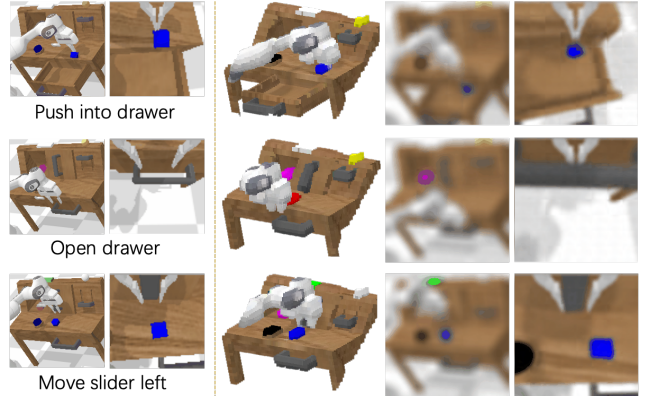Table 3. Performance comparison with OpenVLA on CALVIN.



Figure 7. The illustration demonstrating the auxiliary modality results of *RoboMM* on the Calvin three tasks. The inputs consist of static images, wrist-view images, and text, while the outputs include the current OCC and the predicted next frame's static and wrist-view images.

## 6. Conclusion

This paper presents *RoboMM* and *RoboData*, two innovative solutions designed to address key challenges in robotic learning. We have significantly improved the performance of robotic manipulation tasks by integrating multiple datasets and utilizing advanced model architectures.

*RoboMM* enhances 3D spatial perception through the introduction of camera parameter calibration and occupancy supervision, while also improving fine-grained perception capabilities via multimodal output. In parallel, *RoboData* integrates multiple industry datasets, supplements missing modalities such as depth and camera parameters, and ensures proper alignment of input and output spaces for better model performance.

Although our training objectives in supervised learning have reached only a moderate scale, *RoboMM* and *RoboData* open new avenues for cross-embodiment joint training and evaluation. We hope this work provides the industry with a comprehensive and fair evaluation system, inspires future exploration in foundational robotic model research, and enhances the generality and performance of robotic learning overall.

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ah-mad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 1, 3

[2] 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, 5

[3] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*, 2023. 1, 2, 4

[4] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. 5

[5] Kevin Black, Mitsuhiko Nakamoto, Pranav Atreya, Homer Walke, Chelsea Finn, Aviral Kumar, and Sergey Levine. Zero-shot robotic manipulation with pretrained image-editing diffusion models. *arXiv preprint arXiv:2310.10639*, 2023. 7, 5

[6] 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, 3, 7, 5

[7] 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. 1, 7, 5

[8] Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020. 1

[9] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023. 1

[10] Chi-Lam Cheang, Guangzeng Chen, Ya Jing, Tao Kong, Hang Li, Yifeng Li, Yuxiao Liu, Hongtao Wu, Jiafeng Xu, Yichu Yang, Hanbo Zhang, and Minzhao Zhu. Gr-2: A generative video-language-action model with web-scale knowledge for robot manipulation. *arXiv preprint arXiv:2410.06158*, 2024. 3

[11] Guiming Hardy Chen, Shunian Chen, Ruifei Zhang, Junying Chen, Xiangbo Wu, Zhiyi Zhang, Zhihong Chen, Jianquan Li, Xiang Wan, and Benyou Wang. Allava: Harnessing gpt4v-synthesized data for a lite vision-language model. *arXiv preprint arXiv:2402.11684*, 2024. 1

[12] Lili Chen, Shikhar Bahl, and Deepak Pathak. Playfusion: Skill acquisition via diffusion from language-annotated play. In *CoRL*, 2023. 7

[13] Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023. 1

[14] Lei Chen, Feng Yan, Yujie Zhong, Shaoxiang Chen, Zequn Jie, and Lin Ma. Mindbench: A comprehensive benchmark for mind map structure recognition and analysis. *arXiv preprint arXiv:2407.02842*, 2024. 1

[15] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, page 02783649241273668, 2023. 7, 5

[16] KR1442 Chowdhary and KR Chowdhary. Natural language processing. *Fundamentals of artificial intelligence*, pages 603–649, 2020. 1

[17] Ria Doshi, Homer Walke, Oier Mees, Sudeep Dasari, and Sergey Levine. Scaling cross-embodied learning: One policy for manipulation, navigation, locomotion and aviation. *arXiv preprint arXiv:2408.11812*, 2024. 3

[18] Hao-Shu Fang, Hongjie Fang, Zhenyu Tang, Jirong Liu, Junbo Wang, Haoyi Zhu, and Cewu Lu. Rh20t: A robotic dataset for learning diverse skills in one-shot. In *RSS 2023 Workshop on Learning for Task and Motion Planning*, 2023. 1, 3

[19] Pete Florence, Corey Lynch, Andy Zeng, Oscar Ramirez, Ayzaan Wahid, Laura Downs, Adrian Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. *Conference on Robot Learning (CoRL)*, 2021. 7, 5

[20] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022. 3

[21] Jiayuan Gu, Fanbo Xiang, Xuanlin Li, Zhan Ling, Xiqiang Liu, Tongzhou Mu, Yihe Tang, Stone Tao, Xinyue Wei, Yunchao Yao, et al. Maniskill2: A unified benchmark for generalizable manipulation skills. *arXiv preprint arXiv:2302.04659*, 2023. 2, 5, 6

[22] Meng-Hao Guo, Tian-Xing Xu, Jiang-Jiang Liu, Zheng-Ning Liu, Peng-Tao Jiang, Tai-Jiang Mu, Song-Hai Zhang, Ralph R Martin, Ming-Ming Cheng, and Shi-Min Hu. Attention mechanisms in computer vision: A survey. *Computational visual media*, 8(3):331–368, 2022. 1

[23] Yanjiang Guo, Yucheng Hu, Jianke Zhang, Yen-Jen Wang, Xiaoyu Chen, Chaochao Lu, and Jianyu Chen. Prediction

with action: Visual policy learning via joint denoising process. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 7, 5

[24] Huy Ha, Pete Florence, and Shuran Song. Scaling up and distilling down: Language-guided robot skill acquisition. In *Conference on Robot Learning*, pages 3766–3777. PMLR, 2023. 7, 5

[25] Yiyang Huang, Yuhui Hao, Bo Yu, Feng Yan, Yuxin Yang, Feng Min, Yinhe Han, Lin Ma, Shaoshan Liu, Qiang Liu, and Yiming Gan. Corki: Enabling real-time embodied ai robots via algorithm-architecture co-design, 2024. 1, 3

[26] Stephen James, Zicong Ma, David Rovick Arrojo, and Andrew J Davison. Rlbench: The robot learning benchmark & learning environment. *IEEE Robotics and Automation Letters*, 5(2):3019–3026, 2020. 2, 5, 6

[27] Xiaogang Jia, Qian Wang, Atalay Donat, Bowen Xing, Ge Li, Hongyi Zhou, Onur Celik, Denis Blessing, Rudolf Lioutikov, and Gerhard Neumann. Mail: Improving imitation learning with mamba. *arXiv preprint arXiv:2406.08234*, 2024. 7, 5

[28] Alexander Khazatsky, Karl Pertsch, Suraj Nair, Ashwin Balakrishna, Sudeep Dasari, Siddharth Karamcheti, Soroush Nasiriany, Mohan Kumar Srirama, Lawrence Yunliang Chen, Kirsty Ellis, et al. Droid: A large-scale in-the-wild robot manipulation dataset. *arXiv preprint arXiv:2403.12945*, 2024. 3

[29] Ji Woong Kim, Tony Z Zhao, Samuel Schmidgall, Anton Deguet, Marin Kobilarov, Chelsea Finn, and Axel Krieger. Surgical robot transformer (srt): Imitation learning for surgical tasks. *arXiv preprint arXiv:2407.12998*, 2024. 6

[30] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Foster, Grace Lam, Pannag Sanketi, et al. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 1, 3

[31] Seungjae Lee, Yibin Wang, Haritheja Etukuru, H. Jin Kim, Nur Muhammad Mahi Shafiullah, and Lerrel Pinto. Behavior generation with latent actions. *arXiv preprint arXiv:2403.03181*, 2024. 7, 5

[32] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 3

[33] Xinghang Li, Minghuan Liu, Hanbo Zhang, Cunjun Yu, Jie Xu, Hongtao Wu, Chilam Cheang, Ya Jing, Weinan Zhang, Huaping Liu, et al. Vision-language foundation models as effective robot imitators. *arXiv preprint arXiv:2311.01378*, 2023. 3, 7

[34] Bo Liu, Yifeng Zhu, Chongkai Gao, Yihao Feng, Qiang Liu, Yuke Zhu, and Peter Stone. Libero: Benchmarking knowledge transfer for lifelong robot learning. *Advances in Neural Information Processing Systems*, 36, 2024. 2, 5, 6, 7

[35] Fanfan Liu, Feng Yan, Liming Zheng, Chengjian Feng, Yiyang Huang, and Lin Ma. Robouniview: Visual-language model with unified view representation for robotic manipulaiton. *arXiv preprint arXiv:2406.18977*, 2024. 1, 3, 4

[36] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. 1, 4

[37] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024. 1, 3

[38] Arjun Majumdar, Karmesh Yadav, Sergio Arnaud, Jason Ma, Claire Chen, Sneha Silwal, Aryan Jain, Vincent-Pierre Berges, Tingfan Wu, Jay Vakil, et al. Where are we in the search for an artificial visual cortex for embodied intelligence? *Advances in Neural Information Processing Systems*, 36:655–677, 2023. 3

[39] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In *Conference on Robot Learning (CoRL)*, 2021. 2, 5, 6, 7

[40] Oier Mees, Lukas Hermann, and Wolfram Burgard. What matters in language conditioned robotic imitation learning over unstructured data. *IEEE Robotics and Automation Letters (RA-L)*, 7(4):11205–11212, 2022. 7

[41] 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. 2, 3, 5, 6, 7

[42] 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. 1, 3

[43] Oier Mees, Jessica Borja-Diaz, and Wolfram Burgard. Grounding language with visual affordances over unstructured data. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, London, UK, 2023. 7

[44] Meta. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edgemobile-devices/, 2024. Accessed on 2024-10-01. 3

[45] Atharva Mete, Haotian Xue, Albert Wilcox, Yongxin Chen, and Animesh Garg. Quest: Self-supervised skill abstractions for learning continuous control, 2024. 7, 5

[46] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. *arXiv preprint arXiv:2203.12601*, 2022. 3

[47] Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, Adeet Parikh, Aaron Lo, Abhishek Joshi, Ajay Mandlekar, and Yuke Zhu. Robocasa: Large-scale simulation of everyday tasks for generalist robots. *arXiv preprint arXiv:2406.02523*, 2024. 2, 5, 6, 7

[48] Soroush Nasiriany, Abhiram Maddukuri, Lance Zhang, Adeet Parikh, Aaron Lo, Abhishek Joshi, Ajay Mandlekar, and Yuke Zhu. Robocasa: Large-scale simulation of every-

day tasks for generalist robots. In *Robotics: Science and Systems (RSS)*, 2024. 7, 2, 5

[49] Anibal Ollero, Marco Tognon, Alejandro Suarez, Dongjun Lee, and Antonio Franchi. Past, present, and future of aerial robotic manipulators. *IEEE Transactions on Robotics*, 38(1): 626–645, 2021. 1

[50] OpenAI. Chatgpt, 2022. Computer software. 5

[51] 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. 1, 2, 3, 5

[52] Wilbert Pumacay, Ishika Singh, Jiafei Duan, Ranjay Krishna, Jesse Thomason, and Dieter Fox. The colosseum: A benchmark for evaluating generalization for robotic manipulation. *arXiv preprint arXiv:2402.08191*, 2024. 2, 5, 6

[53] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 3

[54] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024. 1

[55] Moritz Reuss, Ömer Erdinç Yağmurlu, Fabian Wenzel, and Rudolf Lioutikov. Multimodal diffusion transformer: Learning versatile behavior from multimodal goals. In *Robotics: Science and Systems*, 2024. 7, 5

[56] Rutav Shah, Roberto Martín-Martín, and Yuke Zhu. Mutex: Learning unified policies from multimodal task specifications. In *7th Annual Conference on Robot Learning*, 2023. 7, 5

[57] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Proceedings of The 6th Conference on Robot Learning*, pages 785–799. PMLR, 2023. 1

[58] Octo Model Team, Dibya Ghosh, Homer Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, et al. Octo: An open-source generalist robot policy. *arXiv preprint arXiv:2405.12213*, 2024. 1, 3

[59] 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. 1, 3, 5

[60] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 1

[61] Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaiming He. Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers. *arXiv preprint arXiv:2409.20537*, 2024. 1, 3

[62] Zhiqiang Wang, Hao Zheng, Yunshuang Nie, Wenjun Xu, Qingwei Wang, Hua Ye, Zhe Li, Kaidong Zhang, Xuewen Cheng, Wanxi Dong, Chang Cai, Liang Lin, Feng Zheng, and Xiaodan Liang. All robots in one: A new standard and unified dataset for versatile, general-purpose embodied agents, 2024. 5

[63] 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. 7, 5

[64] Hongtao Wu, Ya Jing, Chilam Cheang, Guangzeng Chen, Jiafeng Xu, Xinghang Li, Minghuan Liu, Hang Li, and Tao Kong. Unleashing large-scale video generative pre-training for visual robot manipulation. In *International Conference on Learning Representations*, 2024. 7, 5

[65] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pages 1094–1100. PMLR, 2020. 2, 5, 6, 7, 3

[66] Edwin Zhang, Yujie Lu, William Wang, and Amy Zhang. Language control diffusion: Efficiently scaling through space, time, and tasks. *arXiv preprint arXiv:2210.15629*, 2022. 7

[67] 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. 1, 3, 7, 5

[68] Liming Zheng, Feng Yan, Fanfan Liu, Chengjian Feng, Zhuoliang Kang, and Lin Ma. Robocas: A benchmark for robotic manipulation in complex object arrangement scenarios. *arXiv preprint arXiv:2407.06951*, 2024. 2, 3, 5

[69] Ruijie Zheng, Ching-An Cheng, Hal Daumé III au2, Furong Huang, and Andrey Kolobov. Prise: Learning temporal action abstractions as a sequence compression problem, 2024. 7, 5

[70] Hongkuan Zhou, Zhenshan Bing, Xiangtong Yao, Xiaojie Su, Chenguang Yang, Kai Huang, and Alois Knoll. Language-conditioned imitation learning with base skill priors under unstructured data. *IEEE Robotics and Automation Letters*, pages 1–8, 2024. 7

# *RoboMM*: All-in-One Multimodal Large Model for Robotic Manipulation

## Supplementary Material

## 7. *RoboMM* detailed information

In summary, the parameters used in the study are as follows: $H = 12$, $N = 3$, $H = 256$, $W = 256$, $L = 80$, $B = 80$, $P = 40$, $C = 1024$, $\lambda_{\text{image}} = 0.1$, $\lambda_{\text{occ}} = 0.1$, $\lambda_g = 0.01$, and $\lambda_{\text{rgb}} = 0.5$. The optimization strategy employs AdamW, while the learning rate schedule utilizes cosine annealing, with an initial learning rate of $10^{-4}$ and a termination rate of $10^{-6}$. The model is trained for a total of 10 epochs unless otherwise specified.

## 8. Action Representation

Different datasets have different methods for obtaining actions. For example, given the poses at two consecutive time steps $P^t = (p^t, r^t_{quat})$ and $P^{t+1} = (p^{t+1}, r^{t+1}_{quat})$, which are represented by 3D coordinates and quaternions, respectively.

### 8.1. Euler Angle Difference Method (EADM)

The Euler Angle Difference Method is a way to describe rotational transformations by calculating the difference in Euler angles between two poses (or orientations). The specific steps are as follows:

1. Convert the quaternions $r^t_{quat}$ and $r^{t+1}_{quat}$ to Euler angles $r^t_{euler}$ and $r^{t+1}_{euler}$, respectively.
2. Compute the differences in the 3D coordinates and Euler angles to obtain the action:

$$A_t = (p^{t+1} - p^t, r^{t+1}_{euler} - r^t_{euler}). \tag{12}$$

This method is intuitive and easy to understand, but it may encounter gimbal lock issues when dealing with large-angle rotations or multiple rotations.

### 8.2. Composite Rotation Matrix Method (CRMM)

The Composite Rotation Matrix Method describes complex rotational transformations by multiplying multiple rotation matrices. A rotation matrix is a linear algebra tool used to represent rotations in three-dimensional space. The specific steps are as follows:

1. Convert the quaternions $r^t_{quat}$ and $r^{t+1}_{quat}$ to rotation matrices $r^t_{matrix}$ and $r^{t+1}_{matrix}$, respectively.
2. Compute the composite rotation by multiplying the rotation matrices to obtain the action:

$$A_t = (p^{t+1} - p^t, r^{t+1}_{matrix} \cdot Inv(r^t_{matrix})) \tag{13}$$

This method is advantageous because it can conveniently handle any complex combination of rotations and avoids the gimbal lock problem.

### 8.3. Pose Composition Method (PCM)

The pose composition method is a way to describe the position and orientation of an object in space. By combining the poses at two consecutive time steps, complex motions can be described. The specific steps are as follows:

1. Convert the quaternions $r^t_{quat}$ and $r^{t+1}_{quat}$ to rotation matrices $r^t_{matrix}$ and $r^{t+1}_{matrix}$, respectively.
2. Combine the poses to obtain the action:

$$A_t = \left(Inv(R^t_{matrix}) \cdot (p^{t+1} - p^t), Inv(R^t_{matrix}) \cdot R^{t+1}_{matrix}\right) \tag{14}$$

This method is advantageous because it can conveniently describe and compute complex motions of objects in space and is widely used in robotics and computer vision.

## 9. *RoboData* detailed information

### 9.1. CALVIN Dataset

CALVIN is an open-source simulated benchmark specifically designed for learning long-horizon language-conditioned tasks in robotics. The dataset features four distinct environment splits, labeled A, B, C, and D. Each environment contains 6 hours of human-teleoperated recording data, resulting in over 2 million trajectories. However, only 1% of this data is annotated with language instructions, amounting to approximately 24,000 trajectories. Each environment split is uniquely configured with various objects and scenarios, allowing for comprehensive validation of the performance, robustness, and generality of the trained policies across different data combinations.

The benchmark utilizes a 7-degree-of-freedom (7-DOF) Franka Emika Panda robotic arm equipped with a parallel gripper. This robotic platform is enhanced with onboard sensors and captures images from two camera perspectives, enabling it to effectively execute complex sequences of language instructions. The coordinate system is based on the robot's body, represented as Right-Forward-Up, where the X-axis represents the right direction, the Y-axis denotes the forward direction, and the Z-axis indicates the upward direction.

For action representation, CALVIN employs EADM. To ensure that actions are appropriately scaled for network predictions, specific scaling factors are applied: 0.02 for the X, Y, and Z axes, and 0.05 for the pitch, roll, and yaw angles. The states of the gripper are represented using -1 for open and 1 for closed, facilitating clear action commands.

**Space Alignment:** *RoboData* includes all 34 distinct tasks, providing 20,000 episodes with language instructions

| Platform | Physics Engine | Robot | Coordinate (X-Y-Z) | Views | Camera Parameters | Action[a] Representation | Tasks | Episodes |
|---|---|---|---|---|---|---|---|---|
| CALVIN [41] | PyBullet | 7-DOF Franka | Right-Forward-Up | Static, Gripper | No | EADM | 34 | 20K |
| Meta-World [65] | MuJoCo | 4-DOF Sawyer | Right-Forward-Up | behindGripper, corner, corner2, corner3, topview, gripperPOV | No | None | 50 | 5K |
| Libero [34] | MuJoCo | 7-DOF Franka | Forward-Left-Up | frontview, birdview, agentview, sideview | No | CRMM | 130 | 6.5K |
| RoboMimic [39] | MuJoCo | 7-DOF Franka | Forward-Left-Up | agentview, robot0_eye_in_hand | No | CRMM | 8 | 1.6K |
| RoboCasa [48] | MuJoCo | 12-DOF Franka | Forward-Left-Up | center, left, right, frontview, eye_in_hand | No | CRMM | 100 | 5K |
| ManiSkill2 [21] | SAPIEN | 7-DOF Franka | Forward-Right-Down | base_camera, hand_camera | No | PCM | 20 | 30K |
| RoboCAS [68] | SAPIEN /Isaac | 7-DOF Franka | Forward-Left-Up | gripper_camera, base_camera, static_camera | Yes | Absolute | 3 | 7.3K |
| RLBench [26] | V-REP | 7-DOF Franka | Forward-Left-Up | left_shoulder, right_shoulder, wrist, front | Yes | Absolute | 18 | 1.8K |
| Colosseum [52] | PyRep | 7-DOF Franka | Forward-Left-Up | left_shoulder, right_shoulder, wrist, front | Yes | Absolute | 20 | 2K |

| Platform | Workspace $\mathbf{Min}_{[X,Y,Z]}$, $\mathbf{Max}_{[X,Y,Z]}$ | Action Space $\mathbf{Min}_{[X,Y,Z,Pitch,Roll,Yaw]}$, $\mathbf{Max}_{[X,Y,Z,Pitch,Roll,Yaw]}$ | Gripper (Open/Close) |
|---|---|---|---|
| CALVIN [41] | [-0.43, -0.57, 0.43], [0.37, -0.00, 0.80] | [-0.03, -0.03, -0.03, -6.28, -0.07, -6.27], [0.04, 0.02, 0.02, 6.28, 0.06, 6.28] | -1/1 |
| Meta-World [65] | [-0.50, -0.10, 0.12], [0.48, 0.41, 0.60] | [-1.00, -1.00, -1.00], [1.00, 1.00, 1.00] | 0.5/-0.5 |
| Libero [34] | [-0.24, -0.43, 0.01], [0.86, 0.57, 0.90] | [-0.93, -0.93, -0.93, -0.33, -0.37], [0.93, 0.93, 0.37, 0.37, 0.37] | 1/-1 |
| RoboMimic [39] | [-0.17, -0.40, 0.90], [0.33, 0.33, 1.29] | [-1.0, -1.0, -1.0, -0.55, -1.0, -1.0], [1.0, 1.0, 1.0, 0.72, 0.45, 1.0] | 1/-1 |
| RoboCasa [48] | [-0.81, -1.35, 0.70], [0.85, 0.75, 1.83] | [-1.0, -1.0, -1.0, -1.0, -1.0, -1.0], [1.0, 1.0, 1.0, 1.0, 1.0, 0.89] | 1/-1 |
| ManiSkill2 [21] | [-0.26, -0.79, -1.17], [0.85, 0.76, 0.00] | [-0.14, -0.15, -0.16, -0.09, -0.09, -0.09,], [0.17, 0.16, 0.15, 0.09, 0.09, 0.09] | -1/1 |
| RoboCAS [68] | [-0.70, -0.82, 0.062], [0.85, 0.67, 0.92] | [-0.04, -0.04, -0.04, -0.12, -0.10, 0.15], [0.03, 0.04, 0.03, 0.07, 0.09, 0.16, 0.08] | 0/0.08 |
| RLBench [26] | [-0.89, -0.72, 0.80], [0.56, 0.69, 1.89] | [-0.05, -0.04, -0.03, -1.0, -0.15, -1.0], [0.05, 0.06, 0.03, 1.0, 0.15, 1.0] | 0/1 |
| Colosseum [52] | [-0.68, -0.68, 0.83], [0.54, 0.70, 1.85] | [-0.04, -0.04, -0.04, -0.79, -0.12, -0.76], [0.03, 0.03, 0.04, 0.79, 0.35, 0.79] | 0/1 |

Table 4. Detailed information of CALVIN [41], Meta-World [65], LIBERO [34], Robomimic [39], RoboCAS [68], ManiSkill2 [21], RoboCasa [47], RLBench [26], and Colosseum [52].

for training. Action representations are regenerated using CRMM, and camera parameters are obtained through replay. Since the other input spaces are consistent with those predefined by *RoboData*, no alignment adjustments are necessary.

The dataset evaluates 1,000 unique instruction chains, focusing primarily on sequential task execution. In each experiment, the robotic agent successfully completes a series of up to five language instructions in succession. The agent can only proceed to the next instruction after successfully achieving the current task, establishing a clear dependency on the completion of prior actions.

## 9.2. Meta-World Dataset

Meta-World is a tabletop manipulation benchmark designed to facilitate the training and evaluation of robotic manipulation policies in a simulated environment. This dataset focuses on the reinforcement learning domain and does not release training data. The simulator includes six perspectives: behindGripper, corner, corner2, corner3, topview, gripperPOV.

The benchmark utilizes a 4-degree-of-freedom (4-DOF) Franka Emika Panda robotic arm equipped with a parallel gripper, which does not allow end rotation. The gripper states are represented by the numbers 0.5/-0.5 for open/close, and the coordinate system is consistent with that of the CALVIN dataset.

**Space Alignment:** *RoboData* includes the ML-45 version, which consists of 45 distinct tasks. To address the lack of training data for simulation learning, we adopt the scripted policies from Yu et al. [65] and introduce Gaussian noise $N(0, 0.1)$ to the generated actions at each step, resulting in a total of 22,500 trajectories, with each task producing 500 successful trajectories. To align with *RoboData*'s predefined settings, we extract observations from the corner2 and gripperPOV perspectives. The rotational variables in the actions are zero-padded, and the gripper states are represented using -1 for open and 1 for closed, while other parameters remain unchanged.

For performance evaluation, we test on 20 unseen start and goal configurations for each task, totaling 900 unseen configurations. We report the average performance over these 900 trajectories, providing a comprehensive measure of the model's ability to generalize to new tasks and configurations.

## 9.3. LIBERO Dataset

LIBERO is a lifelong learning benchmark that includes multiple task suites involving various language-labeled rigid and articulated-body manipulation tasks. The dataset consists of a total of 130 tasks and 6,500 trajectories. The simulator includes four perspectives: frontview, birdview, agentview, sideview, all with a resolution of 256 × 256

pixels. The action representation differs from that used in CALVIN, employing CRMM to define actions.

**Space Alignment:** *RoboData* includes the LIBERO-90 suite, which consists of 90 manipulation tasks, each with 50 demonstration trajectories collected through human teleoperation, providing a rich set of examples for training and evaluation. We select frontview and birdview as the observation perspectives, and camera parameters are obtained through replay. The coordinate system is defined as Forward-Left-Up. Due to differences in the coordinate system and workspace compared to the predefined settings in *RoboData*, we align them through rotation and translation:

$$W_{LIBERO} = \begin{bmatrix} 0 & 1 & 0 & 0.3 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & -0.1 \\ 0 & 0 & 0 & 1 \end{bmatrix} W_{LIBERO}^{ori}.$$

The states of the gripper are similarly represented using -1 for open and 1 for closed.

During evaluation, we test on 20 unseen start and goal configurations for each task, totaling 1,800 unseen configurations. This approach allows for a comprehensive assessment of the agent's performance and generalization capabilities, ensuring that the evaluation reflects the agent's ability to adapt to new situations and previously unencountered scenarios.

## 9.4. RoboMimic Dataset

RoboMimic is a large-scale robotic manipulation benchmark designed to study imitation learning and offline reinforcement learning. The dataset includes 5 distinct manipulation tasks, each with a dataset of demonstrations teleoperated by proficient humans. These tasks are designed to enhance the learning effectiveness of robots through real human demonstrations.

**Space Alignment:** *RoboData* includes 3 of these tasks (Lift, Can, Square) and excludes the other two dual-arm tasks. Given that the characteristics of RoboMimic align with those of LIBERO, all alignment methods can refer to LIBERO.

During evaluation, we test on 20 unseen start and goal configurations for each task, totaling 600 unseen configurations.

## 9.5. RoboCasa Dataset

RoboCasa is an open-source simulation benchmark designed to study robotic manipulation tasks in household environments, utilizing a 12-DOF Franka robot, where the first 7 degrees of freedom are related to manipulation and the remaining 5 are related to mobility. The dataset includes a simulation environment featuring 120 distinct real-world scenes, thousands of interactive objects, and household appliances, utilizing generative AI tools to create
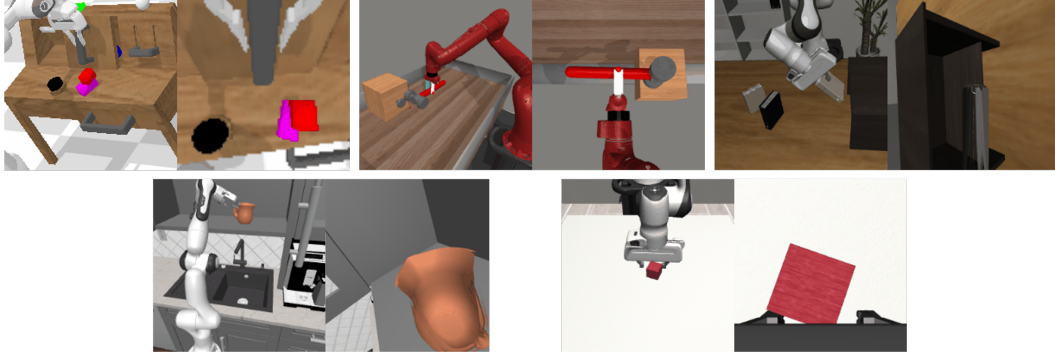
Figure 8. Evaluation Datasets. We evaluate *RoboMM* across five simulation benchmarks and present policy rollout visualizations of the experiments. From left to right, the benchmarks include CALVIN, Meta-World, LIBERO, RoboCasa, and Robomimic. Experiment details can be found in Section 5.1.

environmental textures and 3D objects. The RoboCasa dataset introduces 100 systematic evaluation tasks, consisting of 25 atomic tasks and 75 composite tasks generated with the guidance of large language models. Additionally, RoboCasa provides a large-scale multi-task dataset containing over 100,000 trajectories for model training, showcasing performance improvements achieved through behavior cloning training with synthetic data, as well as the applicability of simulation data in real-world tasks. These features make RoboCasa an important resource for researching and developing language-conditioned robotic technologies, laying a solid foundation for advancing intelligent applications of robots in household environments.

**Space Alignment:** *RoboData* includes 5,000 samples collected through remote control, utilizing two perspectives: front view and eye-in-hand. Only the degrees of freedom related to manipulation are retained. Given that the characteristics of RoboCasa align with those of LIBERO, all alignment methods can refer to LIBERO.

During evaluation, we test on 20 unseen start and goal configurations for each task, totaling 2,000 unseen configurations.

### 9.6. ManiSkill2 Dataset

ManiSkill2 is a unified benchmark designed for learning generalizable robotic manipulation skills, built on the SAPIEN platform. It includes 20 out-of-the-box task families, featuring over 2,000 distinct object models and more than 4 million demonstration frames. The dataset supports fast visual input learning algorithms, enabling a CNN-based policy to collect samples at approximately 2,000 frames per second (FPS) using just one GPU and 16 processes on a workstation. As the next generation of the SAPIEN ManiSkill benchmark, ManiSkill2 addresses critical pain points often encountered by researchers when utilizing benchmarks for developing generalizable manipulation skills, covering various task types, including station-

ary/mobile bases, single/dual-arm, and rigid/soft-body manipulation tasks. This extensive diversity of tasks and objects aims to enhance the robustness and applicability of robotic manipulation algorithms in real-world scenarios, making it an essential resource for advancing research in the field.

**Space Alignment:** *RoboData* includes 20 tasks related to single-arm manipulation from the ManiSkill2 dataset. The coordinate system and workspace are defined as Forward-Right-Down and [-0.26, -0.79, -1.17] to [0.85, 0.76, 0.00]. To ensure spatial consistency and compatibility, the corresponding coordinate transformations are applied:

$$W_{ManiSkill2} = \begin{bmatrix} 0 & 1 & 0 & 0.3 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} W_{ManiSkill2}^{ori}.$$

The action representation uses CRMM to replace PCM. Since the original data did not include the camera's intrinsic and extrinsic parameters, we replayed the data and saved the relevant parameters.

During evaluation, we test on 20 unseen start and goal configurations for each task, totaling 400 unseen configurations.

### 9.7. RoboCAS Dataset

RoboCAS is a benchmark proposed by Meituan's Embodied Intelligence Team, specifically designed for complex object arrangement scenarios in robotic manipulation. It is the first benchmark of its kind for such tasks and the first to employ a flexible and concise scripting strategy to collect samples in a cost-effective and efficient manner. RoboCAS showcases the handling of dispersed, ordered, and stacked objects within a highly realistic physical simulation environment, aiming to enhance robots' operational capabilities and performance across diverse settings. The benchmark provides a variety of proprioceptive observations and visual

data, including RGB images and depth maps captured from the left gripper camera, base camera, and static camera.

**Space Alignment:** *RoboData* includes all samples, utilizing only the base camera and static camera. The coordinate system and workspace are defined as Forward-Left-Up and [-0.70, -0.82, 0.062] to [0.85, 0.67, 0.92]. To ensure spatial consistency and compatibility, the following coordinate transformation is applied:

$$W_{RLBench} = \begin{bmatrix} 0 & 1 & 0 & 0.3 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0.7 \\ 0 & 0 & 0 & 1 \end{bmatrix} W_{RLBench}^{ori}.$$

Since only end-effector positions are provided in the dataset, the research team utilized a composite rotation matrix to generate corresponding action representations, changing the gripper's open/close state from 0/0.04 to -1/1. Notably, the RGB images from this perspective are 480x640 pixels; to maintain consistency across all data in *RoboData*, we only extract the central region of 480x480 pixels.

During evaluation, we test on 20 unseen start and goal configurations for each task.

### 9.8. RLBench Dataset

RLBench is a challenging benchmark and learning environment specifically designed for robot learning. This benchmark features 18 completely unique, hand-designed tasks that range in difficulty from simple target reaching and door opening to more complex multi-stage tasks, such as opening an oven and placing a tray inside. RLBench provides a variety of proprioceptive observations and visual observation data, including RGB images, depth maps, and segmentation masks from the left shoulder, right shoulder, wrist, and front views.

**Space Alignment:** *RoboData* includes all experiments, totaling 1.8 experiments, with visual input extracted from the wrist and front views. The coordinate system in this dataset differs from that of other datasets, defined as Forward-Left-up, with a workspace range from [-0.89, -0.72, 0.80] to [0.56, 0.69, 1.89]. We apply spatial transformations to shift the data into a predefined coordinate system.

$$W_{RLBench} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0.7 \\ 0 & 0 & 0 & 1 \end{bmatrix} W_{RLBench}^{ori}.$$

Additionally, only end-effector positions are provided, so we use CRMM to transform action representations, changing the gripper's open/close state from 0/1 to -1/1.

During evaluation, we test on 20 unseen start and goal configurations for each task, totaling 360 unseen configurations.

### 9.9. Colosseum Dataset

Colosseum is a benchmark that complements RLBench by addressing the limitations of single environmental variables. It features 20 diverse manipulation tasks that enable systematic evaluation of models across 14 axes of environmental perturbations. These perturbations include changes in the color, texture, and size of objects, as well as variations in tabletop surfaces, backgrounds, and the physical properties of objects. Additionally, lighting conditions, distractors, and camera poses are adjusted. All configurations align with those of RLBench, allowing researchers to test and compare the robustness and adaptability of their models under a wider range of environmental conditions.

**Space Alignment:** *RoboData* includes all samples, and the alignment method is consistent with that of RLBench.

## 10. Experiment Details

The success rates of the expert models in the Table 1 are organized from the following sources: The evaluation methods on the CALVIN [41] dataset are sourced from the official CALVIN leaderboard (url: http://calvin.cs.uni-freiburg.de/). In the Meta-World [65] dataset, the results of PAD [23], GR-1 [64], SuSIE [5], RT-2* [7], and RT-1 [6] come from PAD [23], while the results of PRISE [69] are derived from related papers. In the Libero [34] dataset, the results of QueST [45], VQ-BeT [31], PRISE [69], DiffusionPolicy [15], ACT [67], and ResNet-T [34] all come from QueST [45], while the results of MDT [55] and Distill-D [24] are sourced from MDT [55]; the results of MaIL [27], ATM [63], and MUTEX [56] come from their respective research papers. The results of RoboCasa [48] in the RoboCasa [48] dataset are sourced from related papers. In the Robomimic [39] dataset, the results of IBC [19] come from personal replication.

The success rates of each task for Robomimic on various datasets are shown in the Table 5, 6, 7, 8, 9.

| Task | Success Rate |
|---|---|
| Rotate Blue Block Right | 87.9% |
| Move Slider Right | 86.8% |
| Turn Off LED | 87.7% |
| Push Into Drawer | 45.2% |
| Lift Blue Block Drawer | 93.3% |
| Lift Pink Block Slider | 56.2% |
| Place In Slider | 43.4% |
| Open Drawer | 90.1% |
| Rotate Red Block Right | 81.2% |
| Lift Pink Block Table | 82.9% |
| Push Blue Block Left | 88.3% |
| Close Drawer | 86.9% |
| Turn On LED | 92.4% |
| Stack Block | 39.2% |
| Push Pink Block Right | 75.4% |
| Rotate Blue Block Left | 89.7% |
| Lift Blue Block Table | 86.7% |
| Place In Drawer | 86.2% |
| Turn On Lightbulb | 81.9% |
| Move Slider Left | 84.8% |
| Rotate Red Block Left | 90.0% |
| Lift Red Block Slider | 62.5% |
| Push Pink Block Left | 86.4% |
| Push Red Block Left | 70.0% |
| Lift Blue Block Slider | 64.8% |
| Push Red Block Right | 69.7% |
| Lift Red Block Table | 82.3% |
| Turn Off Lightbulb | 90.7% |
| Lift Pink Block Drawer | 70.0% |
| Rotate Pink Block Right | 56.7% |
| Rotate Pink Block Left | 83.3% |
| Push Blue Block Right | 41.0% |
| Unstack Block | 78.9% |
| Lift Red Block Drawer | 85.7% |

Table 5. *RoboMM* Success Rates on Various Tasks in CALVIN [41].

| Task | Success Rate (%) |
|---|---|
| Assembly | 65% |
| Basketball | 100% |
| Bin Picking | 75% |
| Box Close | 75% |
| Button Press Topdown | 100% |
| Button Press Topdown Wall | 100% |
| Button Press | 100% |
| Button Press Wall | 90% |
| Coffee Button | 70% |
| Coffee Pull | 40% |
| Coffee Push | 70% |
| Dial Turn | 90% |
| Disassemble | 50% |
| Door Close | 100% |
| Door Lock | 100% |
| Door Open | 100% |
| Door Unlock | 100% |
| Hand Insert | 70% |
| Drawer Close | 100% |
| Drawer Open | 100% |
| Faucet Open | 0% |
| Faucet Close | 100% |
| Hammer | 35% |
| Handle Press Side | 100% |
| Handle Press | 100% |
| Handle Pull Side | 20% |
| Handle Pull | 50% |
| Lever Pull | 100% |
| Peg Insert Side | 65% |
| Pick Place Wall | 100% |
| Pick Out of Hole | 50% |
| Reach | 50% |
| Push Back | 80% |
| Push | 95% |
| Pick Place | 85% |
| Plate Slide | 100% |
| Plate Slide Side | 100% |
| Plate Slide Back | 100% |
| Plate Slide Back Side | 100% |
| Peg Unplug Side | 25% |
| Soccer | 45% |
| Stick Push | 100% |
| Stick Pull | 85% |
| Push Wall | 95% |
| Reach Wall | 90% |
| Shelf Place | 20% |
| Sweep Into | 95% |
| Sweep | 65% |
| Window Open | 85% |
| Window Close | 100% |

Table 6. *RoboMM* Success Rates on Various Tasks in Meta-World [65]

| Task | Success Rate |
|---|---|
| Lift | 45% |
| Can | 0 |
| Square | 0 |

Table 7. *RoboMM* Success Rates on Various Tasks in RoboMimic [39].

| Task Index | Success Rate (%) | Task Index | Success Rate (%) |
|---|---|---|---|
| 0 | 100% | 1 | 40% |
| 2 | 100% | 3 | 95% |
| 4 | 95% | 5 | 100% |
| 6 | 85% | 7 | 100% |
| 8 | 100% | 9 | 100% |
| 10 | 100% | 11 | 100% |
| 12 | 75% | 13 | 95% |
| 14 | 100% | 15 | 100% |
| 16 | 90% | 17 | 85% |
| 18 | 100% | 19 | 85% |
| 20 | 100% | 21 | 100% |
| 22 | 100% | 23 | 55% |
| 24 | 100% | 25 | 95% |
| 26 | 90% | 27 | 95% |
| 28 | 100% | 29 | 100% |
| 30 | 65% | 31 | 90% |
| 32 | 30% | 33 | 75% |
| 34 | 100% | 35 | 100% |
| 36 | 90% | 37 | 100% |
| 38 | 90% | 39 | 100% |
| 40 | 80% | 41 | 100% |
| 42 | 100% | 43 | 85% |
| 44 | 100% | 45 | 100% |
| 46 | 100% | 47 | 100% |
| 48 | 95% | 49 | 100% |
| 50 | 100% | 51 | 20% |
| 52 | 100% | 53 | 85% |
| 54 | 100% | 55 | 100% |
| 56 | 100% | 57 | 100% |
| 58 | 100% | 59 | 100% |
| 60 | 100% | 61 | 95% |
| 62 | 95% | 63 | 100% |
| 64 | 70% | 65 | 100% |
| 66 | 100% | 67 | 85% |
| 68 | 100% | 69 | 100% |
| 70 | 100% | 71 | 100% |
| 72 | 90% | 73 | 75% |
| 74 | 100% | 75 | 55% |
| 76 | 90% | 77 | 95% |
| 78 | 90% | 79 | 100% |
| 80 | 75% | 81 | 25% |
| 82 | 100% | 83 | 80% |
| 84 | 90% | 85 | 75% |
| 86 | 100% | 87 | 100% |
| 88 | 100% | 89 | 95% |

Table 8. *RoboMM* Success Rates on Various Tasks in Libero [34].

| Task Name | Success Rate (%) |
|---|---|
| CoffeePressButton | 90% |
| CoffeeServeMug | 55% |
| CoffeeSetupMug | 25% |
| CloseDoubleDoor | 5% |
| CloseSingleDoor | 90% |
| OpenDoubleDoor | 0% |
| OpenSingleDoor | 45% |
| CloseDrawer | 95% |
| OpenDrawer | 50% |
| TurnOffMicrowave | 35% |
| TurnOnMicrowave | 70% |
| PnPCabToCounter | 30% |
| PnPCounterToCab | 30% |
| PnPCounterToMicrowave | 15% |
| PnPCounterToSink | 30% |
| PnPCounterToStove | 30% |
| PnPMicrowaveToCounter | 30% |
| PnPSinkToCounter | 5% |
| PnPStoveToCounter | 20% |
| TurnOffSinkFaucet | 70% |
| TurnOnSinkFaucet | 70% |
| TurnSinkSpout | 90% |
| TurnOffStove | 30% |
| TurnOnStove | 60% |
| PrepareCoffee | 0% |
| ArrangeVegetables | 0% |
| MicrowaveThawing | 0% |
| RestockPantry | 0% |
| PreSoakPan | 0% |

Table 9. *RoboMM* Success Rates on Various Tasks in Robo-Casa [48].