

Rethinking Bimanual Robotic Manipulation: Learning with Decoupled Interaction Framework

Jian-Jian Jiang¹, Xiao-Ming Wu¹, Yi-Xiang He¹, Ling-An Zeng¹,
Yi-Lin Wei¹, Dandan Zhang², Wei-Shi Zheng^{1†}

¹ School of Computer Science and Engineering, Sun Yat-sen University, China

² Imperial-X Initiative and Department of Bioengineering, Imperial College London, U.K.

Abstract

Bimanual robotic manipulation is an emerging and critical topic in the robotics community. Previous works primarily rely on integrated control models that take the perceptions and states of both arms as inputs to directly predict their actions. However, we think bimanual manipulation involves not only coordinated tasks but also various uncoordinated tasks that do not require explicit cooperation during execution, such as grasping objects with the closest hand, which integrated control frameworks ignore to consider due to their enforced cooperation in the early inputs. In this paper, we propose a novel decoupled interaction framework that considers the characteristics of different tasks in bimanual manipulation. The key insight of our framework is to assign an independent model to each arm to enhance the learning of uncoordinated tasks, while introducing a selective interaction module that adaptively learns weights from its own arm to improve the learning of coordinated tasks. Extensive experiments on seven tasks in the RoboTwin dataset demonstrate that: (1) Our framework achieves outstanding performance, with a 23.5% boost over the SOTA method. (2) Our framework is flexible and can be seamlessly integrated into existing methods. (3) Our framework can be effectively extended to multi-agent manipulation tasks, achieving a 28% boost over the integrated control SOTA. (4) The performance boost stems from the decoupled design itself, surpassing the SOTA by 16.5% in success rate with only 1/6 of the model size.

1. Introduction

Bimanual robotic manipulation [2, 8, 11, 13, 38, 41], which has strong capabilities to handle a wide range of complex tasks such as household services [38], medical surgery [13], health care [41] and industrial assembly [11], is an emerging and critical topic in the robotics community.

Recently, with the help of imitation learning [9, 31],

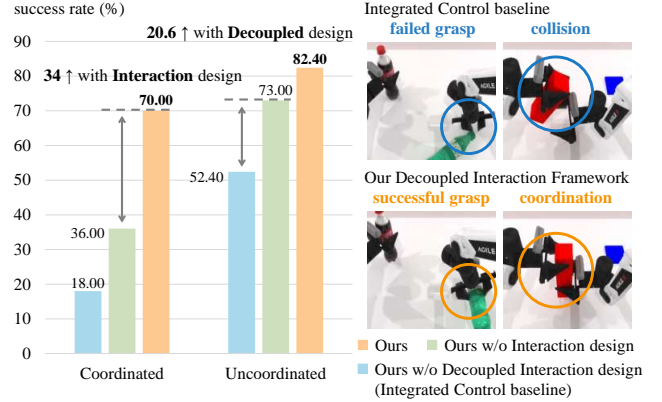


Figure 1. **Integrated Control vs. Decoupled Interaction.** The blue bar represents the success rate of coordinated and uncoordinated tasks for the integrated control baseline built upon our framework without decoupled interaction design. The green and orange bars represent the success rate of coordinated and uncoordinated tasks for our framework without interaction design and our framework respectively. Our experiments are conducted on two coordinated tasks and five uncoordinated tasks in the RoboTwin dataset [27]. It can be observed that adding the decoupled design to the integrated control baseline promotes the learning of uncoordinated tasks. Furthermore, incorporating the interaction module on top of this design facilitates the learning of coordinated tasks.

bimanual robotic manipulation [2, 8, 16, 21, 40] achieves significant progress. Previous works typically use an integrated control model that takes observations and states of both arms as inputs and predicts actions for both arms simultaneously. However, based on taxonomy research on bimanual manipulation [15], we think that bimanual robotic manipulation is more complex, involving not only coordinated tasks but also various uncoordinated tasks. In **uncoordinated tasks**, the two arms are neither spatially nor temporally coordinated and do not share directly connected goals. Each arm fulfills its task-specific constraints, with spatial coordination limited to avoiding collisions and no temporal coupling. For example, one arm holds a coffee

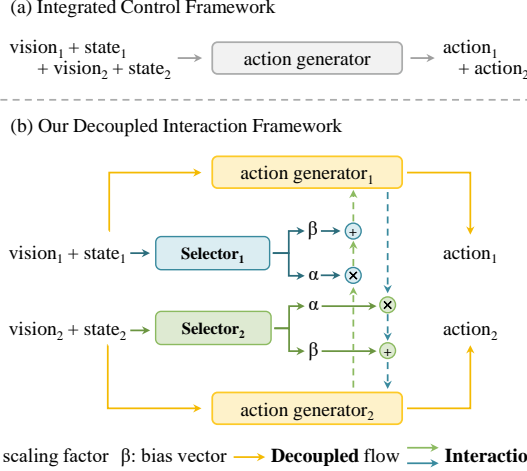


Figure 2. **Comparisons of our Decoupled Interaction Framework with integrated control frameworks.** Integrated control frameworks (a) mainly use a single model that takes the observations and states of both arms as inputs and directly outputs their actions. Our Decoupled Interaction Framework (b) first assigns an independent model to each arm to solely handle the inputs of the current arm (the yellow lines). Then, different from the naive interaction modeling in integrated control frameworks, a selective interaction module is proposed which learns its weights from its own arm to perform explicit modeling (the green and blue lines) on the exchanged state features (the green and blue dashed lines).

cup while the other takes notes. In contrast, in **coordinated tasks**, both arms require spatial or temporal coordination, defined by specific task constraints. For instance, one arm places a block at a designated location, and the other positions another block on top afterward. Due to the unique characteristics of these tasks and the neglect of integrated control models to account for these differences, these models struggle to effectively learn uncoordinated and coordinated tasks. We conduct some confirmatory experiments using an integrated control baseline, as illustrated in Fig. 1, which verifies our claim.

To address this, we propose learning bimanual manipulation via a precise task taxonomy [15] and introduce a novel Decoupled Interaction Framework for bimanual robotic manipulation. Compared to integrated control frameworks (as illustrated in Fig. 2 (a)), which directly predict actions for both arms, our framework (as illustrated in Fig. 2 (b)) first decouples the action dimensions the model needs to predict. Subsequently, unlike integrated control frameworks that primarily use the states of both arms for joint encoding to model the interaction between the two arms, our framework performs explicit interaction modeling between the two arms. Moreover, cooperation requirements vary at different stages, our framework selectively modulates interaction information during execution, enabling more effective use of this information. As illustrated in Fig. 1, the integrated control baseline struggles to learn both uncoordi-

nated and coordinated tasks, while incorporating the decoupled design and interaction design into the baseline effectively promotes learning of uncoordinated and coordinated tasks respectively. Notably, the decoupled design also facilitates learning of coordinated tasks, as it enhances the success rate of non-cooperative phases in coordinated tasks. Meanwhile, the interaction design also facilitates learning of uncoordinated tasks, as uncoordinated tasks still require minimal cooperation to avoid collisions between arms.

Specifically, our framework first assigns an independent model to each arm and takes the observation and state of a single arm as inputs to generate actions only for that arm. This design is beneficial for uncoordinated tasks, where the action intentions and execution of each arm are relatively independent, enabling each arm to complete its own operations more effectively. Building upon this strong capability to handle uncoordinated tasks, our framework further employs explicit interaction modeling to improve cooperation, thereby facilitating the learning of coordinated tasks. To this end, our framework selects state features as interactive information and introduces a selective interaction module to adaptively modulate the exchanged features. First, our method uses a selective scaling module to adjust the intensity of the received state features. Then, to align the received state features with the current arm state space, a selective alignment module is introduced to adaptively generate a bias vector for the received state features. With our selective interaction module, our method satisfies the cooperation requirements of coordinated tasks.

Extensive experiments on seven tasks in the RoboTwin dataset [27] demonstrate the following: (1) **Effectiveness**: Our framework achieves outstanding performance, obtaining a **23.5%** improvement over the previous SOTA method [37]. (2) **Flexibility**: Our framework can be seamlessly integrated into existing methods, such as DP3 [37] and Point Flow Matching (PFM) [4]. (3) **Extensibility**: Our framework can be effectively extended to multi-agent manipulation tasks, achieving a **28%** improvement over the SOTA method [37]. (4) **Scalability**: Experiments with an increased number of expert demonstrations confirm that our framework adheres to the scaling law. It should be emphasized that the performance improvement stems from the decoupled interaction design itself, surpassing the SOTA by **16.5%** in performance with only **1/6** of the model size. We also conduct real-world experiments to further validate the effectiveness of our approach. Upon the publication of this paper, we will open-source the code.

2. Revisiting Bimanual Robotic Manipulation

Robotic manipulation aims to utilize extensive observation-action pairs from expert demonstrations to acquire human-like skills. Previous research mainly focuses on single-arm manipulation, where the modeling of policies is primarily

divided into reinforcement learning [6, 10, 24, 33, 35, 39], imitation learning [1, 3–5, 12, 18, 23, 25, 29, 32, 36, 37] and Vision-Language-Action (VLA) models [7, 14, 17, 42]. Concretely, imitation learning methods mainly focus on regression-based [1, 5, 23, 29, 32, 36] and generation-based methods [3, 4, 12, 18, 25, 37], which lay the foundation for policy designs of bimanual manipulation.

Recently, bimanual manipulation becomes an emerging hotspot in the robotics community due to its ability to handle a wide range of tasks. First, we define this task. Given the observations O_{pcd} (point clouds as used in this paper) captured by depth cameras and the states of each arm S_{arm} provided by the ROS interface, bimanual manipulation aims to generate actions A_{arm} for the left and right arms to accomplish different categories of tasks.

Benefiting from the development of imitation learning [9, 31], bimanual manipulation [2, 8, 16, 20, 21, 40] methods achieve great progress. Previous works [2, 8, 16, 21, 40] mainly employ an integrated control framework that takes the observations and states of both arms as inputs and directly outputs their actions, i.e.:

$$\begin{aligned} A_{ic} &= \pi_{ic}(O_{pcd1}, O_{pcd2}, S_{arm1}, S_{arm2}), \\ A_{ic} &\in \mathbb{R}^{1 \times 14}, O_{pcd} \in \mathbb{R}^{N \times 3}, S_{arm} \in \mathbb{R}^{1 \times 7}. \end{aligned} \quad (1)$$

However, the high-dimensional task spaces hinder each arm from learning its own actions. This affects the efficiency for learning uncoordinated tasks, which focus more on whether each arm can complete its own actions. Moreover, the interaction modeling of integrated control frameworks is simple, where they typically perform joint encoding of the states of both arms. The implicit interaction modeling makes it difficult to learn coordinated tasks, where the cooperation requirements vary at different stages during execution.

A few works like Voxact-b [20] utilize a fully decoupled framework for bimanual robotic manipulation. These works often take the observation and state of a single arm as inputs and output the actions of the current arm, i.e.:

$$\begin{aligned} A_{arm1} &= \pi_{dec}(O_{pcd1}, S_{arm1}), \\ A_{arm2} &= \pi_{dec}(O_{pcd2}, S_{arm2}), \\ A_{arm} &\in \mathbb{R}^{1 \times 7}, O_{pcd} \in \mathbb{R}^{N \times 3}, S_{arm} \in \mathbb{R}^{1 \times 7}. \end{aligned} \quad (2)$$

These methods are usually applied to uncoordinated bimanual manipulation tasks because they do not explicitly model the collaboration between arms, making them less effective in handling coordinated tasks.

In this paper, we propose to learn bimanual robotic manipulation through precise task taxonomy [15] and further introduce a Decoupled Interaction Framework. The insight of our framework is to utilize the decoupled design to reduce the action dimensions, thereby promoting the learning of uncoordinated tasks, while introducing a selective interaction module that adaptively modulates the interaction

information to enhance the learning of coordinated tasks. Our experiments also comprehensively validate the excellent performance of our framework.

3. A Decoupled Interaction Framework for Bimanual Robotic Manipulation

In this section, we describe how our Decoupled Interaction Framework addresses the challenges in integrated control frameworks for bimanual manipulation. First, we revisit the motivation of our method, emphasizing that integrated control frameworks face issues of high-dimensional actions and a lack of explicit interaction modeling (Sec. 3.1). Next, we explain how our framework resolves these problems (Sec. 3.2 and Sec. 3.3). Finally, we summarize the overall structure of our framework (Sec. 3.4).

3.1. Challenges in Integrated Control Frameworks

The action intentions and execution of each arm are relatively independent in uncoordinated tasks. Therefore, it is important for methods to help each robotic arm effectively learn its own actions. However, most previous works do not adequately address this issue, since they mainly use an integrated control model to generate actions for both arms and combine the joint angle-based action spaces of each arm. This approach increases the dimensions of the actions the model needs to predict and enlarges the action space the neural network must explore, making the action learning processes for each arm more challenging.

Meanwhile, although cooperation between arms is important for coordinated tasks, this does not mean that both arms are constantly collaborating throughout the entire task. In other words, even in coordinated tasks, there are phases where each arm performs its own actions independently and phases where collaboration takes place. However, previous integrated control methods fail to effectively model the interaction between the arms, as they often rely on the states of both arms for joint encoding. This implicit modeling cannot distinguish different phases of coordinated tasks.

To validate it, we conduct experiments using our decoupled interaction framework. As shown in Fig. 1, the integrated control baseline exhibits low performance, resulting in failed grasps and collisions between targets.

3.2. The Proposed Decoupled Design Contributes to Uncoordinated Tasks

To address the problem of high-dimensional action spaces, in this paper, we propose to decouple the joint action space of both arms to reduce the space the neural network needs to explore, which helps our framework learn uncoordinated tasks more effectively. Specifically, our framework assigns an independent model to each arm, where it takes the observation and state of the current arm as inputs and outputs

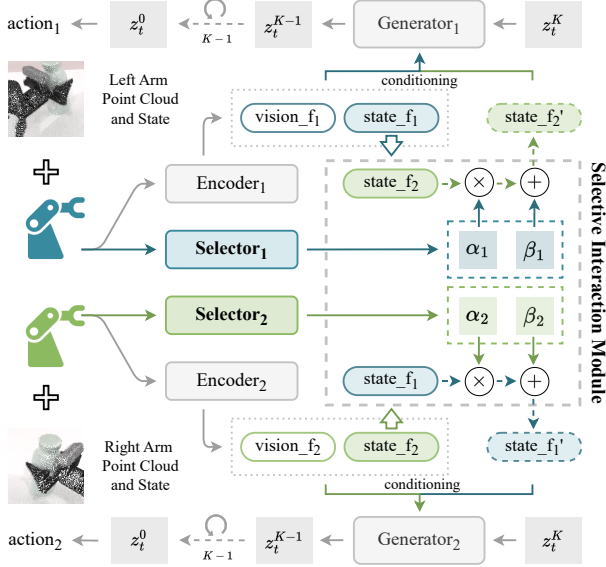


Figure 3. **Architecture of the Decoupled Interaction Framework.** Our framework first assigns a separate model to each arm to process its inputs. Then, we exchange state features between the models and utilize a selective interaction module to modulate them. Specifically, we use a selector to predict a scaling factor α and a bias vector β to adaptively adjust the exchanged features. Finally, we combine the original visual features, state features and exchanged state features as interactive conditions to predict actions using action generators.

actions for the current arm. As shown in Fig. 1, by incorporating the decoupled design, our framework outperforms the integrated control baseline in uncoordinated tasks. Moreover, as illustrated in Fig. 1, our framework can effectively handle basic tasks like single-arm grasping, which is important in uncoordinated tasks. In conclusion, we consider that the decoupled design is beneficial for uncoordinated tasks. The loss of our decoupled design is:

$$\begin{aligned} \mathcal{L}_{total} &= \mathcal{L}_{arm1} + \mathcal{L}_{arm2}, \\ \mathcal{L}_{arm1} &= \mathcal{L}_{action}(\theta_1, (O_{pcd1}, S_{arm1})), \\ \mathcal{L}_{arm2} &= \mathcal{L}_{action}(\theta_2, (O_{pcd2}, S_{arm2})). \end{aligned} \quad (3)$$

where \mathcal{L}_{action} represents the action generation loss, which will be described in Sec. 3.4. θ_1 and θ_2 represent the action generator for the left and right arm respectively.

It should be noted that although uncoordinated tasks require low cooperation, they still necessitate a certain level of interaction, i.e., the ability to perceive the state of the other arm to avoid collisions between the arms.

3.3. The Proposed Interaction Design Contributes to Coordinated Tasks

To improve explicit interaction modeling, we build on the ability of our framework to handle uncoordinated tasks and introduce a novel selective interaction module. This mod-

ule enhances interaction modeling and facilitates the learning of coordinated tasks. The key insight of our interaction module is to selectively modulate the state representations of the other arm by utilizing the observations and states of the current arm. This approach enables our framework to effectively control the influence of the other arm on each arm when each arm learns its own actions. As illustrated in Fig. 1, by incorporating the selective interaction design, our framework outperforms the integrated control baseline in coordinated tasks. Moreover, as illustrated in Fig. 1, our framework performs actions with better coordination.

Scale the Intensity of Exchanged State Features. Specifically, we exchange the encoded states between different arms and introduce a selective scaling module to scale the intensity of the received state representations. This module first processes the inputs of the current arm using modulation visual and state encoders. The encoded representations are then passed into an MLP, which predicts a scaling factor α ranging from 0 to 1. Through this scaling factor α , our framework adaptively adjusts the influence of the state information from the other arm on the current arm across different phases of coordinated tasks.

Align the Exchanged State Features. Following that, to further align the received state representations with the state space of the current arm, we introduce a selective alignment module, which dynamically generates a bias vector for the received state representations. This module first takes the encoded observations and states from modulation visual and state encoders as inputs and then utilizes an MLP to predict a bias vector β for alignment.

To summarize, with this selective interaction module, the current arm can effectively filter and align the exchanged information from the other arm, thereby meeting the coordination requirements for both coordinated and uncoordinated tasks more effectively. The workflow of our selective interaction module is shown as follows:

$$\begin{aligned} \alpha_1, \beta_1 &= f_1(O_{pcd1}, S_{arm1}), g_1(O_{pcd1}, S_{arm1}), \\ \alpha_2, \beta_2 &= f_2(O_{pcd2}, S_{arm2}), g_2(O_{pcd2}, S_{arm2}), \\ F'_{arm1} &= \alpha_2 \times F_{arm1} + \beta_2, \\ F'_{arm2} &= \alpha_1 \times F_{arm2} + \beta_1. \end{aligned} \quad (4)$$

where F_{arm1} and F_{arm2} denote the exchanged state features, $f(\cdot)$ represents the regression head of scaling factors and $g(\cdot)$ represents the regression head of bias vectors. On the basis of the interaction design, the loss of our decoupled interaction framework can be updated as:

$$\begin{aligned} \mathcal{L}_{total} &= \mathcal{L}_{arm1} + \mathcal{L}_{arm2}, \\ \mathcal{L}_{arm1} &= \mathcal{L}_{action}(\theta_1, (O_{pcd1}, S_{arm1}, F'_{arm2})), \\ \mathcal{L}_{arm2} &= \mathcal{L}_{action}(\theta_2, (O_{pcd2}, S_{arm2}, F'_{arm1})). \end{aligned} \quad (5)$$

Methods	Coordinated Tasks		Uncoordinated Tasks					Average
	Block Handover	Blocks Stack	Dual Bottles Pick (Easy)	Dual Bottles Pick (Hard)	Diverse Bottles Pick	Empty Cup Place	Block Hammer Beat	
Voxact-b [20]	0.440	0.290	0.820	0.490	0.400	0.770	0.650	0.551 (↓ 0.238)
ACT [40]	0.070	0.020	0.340	0.040	0.020	0.310	0.310	0.159 (↓ 0.630)
DP [3]	0.280	0.020	0.540	0.280	0.000	0.200	0.000	0.189 (↓ 0.600)
DP3 [37]	0.700	0.230	0.780	0.460	0.380	0.730	0.600	0.554 (↓ 0.235)
Ours	1.000	0.400	0.990	0.630	0.700	0.900	0.900	0.789

Table 1. **Performance comparison of seven simulation tasks in the RoboTwin dataset.** Best results are highlighted in bold. Important comparison metrics are marked with gray cells. The red arrows indicate the performance difference between each baseline and our method.

3.4. A Decoupled Interaction Framework

Loss Design. Our Decoupled Interaction Framework can be seamlessly integrated into existing methods, such as DP3 [37] and Point Flow Matching [4]. For DP3, we use DDIM [30] as the action generator, and the ground truth is represented by the sample. For Point Flow Matching, we use Flow Matching [19] as the action generator, and the ground truth is represented by the vector field. When employing DDIM, \mathcal{L}_{action} is represented as:

$$\mathcal{L}_{action}(\theta, I) = \mathbb{E}_{t, z \sim \mathcal{D}(x)} \|z - \theta(z_t, t|I)\|^2. \quad (6)$$

where $z \sim \mathcal{D}(x)$ denotes the distribution of the sample in demonstrations, I denotes the interactive condition and θ denotes the DDIM action generator of each arm. When applying Flow Matching, \mathcal{L}_{action} is represented as:

$$\begin{aligned} \mathcal{L}_{action}(\theta, I) &= \mathbb{E}_t \|\theta(x_t, t|I) - (x_1 - x_0)\|^2, \\ x_0 &\sim p_0, x_1 \sim p_1, x_t = tx_1 + (1-t)x_0. \end{aligned} \quad (7)$$

where p_0 represents a simple base density at time $t = 0$, p_1 represents the target complicated distribution at time $t = 1$, while x_0 and x_1 are the corresponding samplings. x_t is defined as the linear interpolation between x_0 and x_1 , following the Optimal Transport theory [28]. θ denotes the Flow Matching action generator of each arm, and I denotes the interactive condition.

Overall Framework. Bringing all together, our Decoupled Interaction Framework is developed. As shown in Fig. 3, we first decouple the joint action space used in the integrated control framework by assigning an independent model to each arm, thereby promoting the learning of uncoordinated tasks. Then, we propose a brand new selective interaction module, which takes the inputs of the current arm to predict the modulation factors and modulates the exchanged state representations from the other arm, thereby promoting the learning of coordinated tasks. Finally, we combine the original observation and state representations with the modulated state representations as interactive conditions and utilize the action generator to predict actions for each arm. Benefiting from our decoupled interaction design, our framework effectively accommodates the unique characteristics of different tasks in bimanual manipulation, achieving significant improvements across various tasks.

Implementation Details. In this paper, we adopt the same point cloud backbone as DP3 [37]. Following the setup in RoboTwin [27], we utilize joint angles for proprioception and predict joint angles as actions. In the loss function, the weights for the left arm and right arm are both set to 1. Our framework is implemented using PyTorch and trained on a single NVIDIA RTX 4090 GPU for 3000 epochs with the AdamW optimizer and a batch size of 120. Additional implementation details are provided in Appendix E.

4. Experiments

4.1. Experiment Setups

Simulation Benchmark. In this paper, we evaluate our framework on seven distinct manipulation tasks from the RoboTwin benchmark [27]. Based on the taxonomy of bimanual manipulation [15], we define two tasks as coordinated tasks and five tasks as uncoordinated tasks. Each task is designed to assess specific aspects and detailed task descriptions can be found in Appendix F.

Coordinated Tasks. (1) *Block Handover*: Transferring the long block from the left arm to the right arm and placing it at the designated location. (2) *Blocks Stack*: Stacking blocks of different colors in a specific order (the red block first and then the black block) to the designated location.

Uncoordinated Tasks. (1) *Dual Bottles Pick (Easy)*: Lifting the bottles that are positioned randomly and standing upright simultaneously. (2) *Dual Bottles Pick (Hard)*: Lifting the bottles that are positioned with random 6D poses simultaneously. (3) *Diverse Bottles Pick*: Lifting the bottles that vary in random shapes and do not repeat in the training and testing sets simultaneously. (4) *Block Hammer Beat*: Autonomously determining which arm picks up the hammer and strikes the block based on the position of the block. (5) *Empty Cup Place*: Autonomously determining which arm picks up the cup and places it at the designated location based on the position of the cup.

Baselines. We select several representative approaches as our baselines. For regression-based integrated control approaches, we choose ACT [40] as a baseline, which takes multi-view images as inputs and directly regresses actions of both arms with an action chunking mechanism. For generation-based integrated control approaches, we choose

Decoupled	Scaling Factors	Bias Vectors	Average
			0.426
✓			0.624 (↑ 0.198)
✓	✓		0.754 (↑ 0.328)
✓		✓	0.747 (↑ 0.321)
✓	✓	✓	0.789 (↑ 0.363)

Table 2. **The ablation study.** Important metrics are in gray cells. The green arrows indicate the performance difference between each line and the first line (the baseline).

Task Name	Demo-100	Demo-150	Demo-200
Blocks Stack	0.470	0.530	0.650
Dual Bottles Pick (Hard)	0.660	0.740	0.850

Table 3. **The scaling experiment on demonstration quantity and performance.** “Demo-X” indicates that X expert demonstrations are used during model training.

and modify DP [3], which takes multi-view images as inputs and directly generates actions of both arms with receding horizon control [26], and DP3 [37], which takes 3D point clouds as inputs and directly generates actions like DP, as baselines. For decoupled approaches, we refer to the architecture of Voxact-b [20] and construct a decoupled baseline that takes point clouds as inputs and outputs the actions of each arm separately. More detailed implementations of all methods can be seen in Appendix E.

4.2. Evaluation on RoboTwin Dataset

Quantitative Experiments. We test our method and other representative methods on the RoboTwin [27] dataset and report the results in Tab. 1. The models for different tasks are trained separately. From the table, we observe that: (1) Our method significantly outperforms previous methods. Compared to DP3 [37], the previous SOTA method, our method achieves an improvement of **23.5%** in average success rate. (2) Specifically, our method surpasses DP3 by **23.5%** in coordinated tasks and **23.5%** in uncoordinated tasks respectively, which verifies the effectiveness of the decoupled interaction design in our framework. (3) In the “Diverse Bottles Pick” task, the bottles used in the training and testing sets exhibit obvious differences in shape and texture. Achieving SOTA on this task verifies that our framework has great **intra-class generalization** abilities. (4) “Block Hammer Beat” and “Empty Cup Place” are single-arm manipulation tasks, where the challenge lies in the bimanual robot autonomously deciding which arm to use for manipulation based on the position of the object. Achieving a 90% success rate on these tasks verifies that our model possesses great **decision-making** abilities.

Qualitative Experiments. Moreover, we also visualize the execution process of DP3 [37] and our framework in the “Blocks Stack” task. Both models are evaluated within the same scene configuration for consistency. As shown in Fig.

Method	Coordinated	Uncoordinated	Average
PFM [4]	0.180	0.524	0.426
Ours (FM)	0.700	0.824	0.789 (↑ 0.363)
DP3 [37]	0.465	0.590	0.554
Ours (DDIM)	<u>0.645</u>	<u>0.686</u>	<u>0.674</u> (↑ 0.120)

Table 4. **Integrating into existing methods experiments.** Best results are highlighted in bold. The underlined values indicate the second-best results. Important comparison metrics are marked with gray cells. “FM” denotes Flow Matching [19].

Method	Model Size	Coordinated	Uncoordinated	Average
DP3	262.43M	0.465	0.590	0.554
	576.43M	0.150	0.446	0.361
Ours	42.95M	0.640	0.750	0.719
	146.37M	<u>0.665</u>	<u>0.788</u>	<u>0.753</u>
	536.07M	0.700	0.824	0.789

Table 5. **Model size and performance experiments.** Best results are highlighted in bold. The underlined values indicate the second-best results. Important metrics are marked with gray cells.

4, DP3 is more prone to failed grasps, while our model achieves more precise grasping. Additionally, for complex long-horizon tasks, DP3 often suffers from prolonged freezing behavior, where the arm remains stationary, as well as execution order errors. In contrast, our model can complete tasks in a more orderly manner. This is attributed to our decoupled interaction design, where the decoupled design decreases the action dimensions the model needs to predict, reducing the difficulty of network learning, and the interaction design modulates the exchanged state features of the other arm, effectively meeting the requirements for perceiving the other arm at different stages of coordinated tasks.

4.3. Ablation Study

Furthermore, we conduct a series of ablation studies to verify the effectiveness of each design in our framework. All experiments are trained on seven tasks in the RoboTwin dataset [27]. In the ablation study, we choose Flow Matching [19] as our action generator.

As illustrated in the second line of Tab. 2, the decoupled design enhances performance by 19.8% in average success rate, demonstrating that the decoupled design is effective as it reduces the dimensions of actions and simplifies network learning. Moreover, as illustrated in the fifth line of Tab. 2, adding the interaction module, which includes scaling factors and bias vectors, to the decoupled design further improves performance by 16.5% in average success rate. This demonstrates that the selective interaction design is beneficial, as it adaptively adjusts the exchanged information at different temporal stages during task execution. It should be noted that, from the third and the fourth lines of Tab. 2, although using scaling factors or bias vectors alone can effectively improve performance, they cannot achieve SOTA results. This is because they are complementary, i.e., only

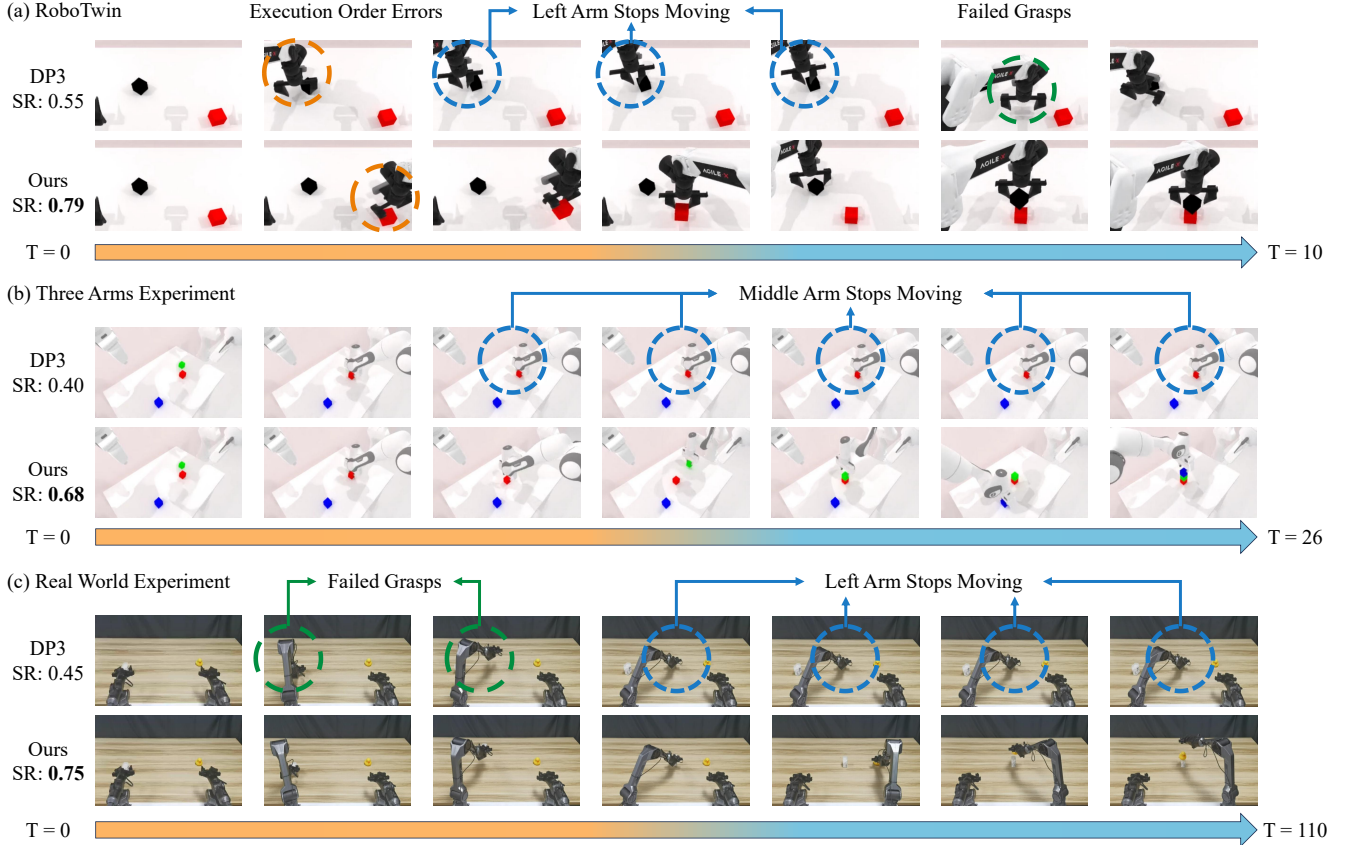


Figure 4. **Qualitative experiments.** In Fig. (a), we visualize the execution process of the “Blocks Stack” task for DP3 and our framework. In Fig. (b), we visualize the execution process of the three-arm experiment for DP3 and our framework. In Fig. (c), we visualize the execution process of DP3 and our framework in the real-world experiment. Dashed circles of different colors highlight common issues that typically arise in integrated control frameworks. “SR” denotes the success rate, which represents the average success rate of the model across all tasks under different experimental settings. Zoom in for a better view.

by scaling the intensity of the exchanged information and aligning it to the current state space together can the interaction information be fully modulated.

4.4. Model Analysis

Adhering to the Scaling Law. We analyze how the number of demonstrations affects the performance. We choose Flow Matching [19] as the action generator and test different numbers of demonstrations used in the model training. The results are recorded in Tab. 3. From the table, it can be observed that as the number of demonstrations increases, the performance of our model continues to improve, indicating that our framework adheres to the scaling law.

Integrating into Existing Methods Experiments. As mentioned in Sec. 1, our framework can be seamlessly integrated into existing methods. In this paper, we conduct experiments with two point cloud-based methods: DP3 [37] and Point Flow Matching (PFM) [4]. As illustrated in Tab. 4, applying the proposed decoupled interaction design to integrated control frameworks using DDIM and Flow Match-

ing significantly improves the success rate.

It should be noted that in DP3 [37], the performance of Flow Matching [19] implemented with DPM Solver++ [22] decreases compared to DDIM [30]. The explanation provided in the DP3 paper is that the high dimensions of the tasks pose a challenge to the learning of DPM Solver++. For integrated control frameworks, as seen in the first and the third lines of Tab. 4, the high dimensions of the tasks result in worse performance for PFM compared to DP3. From the second and fourth lines of Tab. 4, it can be seen that our decoupled design effectively reduces the action space that Flow Matching needs to explore, thereby fully leveraging the generative abilities of Flow Matching.

Model Size and Performance Experiments. We analyze the relationship between model size and performance, and record the results in Tab. 5. Here, the action generator we adopt is Flow Matching [19]. It can be concluded that: (1) From the first and second lines of Tab. 5, it can be observed that when we increase the size of the action generator in DP3 [37] to match the size of the action generator

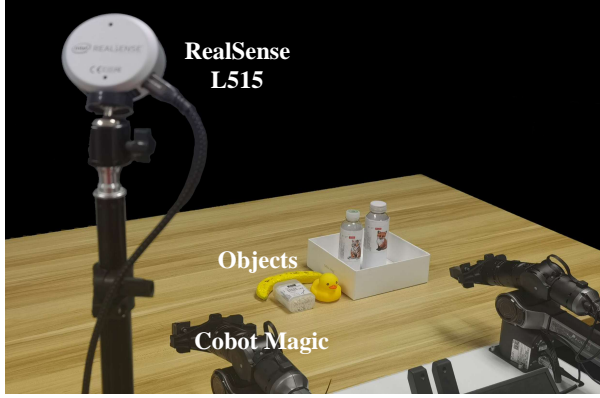


Figure 5. **Illustration of our real-world manipulation experimental settings.** We use the Cobot Magic as our bimanual robots and include everyday objects in our manipulation tasks. A RealSense L515 camera is applied to capture 3D point clouds.

in our framework, the performance of DP3 significantly decreases by **19.3%** compared to its original model size. (2) To test the performance of our framework with a smaller action generator, we compress the size of the action generator in our framework to match that of the action generator in DP3. As shown in the fourth and fifth lines of Tab. 5, our model experiences only a 3.6% performance drop despite a **73%** reduction in the number of parameters. (3) As shown in the third line of Tab. 5, when we further compress our model size to **1/6** of DP3 [37], our framework still outperforms DP3 by **16.5%** in success rate. These experiments demonstrate that the primary reason for the improvement of our framework is not the increase in model parameters but rather the decoupled interaction design itself.

Extensibility for Multi-Agent Manipulation. Bimanual manipulation is the simplest form of multi-agent manipulation. To demonstrate the extensibility of our method, i.e., its applicability to broader tasks, we build a three-arm block stacking task based on the Sapien [34] simulation. In this task, three robotic arms need to sequentially stack blocks of different colors in a specific order. In this experiment, a total of 50 demonstrations are utilized to train the two models independently. For evaluation, the success rate is measured for both models across 100 scenes.

As shown in Fig. 4, our framework achieves a **28%** improvement compared to DP3. At the same time, the visualized trajectories indicate that as the dimensionality of the action space in integrated control frameworks further increases, the prolonged stationary state of the robotic arm after completing an action becomes more serious.

4.5. Real-World Experiments.

Settings. To verify the practical capabilities of our Decoupled Interaction Framework, we conduct real-world manipulation experiments across diverse tasks. For evaluation, we

Method	Coordinated		Uncoordinated		Total
	Banana Handover	Items Stack	Dual Bottles Pick	Banana Place	
DP3	6/15	5/15	6/15	10/15	27/60
Ours	10/15	10/15	13/15	12/15	45/60

Table 6. **Real world experiments.** Best is in bold face.

compare our framework with the integrated control baseline DP3 [37] using 50 high-quality demonstrations collected via teleoperation similar to [40].

Coordinated Tasks. (1) *Banana Handover*: Picking the banana randomly placed on the left side of the table with the left arm, transferring it to the right arm, and placing it in the box. (2) *Items Stack*: Placing the box first and then the toy on top of the box in the designated location in order.

Uncoordinated Tasks. (1) *Dual Bottles Pick*: Lifting the bottles that are positioned randomly and standing upright simultaneously. (2) *Banana Pick and Place*: Autonomously determining which arm picks up the banana that is positioned randomly on the table and places it into the box.

Quantitative Experiments. The experiment setup and objects are shown in Fig. 5, and the testing results are shown in Tab. 6. It can be observed that our method achieves robust real-world manipulation ability, surpassing the SOTA method DP3 [37] in real-world experiments, showing great potential for our Decoupled Interaction Framework.

Qualitative Experiments. We visualize the execution process of DP3 [37] and our framework in “Items Stack” task. As shown in Fig. 4, DP3 often appears failed grasps, while our framework achieves more stable grasping. More examples and details can be seen in the video demo provided in the supplementary materials.

5. Conclusion

In this paper, we propose learning bimanual manipulation through precise task categorization and introduce a novel Decoupled Interaction Framework for bimanual robotic manipulation. The key insight of our framework is to decouple the action space the network needs to explore by assigning an independent model to each arm, thereby enhancing the learning of uncoordinated tasks. Additionally, we propose a brand new selective interaction module that learns its weights from the inputs of its own arm to adaptively modulate the exchanged state features, thereby improving the learning of coordinated tasks. Extensive experiments on seven tasks in the RoboTwin dataset demonstrate that our framework exhibits **effectiveness, flexibility, extensibility, and scalability**. We believe our decoupled interaction framework provides insights for the community to further explore the modeling of bimanual and even multi-agent manipulation tasks. To advance the community, we will release our code upon the publication of the paper.

References

- [1] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil J. Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael S. Ryoo, Grecia Salazar, Pannag R. Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong T. Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. RT-1: robotics transformer for real-world control at scale. In *Robotics: Science and Systems*, 2023. 3
- [2] Tianxing Chen, Yao Mu, Zhixuan Liang, Zanzin Chen, Shijia Peng, Qiangyu Chen, Mingkun Xu, Ruizhen Hu, Hongyuan Zhang, Xuelong Li, and Ping Luo. G3flow: Generative 3d semantic flow for pose-aware and generalizable object manipulation. *arXiv preprint arXiv:2411.18369*, 2024. 1, 3
- [3] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Robotics: Science and Systems*, 2023. 3, 5, 6, 12
- [4] Eugenio Chisari, Nick Heppert, Max Argus, Tim Welschhold, Thomas Brox, and Abhinav Valada. Learning robotic manipulation policies from point clouds with conditional flow matching. In *Conference on Robot Learning*, 2024. 2, 3, 5, 6, 7
- [5] Haoran Geng, Ziming Li, Yiran Geng, Jiayi Chen, Hao Dong, and He Wang. Partmanip: Learning cross-category generalizable part manipulation policy from point cloud observations. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 3
- [6] Yiran Geng, Boshi An, Haoran Geng, Yuanpei Chen, Yaodong Yang, and Hao Dong. Rlafford: End-to-end affordance learning for robotic manipulation. In *IEEE International Conference on Robotics and Automation*, 2023. 3
- [7] Dibya Ghosh, Homer Rich Walke, Karl Pertsch, Kevin Black, Oier Mees, Sudeep Dasari, Joey Hejna, Tobias Kreiman, Charles Xu, Jianlan Luo, You Liang Tan, Lawrence Yunliang Chen, Quan Vuong, Ted Xiao, Pannag R. Sanketi, Dorsa Sadigh, Chelsea Finn, and Sergey Levine. Octo: An open-source generalist robot policy. In *Robotics: Science and Systems*, 2024. 3
- [8] Markus Grotz, Mohit Shridhar, Tamim Asfour, and Dieter Fox. Peract2: Benchmarking and learning for robotic bimanual manipulation tasks. In *Conference on Robot Learning*, 2024. 1, 3
- [9] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. In *Advances in Neural Information Processing Systems*, 2016. 1, 3
- [10] Tianying Ji, Yongyuan Liang, Yan Zeng, Yu Luo, Guowei Xu, Jiawei Guo, Ruijie Zheng, Furong Huang, Fuchun Sun, and Huazhe Xu. ACE: off-policy actor-critic with causality-aware entropy regularization. In *International Conference on Machine Learning*, 2024. 3
- [11] Daqi Jiang, Hong Wang, and Yanzheng Lu. Mastering the complex assembly task with a dual-arm robot: A novel reinforcement learning method. *IEEE Robotics and Automation Magazine*, 2023. 1
- [12] Tsung-Wei Ke, Nikolaos Gkanatsios, and Katerina Fragkiadaki. 3d diffuser actor: Policy diffusion with 3d scene representations. *arXiv preprint arXiv:2402.10885*, 2024. 3
- [13] 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. In *Conference on Robot Learning*, 2024. 1
- [14] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailov, Ethan Paul Foster, Grace Lam, Pannag Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model. *arXiv preprint arXiv:2406.09246*, 2024. 3
- [15] Franziska Krebs and Tamim Asfour. A bimanual manipulation taxonomy. *IEEE Robotics and Automation Letters*, 2022. 1, 2, 3, 5
- [16] Andrew Lee, Ian T. Chuang, Ling-Yuan Chen, and Iman Soltani. Interact: Inter-dependency aware action chunking with hierarchical attention transformers for bimanual manipulation. In *Conference on Robot Learning*, 2024. 1, 3
- [17] Xiaoli Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 3
- [18] Zhixuan Liang, Yao Mu, Hengbo Ma, Masayoshi Tomizuka, Mingyu Ding, and Ping Luo. Skilldiffuser: Interpretable hierarchical planning via skill abstractions in diffusion-based task execution. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 3
- [19] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *International Conference on Learning Representations*, 2023. 5, 6, 7, 11, 12
- [20] I-Chun Arthur Liu, Sicheng He, Daniel Seita, and Gaurav S. Sukhatme. Voxact-b: Voxel-based acting and stabilizing policy for bimanual manipulation. In *Conference on Robot Learning*, 2024. 3, 5, 6, 12
- [21] Songming Liu, Lingxuan Wu, Bangguo Li, Hengkai Tan, Huayu Chen, Zhengyi Wang, Ke Xu, Hang Su, and Jun Zhu. RDT-1B: a diffusion foundation model for bimanual manipulation. *arXiv preprint arXiv:2410.07864*, 2024. 1, 3
- [22] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 7
- [23] Guanxing Lu, Shiyi Zhang, Ziwei Wang, Changliu Liu, Jiwen Lu, and Yansong Tang. Manigaussian: Dynamic gaus-

- sian splatting for multi-task robotic manipulation. In *European Conference on Computer Vision*, 2024. 3
- [24] Jianlan Luo, Zheyuan Hu, Charles Xu, You Liang Tan, Jacob Berg, Archit Sharma, Stefan Schaal, Chelsea Finn, Abhishek Gupta, and Sergey Levine. SERL: A software suite for sample-efficient robotic reinforcement learning. In *IEEE International Conference on Robotics and Automation*, 2024. 3
- [25] Xiao Ma, Sumit Patidar, Iain Haughton, and Stephen James. Hierarchical diffusion policy for kinematics-aware multi-task robotic manipulation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024. 3
- [26] Hannah Michalska and David Q. Mayne. Robust receding horizon control of constrained nonlinear systems. *IEEE Transactions on Automatic Control*, 1993. 6
- [27] Yao Mu, Tianxing Chen, Shijia Peng, Zhanxin Chen, Zeyu Gao, Yude Zou, Lunkai Lin, Zhiqiang Xie, and Ping Luo. Robotwin: Dual-arm robot benchmark with generative digital twins (early version). *arXiv preprint arXiv:2409.02920*, 2024. 1, 2, 5, 6, 11
- [28] Gabriel Peyré and Marco Cuturi. Computational optimal transport. *Foundations and Trends in Machine Learning*, 2019. 5
- [29] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Conference on Robot Learning*, 2022. 3
- [30] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 5, 7, 11, 12
- [31] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In *International Joint Conference on Artificial Intelligence*, 2018. 1, 3
- [32] Chuan Wen, Xingyu Lin, John Ian Reyes So, Kai Chen, Qi Dou, Yang Gao, and Pieter Abbeel. Any-point trajectory modeling for policy learning. In *Robotics: Science and Systems*, 2024. 3
- [33] Tianhao Wu, Mingdong Wu, Jiyao Zhang, Yunchong Gan, and Hao Dong. Learning score-based grasping primitive for human-assisting dexterous grasping. In *Advances in Neural Information Processing Systems*, 2023. 3
- [34] Fanbo Xiang, Yuzhe Qin, Kaichun Mo, Yikuan Xia, Hao Zhu, Fangchen Liu, Minghua Liu, Hanxiao Jiang, Yifu Yuan, He Wang, Li Yi, Angel X. Chang, Leonidas J. Guibas, and Hao Su. SAPIEN: A simulated part-based interactive environment. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 8
- [35] Guowei Xu, Ruijie Zheng, Yongyuan Liang, Xiyao Wang, Zhecheng Yuan, Tianying Ji, Yu Luo, Xiaoyu Liu, Jiaxin Yuan, Pu Hua, Shuzhen Li, Yanjie Ze, Hal Daumé III, Furong Huang, and Huazhe Xu. Drm: Mastering visual reinforcement learning through dormant ratio minimization. In *International Conference on Learning Representations*, 2024. 3
- [36] Yanjie Ze, Ge Yan, Yueh-Hua Wu, Annabella Macaluso, Yuying Ge, Jianglong Ye, Nicklas Hansen, Li Erran Li, and Xiaolong Wang. Gnfactor: Multi-task real robot learning with generalizable neural feature fields. In *Conference on Robot Learning*, 2023. 3
- [37] Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, Muhan Wang, and Huazhe Xu. 3d diffusion policy: Generalizable visuomotor policy learning via simple 3d representations. In *Robotics: Science and Systems*, 2024. 2, 3, 5, 6, 7, 8, 12
- [38] Tianle Zhang, Dongjiang Li, Yihang Li, Zecui Zeng, Lin Zhao, Lei Sun, Yue Chen, Xuelong Wei, Yibing Zhan, Lusong Li, and Xiaodong He. Empowering embodied manipulation: A bimanual-mobile robot manipulation dataset for household tasks. *arXiv preprint arXiv:2405.18860*, 2024. 1
- [39] Tony Z. Zhao, Jianlan Luo, Oleg Sushkov, Rugile Pevcevičiute, Nicolas Heess, Jon Scholz, Stefan Schaal, and Sergey Levine. Offline meta-reinforcement learning for industrial insertion. In *International Conference on Robotics and Automation*, 2022. 3
- [40] Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual manipulation with low-cost hardware. In *Robotics: Science and Systems*, 2023. 1, 3, 5, 8
- [41] Jihong Zhu, Michael Gienger, Giovanni Franzese, and Jens Kober. Do you need a hand? - A bimanual robotic dressing assistance scheme. *IEEE Transactions on Robotics*, 2024. 1
- [42] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, Quan Vuong, Vincent Vanhoucke, Huong T. Tran, Radu Soricut, Anikait Singh, Jaspiar Singh, Pierre Sermanet, Pannag R. Sanketi, Grecia Salazar, Michael S. Ryoo, Krista Reymann, Kanishka Rao, Karl Pertsch, Igor Mordatch, Henryk Michalewski, Yao Lu, Sergey Levine, Lisa Lee, Tsang-Wei Edward Lee, Isabel Leal, Yuheng Kuang, Dmitry Kalashnikov, Ryan Julian, Nikhil J. Joshi, Alex Irpan, Brian Ichter, Jasmine Hsu, Alexander Herzog, Karol Hausman, Keerthana Gopalakrishnan, Chuyuan Fu, Pete Florence, Chelsea Finn, Kumar Avinava Dubey, Danny Driess, Tianli Ding, Krzysztof Marcin Choromanski, Xi Chen, Yevgen Chebotar, Justice Carbajal, Noah Brown, Anthony Brohan, Montserrat Gonzalez Arenas, and Kehang Han. RT-2: vision-language-action models transfer web knowledge to robotic control. In *Conference on Robot Learning*, 2023. 3

Appendix

A. Video Demo

A video demo is provided for both simulation and real-world manipulation experiments using our Decoupled Interaction Framework, show the **effectiveness**, **extensibility**, and **practical applicability** of our framework. In simulation experiments, we first collect 50 demonstrations for various tasks in the Sapien simulation environment. We then train our framework separately for each task. For real-world experiments, we utilize the Cobot Magic robotic arm and the RealSense L515 camera, employing single-view, third-person point clouds for model training and inference. Watch the video for more details. Enjoy and have fun !

B. Analysis of Selective Interaction Module

We further explore various interaction design approaches. Specifically, we evaluate concatenation and MLP, which are widely adopted in the Computer Vision community for interaction modeling. The results are presented in Tab. 7. As shown in Tab. 7, compared to the MLP and concatenation-based interaction modeling, our model achieves a performance improvement of at least 5.3%, demonstrating the effectiveness of our selective interaction module.

C. Visualization of the Generated Manipulation Trajectories in RoboTwin

In this section, we visualize some manipulation trajectories in RoboTwin [27] generated by our Decoupled Interaction Framework. As illustrated in Fig. 6, it can be concluded that: (1) We visualize two manipulation trajectories with different target objects in the “Diverse Bottles Pick” task to demonstrate the capability of our framework for **intra-class generalization**. (2) Additionally, we visualize trajectories for the “Block Hammer Beat” and “Empty Cup Place” tasks with varying initial positions of the target objects. As shown in Fig. 6, our model autonomously decides which arm to use based on the position of the target object, demonstrating its **decision-making** ability.

D. Visualization of the Generated Manipulation Trajectories in Real World

In this section, we visualize some manipulation trajectories in real world generated by our Decoupled Interaction Framework. As illustrated in 7, our framework effectively handles both coordinated and uncoordinated tasks. This is because our decoupled design effectively reduces the high-dimensional action space, enabling the network to learn actions more efficiently. Additionally, our selective interaction module explicitly models the interaction between the

arms, allowing our framework to better meet the cooperation requirements of various bimanual manipulation tasks.

E. Implementation Details

In this section, we provide a detailed introduction to the implementation details of all baselines and our Decoupled Interaction Framework.

Training Setup. The key training setup for our Decoupled Interaction Framework based on the DDIM [30] and Flow Matching [19] is detailed in Tab. 8.

Baseline Setups. We also outline the training settings for the baseline in Tab. 9. Because of the differences in hyper-parameters between ACT and other baselines, we provide a description of its hyper-parameters here. For the ACT method, we use the AdamW optimizer with an initial learning rate of $1.0e-5$ and a weight decay of $1.0e-4$. The training process employs a batch size of 8, runs for 2000 epochs and uses an action chunking size of 100.

F. Simulation Tasks

We also visualize the distinct manipulation tasks from the RoboTwin benchmark [27], as illustrated in Fig. 8, and provide detailed descriptions of all simulation tasks in Tab. 10, totaling seven tasks.

	Coordinated		Uncoordinated					Average
	Block Handover	Blocks Stack Easy	Dual Bottles Pick Easy	Dual Bottles Pick Hard	Diverse Bottles Pick	Empty Cup Place	Block Hammer Beat	
MLP	1.000	0.330	0.960	0.550	0.590	0.860	0.860	0.736 (\downarrow 0.053)
Concat	1.000	0.370	0.960	0.510	0.590	0.810	0.890	0.733 (\downarrow 0.056)
Ours	1.000	0.400	0.990	0.630	0.700	0.900	0.900	0.789

Table 7. **Illustration of the performance with different interaction designs on seven tasks in the RoboTwin dataset.** Best results are highlighted in bold. Important comparison metrics are marked with gray cells. The red arrows indicate the performance difference between each baseline and our method.

Parameter	Ours (DDIM[30])	Ours (Flow Matching[19])
horizon	8	8
n_obs_steps	3	3
n_action_steps	6	6
num_inference_steps	10	10
dataloader.batch_size	120	120
dataloader.num_workers	8	8
dataloader.shuffle	True	True
dataloader.pin_memory	True	True
dataloader.persistent_workers	False	False
optimizer._target_	torch.optim.AdamW	torch.optim.AdamW
optimizer.lr	1.0e-4	3.0e-5
optimizer.betas	[0.95, 0.999]	[0.95, 0.999]
optimizer.eps	1.0e-8	1.0e-8
optimizer.weight_decay	1.0e-6	1.0e-6
training.lr_scheduler	cosine	cosine
training.lr_warmup_steps	500	10
training.num_epochs	3000	3000
training.gradient_accumulate_every	1	1
training.use_ema	True	True

Table 8. **Model training settings.** Hyper-parameter Settings for Training and Deployment of our Decoupled Interaction Framework.

Parameter	DP [3]	DP3 [37]	Voxact-b [20]
horizon	8	8	8
n_obs_steps	3	3	3
n_action_steps	6	6	6
num_inference_steps	100	10	10
dataloader.batch_size	128	256	256
dataloader.num_workers	0	8	8
dataloader.shuffle	True	True	True
dataloader.pin_memory	True	True	True
dataloader.persistent_workers	False	False	False
optimizer._target_	torch.optim.AdamW	torch.optim.AdamW	torch.optim.AdamW
optimizer.lr	1.0e-4	1.0e-4	1.0e-4
optimizer.betas	[0.95, 0.999]	[0.95, 0.999]	[0.95, 0.999]
optimizer.eps	1.0e-8	1.0e-8	1.0e-8
optimizer.weight_decay	1.0e-6	1.0e-6	1.0e-6
training.lr_scheduler	cosine	cosine	cosine
training.lr_warmup_steps	500	500	500
training.num_epochs	300	3000	3000
training.gradient_accumulate_every	1	1	1
training.use_ema	True	True	True

Table 9. **Baselines settings.** Hyper-parameter Settings for Training and Deployment of DP, DP3 and Voxact-b Algorithms.

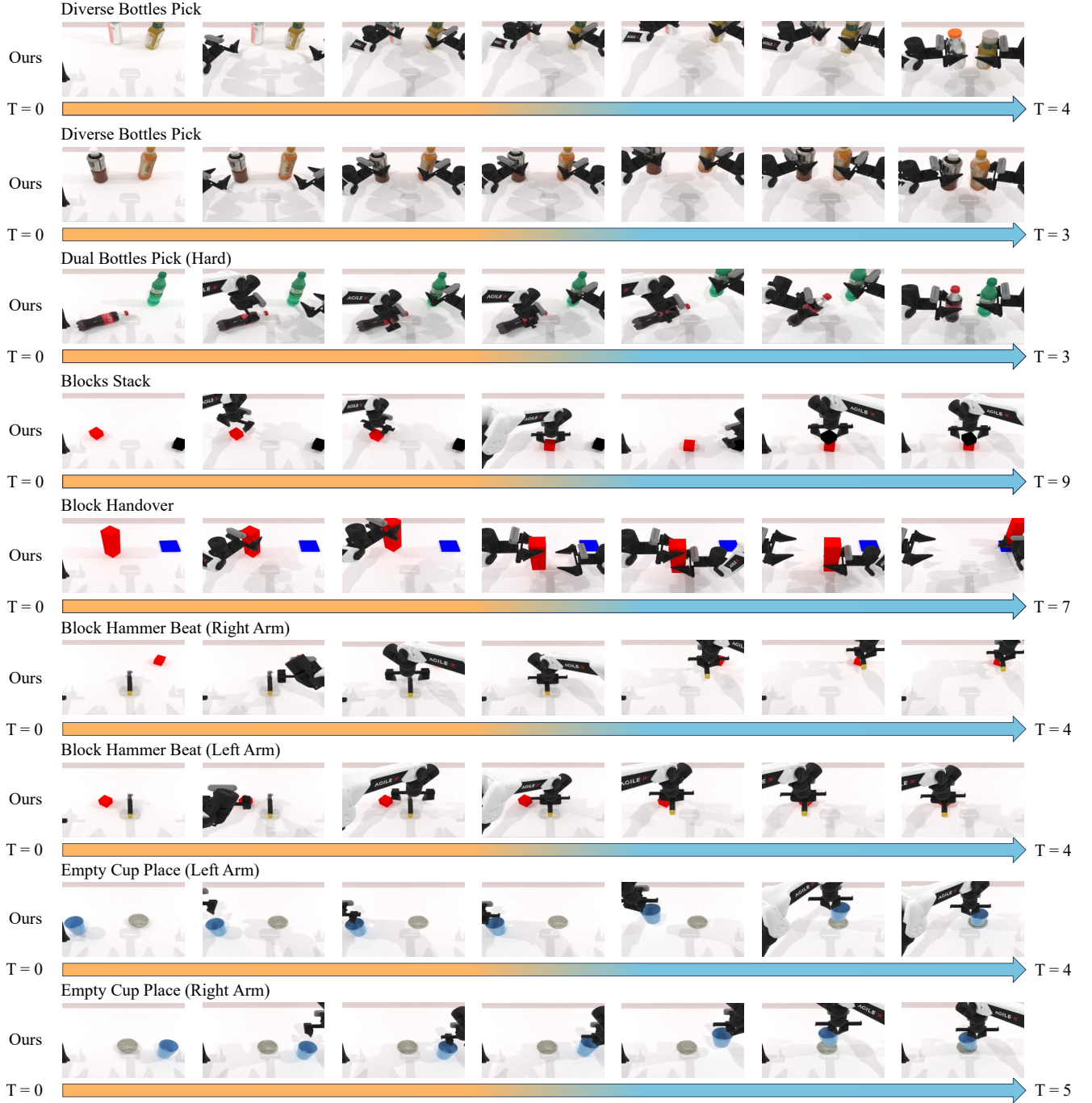


Figure 6. **Visualization of the generated manipulation trajectories of our framework in RoboTwin.** We visualize different coordinated and uncoordinated tasks within various scenes. Zoom in for the best view.

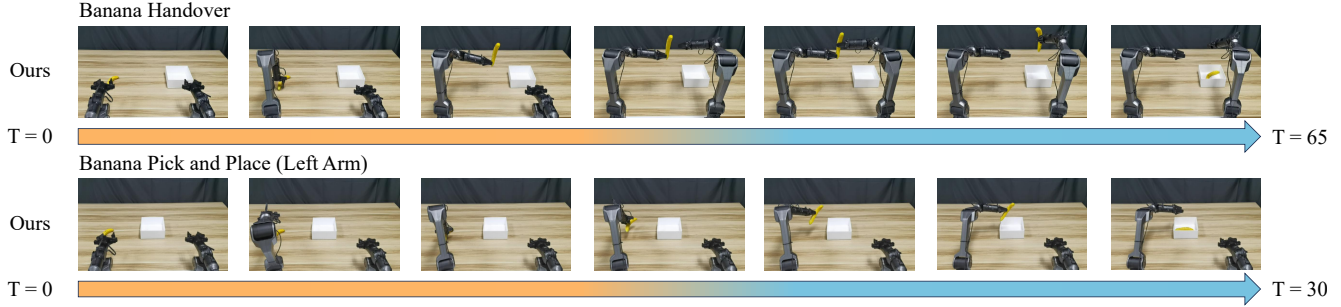


Figure 7. **Visualization of the generated manipulation trajectories of our framework in real-world experiments.** We visualize coordinated and uncoordinated tasks under different configurations. Zoom in for the best view.

<i>Task</i>	<i>Description</i>
<i>Block Handover</i>	A long block is placed on the left side of the table. The left arm grasps the upper side of the block and then hands it over to the right arm, which places the block on the blue mat on the right side of the table.
<i>Blocks Stack Easy</i>	Red and black cubes are placed randomly on the table. The robotic arm stacks the cubes in order, placing the red cubes first, followed by the black cubes, in the designated target location.
<i>Dual Bottles Pick Easy</i>	A red bottle is placed randomly on the left side, and a green bottle is placed randomly on the right side of the table. Both bottles are standing upright. The left and right arms are used simultaneously to lift the two bottles to a designated location.
<i>Dual Bottles Pick Hard</i>	A red bottle is placed randomly on the left side, and a green bottle is placed randomly on the right side of the table. The bottles' postures are random. Both left and right arms are used simultaneously to lift the two bottles to a designated location.
<i>Diverse Bottles Pick</i>	A random bottle is placed on the left and right sides of the table. The bottles' designs are random and do not repeat in the training and testing sets. Both left and right arms are used to lift the two bottles to a designated location.
<i>Empty Cup Place</i>	An empty cup and a cup mat are placed randomly on the left or right side of the table. The robotic arm places the empty cup on the cup mat.
<i>Block Hammer Beat</i>	There is a hammer and a block in the middle of the table. If the block is closer to the left robotic arm, it uses the left arm to pick up the hammer and strike the block; otherwise, it does the opposite.

Table 10. **Task descriptions for RoboTwin platform.**

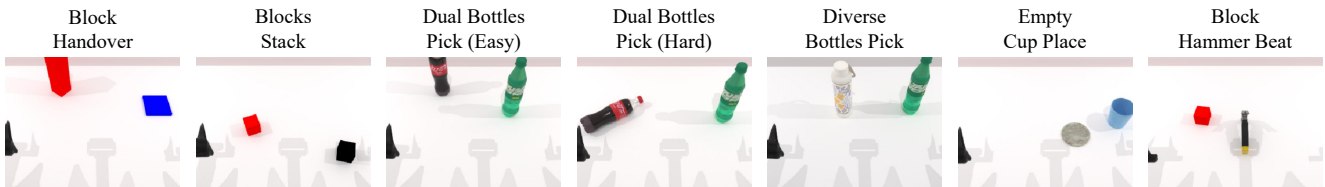


Figure 8. **Seven testing benchmark tasks.** We visualize manipulation tasks used in the RoboTwin benchmark. Zoom in for the best view.