

Spatial-Temporal Graph Diffusion Policy with Kinematic Modeling for Bimanual Robotic Manipulation

Qi Lv^{1,2,3} Hao Li^{1,3} Xiang Deng^{1†} Rui Shao¹ Yinchuan Li³ Jianye Hao³
 Longxiang Gao⁴ Michael Yu Wang² Liqiang Nie^{1†}

¹Harbin Institute of Technology (Shenzhen) ²Great Bay University

³Huawei Noah's Ark Lab ⁴Shandong Computer Science Center

lvqi@stu.hit.edu.cn, dengxiang@hit.edu.cn

Abstract

Despite the significant success of imitation learning in robotic manipulation, its application to bimanual tasks remains highly challenging. Existing approaches mainly learn a policy to predict a distant next-best end-effector pose (NBP) and then compute the corresponding joint rotation angles for motion using inverse kinematics. However, they suffer from two important issues: (1) **rarely considering the physical robotic structure**, which may cause self-collisions or interferences, and (2) **overlooking the kinematics constraint**, which may result in the predicted poses not conforming to the actual limitations of the robot joints. In this paper, we propose *Kinematics enhanced Spatial-Temporal Graph Diffuser (KStar Diffuser)*. Specifically, (1) to incorporate the physical robot structure information into action prediction, KStar Diffuser maintains a dynamic spatial-temporal graph according to the physical bimanual joint motions at continuous timesteps. This dynamic graph serves as the robot-structure condition for denoising the actions; (2) to make the NBP learning objective consistent with kinematics, we introduce the differentiable kinematics to provide the reference for optimizing KStar Diffuser. This module regularizes the policy to predict more reliable and kinematics-aware next end-effector poses. Experimental results show that our method effectively leverages the physical structural information and generates kinematics-aware actions in both simulation and real-world.

1. Introduction

Bimanual manipulation [38, 40, 47, 50] represents a fundamental capability for robotic systems to perform com-

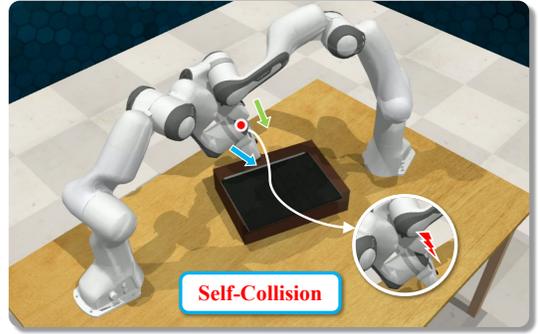


Figure 1. The self-collision problem in bimanual manipulation tasks due to overlooking the robotic structure.

plex tasks requiring two-arm coordination. While imitation learning has demonstrated remarkable success in single-arm manipulation [6, 7, 34, 37], its extension to bimanual scenarios faces unique challenges as robots need to coordinate dual-arm movements while conforming to physical constraints. These challenges significantly impact the reliability and feasibility of predicted actions in real-world applications.

A straightforward approach [9, 10, 13, 60, 61, 65] would be to learn the entire motion trajectory by predicting joint positions at consecutive timesteps. However, this poses significant challenges as trajectories are typically long and difficult to learn. Therefore, mainstream approaches [15, 17, 35, 49] typically adopt a two-stage pipeline: first predicting the next-best end-effector pose (NBP), then computing joint rotations through inverse kinematics [21, 26]. However, the dominant NBP-based approaches, though simplify the learning process, often lead to unreliable motion generation. Concretely, the predicted poses may violate physical structure constraints [27, 48], causing potential arm collisions, or exceeding joint limits due to inadequate consideration of kinematic feasibility, as Figure 1 illustrated.

In real-world bimanual manipulation tasks, existing ap-

*This work was done when Qi Lv and Hao Li were intern at Huawei Noah's Ark Lab. †Corresponding author

proaches [9, 15, 17, 49, 64] exhibit significant limitations in motion execution despite their promising performance in trajectory prediction. Empirical observations indicate that while individual end-effector poses are feasible in isolation, their concurrent execution frequently results in inter-arm collisions [48] and kinematically infeasible configurations [35, 39]. The fundamental challenge lies in the optimization of end-effector poses exclusively in Cartesian space, which introduces a substantial discrepancy between motion prediction and physical constraints. Existing approaches [9, 10, 17] incorporate proprioception information in a simplistic manner, (*e.g.*, gripper rotation angles [17, 49] and end-effector poses [53]), but disregarding the robot’s intrinsic structural properties *e.g.*, the kinematic chain and joint configurations during the prediction phase. This oversimplified representation fails to capture the crucial spatial correlations between robotic arms, joint links, and their potential interactions. Consequently, the predicted dual-arm trajectories may satisfy task objectives in Cartesian space while violating crucial spatial constraints during execution.

Moreover, the conventional paradigm [17, 24, 49, 51] of treating inverse kinematics as a post-processing procedure introduces additional complexities. The disjunction between kinematic constraints and the pose prediction learning objective often yields solutions that exhibit apparent smoothness in Cartesian space while manifesting discontinuous or mechanically infeasible trajectories in the joint space. This limitation becomes particularly pronounced in configurations approaching kinematic singularities or joint limits, where the mapping between Cartesian and joint spaces becomes ill-conditioned. These inherent limitations necessitate the integration of both structural and kinematic constraints into the policy learning framework.

To address the above issues, we propose **Kinematics enhanced Spatial-Temporal gRaph Diffuser (KStar Diffuser)**, a novel framework that explicitly incorporates both robot structures and kinematics into the bimanual motion generation process. Our key insight is that the physical structure and kinematic properties of the robot should guide the learning process of pose prediction, rather than be treated as independent post-processing constraints. Specifically, **(1) to incorporate structural awareness**, we construct a dynamic spatial-temporal graph from the robot’s URDF specifications, where nodes represent joint properties and edges capture both spatial relations and temporal dependencies. This graph structure is encoded via Graph Convolutional Network (GCN) [25] to provide explicit physical constraints for the diffusion process; **(2) for kinematic feasibility**, we regularize the NBP learning objective by incorporating joint-space prediction, where the predicted joint positions are mapped to reference end-effector poses through differentiable forward kinematics. These kinematically-feasible poses then serve as conditions

to guide the diffusion process, ensuring that the generated motions satisfy both structural and kinematic constraints. Our main contributions are as follows:

- Different from existing approaches that optimize end-effector poses solely in Cartesian space, we propose a novel spatial-temporal robot graph that explicitly models the robot physical configuration to guide the generative action denoising procedure.
- We introduce a kinematics regularizer that augments the NBP learning objective by introducing joint-space supervision. This regularizer leverages forward kinematics to provide kinematically-feasible reference poses, effectively guiding the diffusion process to conform to kinematic constraints.
- Extensive experiments show that our proposed KStar Diffuser is superior in both simulation and real-world scenarios, surpassing baselines more than 10% in success rate.

2. Related Work

Diffusion Models in Bimanual Robotic Controls. Diffusion models such as Denoising Diffusion Probabilistic Models (DDPM) [19] have achieved great success in the fields of image generation [43, 44, 46] and video generation [41, 55, 66]. Thus, recent work [8, 45, 52, 57, 67] has been devoted to applying its powerful generation capabilities to action generation for robotic manipulation tasks in a imitation learning mode. With the 2D images [9, 30–32, 35, 45, 52, 54, 62, 63] or 3D point cloud observations [8, 20, 23, 60, 61], the policy is trained to output the action sequence of joint positions or end-effector poses by iterative denoising process.

Bimanual manipulation tasks, which resemble human-like actions more closely, have garnered increasing attention [10, 14, 17, 33, 65]. However, existing research primarily extends single-arm manipulation methods to dual-arm tasks [10, 65], without considering about the distinct challenges specific to bimanual settings. The relative independence of the two arms makes it essential to consider self-collision avoidance [1, 3] when predicting actions. In practice, the robot’s structure significantly influences its motion since the joint types determine movement directions and joint angles constrain the range of motion. Although some studies [15, 49, 53] integrate proprioception in robotic control, they often use simply low-dimensional vectors to represent body information, lacking depth in structural exploration and overlooking critical spatial relationships between dual arms. Therefore, we study the physical structure information of the dual-arm robot system, and improve both the accuracy and adaptability of robotic movements.

Robotic Kinematics Modeling. Robotic kinematics modeling [2, 16, 35, 58, 59] is a classic problem in the robotic control [4, 18, 36, 56] where inverse kinematics (IK) presents a fundamental challenge, *i.e.*, deriving the

joint configuration given the end-effector pose. The inherent non-uniqueness of IK solutions as a significant obstacle for learning algorithms [12, 22], often produce inaccurate models by averaging over nonconvex feasible sets. Recent advancements integrate neural solutions [11, 28, 29] to enhance the adaptability and precision of kinematic models. Bócsi *et al.* [5] leveraged support vector machines to parameterize quadratic programs, aligning solutions with IK in specific workspace regions. Current data-driven imitation learning approaches [10, 13, 15, 17, 65] primarily rely on probabilistic modeling of action distributions to predict subsequent actions, but lack the reliability guarantees from the aspect of kinematics. In this paper, we construct a spatial-temporal graph to learn the robot physical representation, and provide a kinematics-awared latent end-effect embedding for diffusion policy as guidance. In this way, our approach enhances the robot’s ability to perform tasks with superious precision and adaptability.

3. Method

3.1. Preliminary

Diffusion Policy. Chi *et al.* [9] propose the diffusion policy which represents a robot’s visuomotor policy as a conditional denoising diffusion process. It learns a model distribution $\pi_\phi(\mathbf{a}_0|\mathbf{o})$ conditioned on observation \mathbf{o} to approximate joint distribution $q(\mathbf{a}_0)$. The whole procedure consists of a *forward process* and a *reverse process*.

(1) Forward Process: For the Markov chain with Gaussian transitions parameterized, the policy $\pi_\phi(\mathbf{a}_0|\mathbf{o})$ learns a fixed inference procedure $q(\mathbf{a}_k|\mathbf{a}_0)$:

$$\begin{aligned} q(\mathbf{a}_k|\mathbf{a}_0) &= \int q(\mathbf{a}_{1:k}|\mathbf{a}_0) d\mathbf{a}_{1:(k-1)} \\ &= \mathcal{N}(\mathbf{a}_k; \sqrt{\bar{\alpha}_k}\mathbf{a}_0, (1 - \bar{\alpha}_k)\mathbf{I}) \end{aligned} \quad (1)$$

Thus, \mathbf{a}_k can be expressed as a linear combination of \mathbf{a}_0 and a noise variable ϵ :

$$\mathbf{a}_k = \sqrt{\bar{\alpha}_k}\mathbf{a}_0 + \sqrt{1 - \bar{\alpha}_k}\epsilon \quad (2)$$

The training loss is obtained:

$$\mathcal{L} = \mathbb{E}_{\mathbf{a}_0 \sim q(\mathbf{a}_0), \epsilon_k \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), k} [\|\epsilon_k, \pi_\phi(\mathbf{a}_0 + \epsilon_k, \mathbf{o}, k)\|^2] \quad (3)$$

where ϵ_k is the random noise sampling at k iteration.

(2) Reverse Process: Starting from a sample $\mathbf{a}_K \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, the reverse steps are:

$$\mathbf{a}_{k-1} = \sqrt{\bar{\alpha}_{k-1}}\pi_\phi(\mathbf{a}_k, k, \mathbf{o}) + \sqrt{1 - \bar{\alpha}_{k-1}}\epsilon. \quad (4)$$

By the iterative denoising process, the policy generates \mathbf{a}_0 as its next action. The continuous predicted action sequence forms the complete action trajectory.

Robotic Kinematics. Robotics kinematics describes the relationship between the robot joints and its end-effector. It can be divided into *forward kinematics* and *inverse kinematics* for controlling robot movements.

(1) Forward Kinematics (FK): Given a specific configuration of joint angles, i.e., $\theta = [\theta_1, \theta_2, \dots, \theta_n] \in \Theta$, it aims to construct a function $f_{FK} : \mathbb{R}^n \rightarrow SE(3)$, computing the end-effector position and orientation that are represented as a pose matrix $\mathbf{T} \in SE(3)$.

$$FK(\theta) = \mathbf{T} \in \mathcal{T}, \text{ where } \mathbf{T} = \prod_{i=1}^n \mathbf{T}_i(\theta_i), \quad (5)$$

where $\mathbf{T}_i(\theta_i)$ is the homogeneous transformation matrix for joint i , and \mathcal{T} means the task space. $\Theta \in \mathbb{R}^n$ denotes the configuration space, while the configuration refers to the specific arrangement of the robotic joints in a given pose.

(2) Inverse Kinematics (IK): In contrast, IK aims to find the joint configurations which achieve a desired end-effector pose. It is defined as a mapping $f_{IK} : SE(3) \rightarrow \mathbb{R}^n$ from the task space \mathcal{T} to the set of possible configurations Θ that satisfy the given pose matrix \mathbf{T} :

$$IK(\mathcal{T}) = \{\theta \in \Theta | FK(\theta) = \mathbf{T}\}. \quad (6)$$

It is noted that this mapping is often non-unique, especially for redundant manipulators with more than 6 DoF.

3.2. Task Definition

Given the dataset D with language instructions l and the RGB-D observations \mathbf{o} , we aim to learn a policy $\pi_\phi(\mathbf{a}|\mathbf{o}, l)$ which predicts actions \mathbf{a} . Here, \mathbf{a} is composed of a trajectory $\mathbf{a}^{\text{joint}} = \{\mathbf{a}_0, \mathbf{a}_1, \dots, \mathbf{a}_T\}$ and gripper opening or closing action $\mathbf{a}^{\text{grripper}}$, where T denotes the trajectory length and $\mathbf{a}_i \in \mathbb{R}^n$. n represents the number of robot joints. In the bimanual manipulation, n is usually to be 12 or 14 as each robot arm has 6 or 7-DoF. Referring to prior works [15, 17, 49, 64], it is inefficient to train the policy on all trajectory points. Thus, a keyframe discovery method is used to extract a set of K_ξ keyframe indices $\{k_i\}_{i=0}^{K_\xi}$ and the prediction actions are the key set of end-effector poses \mathbf{a}^{EE} .

3.3. KStar Diffuser

3.3.1. Overview

Mainstream methods [17, 23, 65] train a policy to predict actions, but take little consideration about the mechanical robot structure which determines its motion. We thus propose a spatial-temporal graph to model both the static physical structure and dynamic history movement information. In addition, to reduce kinematically infeasible predictions for the end-effector pose, we introduce a differentiable kinematics module that provides kinematics-aware references to the policy network. The overview of our proposed Kinematics enhanced Spatial-Temporal gRaph Diffuser (KStar Diffuser) is shown in Figure B.

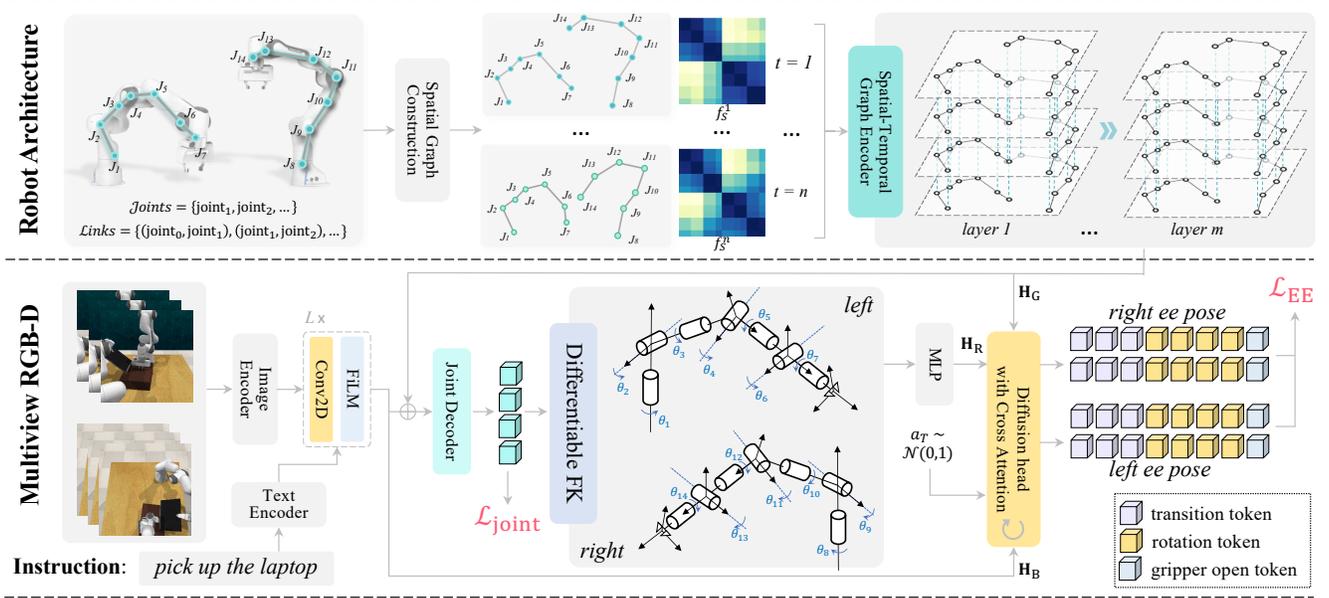


Figure 2. **Overview of KStar Diffuser.** The top part presents the spatial-temporal graph which is constructed according to the robot architecture. The bottom part shows our backbone and the proposed kinematics regularizer. For the backbone, it extracts the multimodal information which consists of multiview RGB-D observations and language instruction, and then generates bimanual 6D end-effector poses. The kinematics regularizer enhances pose learning by incorporating joint-level predictions, which are mapped to reference end-effector poses through differentiable forward kinematics (FK).

3.3.2. Backbone

Given the language instruction l and multiview RGB-D observation images o , we first adopt the Transformer-based encoders to extract their features \mathbf{E}_I and \mathbf{E}_T , respectively. Then the features are fused with the l_{FiLM} -layers FiLM block [42] to obtain the hidden states \mathbf{H}_B . Each layer is combined with an up-sampling 2D convolution layer:

$$\mathbf{H}_B^i = \text{Conv2D}(\mathbf{H}_B^i), \quad (7)$$

$$\mathbf{H}_B^{i+1} = \text{FiLM}^i(\mathbf{H}_B^i, \mathbf{E}_T), \quad (8)$$

where $i \in [0, l_{\text{FiLM}}]$ and the \mathbf{H}_B^0 is initialized with \mathbf{E}_I .

Our backbone π_ϕ uses the last hidden states \mathbf{H}_B as the condition to guide a diffusion head to denoise and generate the bimanual end-effector pose. It is noted that we attach the n historical observation images to provide more information for capturing the motion tendency. Following Chi *et al.* [9], we let the policy predict next m actions during training to alleviate multimodal problems. Here, we set both n and m to 2. The action prediction is as follows:

$$\mathbf{a}_{k-1}^{\text{EE}} \sim \pi_\phi(\mathbf{a}_{k-1}^{\text{EE}} | \mathbf{a}_k^{\text{EE}}, \mathbf{H}_B), \quad (9)$$

$$\begin{aligned} \pi_\phi(\mathbf{a}_{k-1}^{\text{EE}} | \mathbf{a}_k^{\text{EE}}, \mathbf{H}_B) \\ = \mathcal{N}(\mathbf{a}_{k-1}^{\text{EE}}; \sqrt{\bar{\alpha}_{k-1}} \pi_\phi(\mathbf{a}_k^{\text{EE}}, \mathbf{H}_B); (1 - \bar{\alpha}_{k-1})I). \end{aligned} \quad (10)$$

The learning objective is:

$$\mathcal{L}_{\text{EE}} = \mathbb{E}_q \left[\log \frac{\pi_\phi(\mathbf{a}_{0:K}^{\text{EE}} | \mathbf{H}_B)}{q(\mathbf{a}_{1:K}^{\text{EE}} | \mathbf{a}_0^{\text{EE}}, \mathbf{H}_B)} \right]. \quad (11)$$

3.3.3. Spatial-Temporal Robot Graph

The physical architecture impacts the motion of the whole robot, determining whether it can complete the task. Meanwhile, the historical spatial information is also important to future movements. Thus, we propose a spatial-temporal graph method to model the robot architecture at each step and the robot motion at continuous timesteps, representing the static spatial information and dynamic motion features.

Spatial Structure Graph Construction. To represent the robot structure, we first parse the Unified Robot Description Format (URDF) file which is usually used to describe the static physical structure of robot such as joint types, joint limits, and link lengths. Then, we define a dual-arm system as an undirected graph $G_S = \langle V_S, E_S \rangle$ based on the joint and link configuration. Here, $V_S = \{v_{j_1}, v_{j_2}, \dots, v_{j_m}\}$ and $E_S = \{e_{ij} | (\text{joint}_i, \text{joint}_j) \in \text{links}\}$ represents the node set of joints and edge set of links respectively, where m denotes the number of joint. We use $\mathbf{f}_i \in \mathbb{R}^D$ to denote the value of i -th node feature. It consists of three attributes as follows:

- **Joint Coordinate:** we use a vector in Cartesian coordinate system, $\mathbf{f}_i^{\text{JC}} = [x_i, y_i, z_i] \in \mathbb{R}^3$, to denote the absolute coordinates of the i -th joint. The vector is normalized according to the workspace boundary for stable convergence of model training.
- **Joint Distance:** To measure the spatial relationship between the node v_i and the other nodes v_j , we compute the Euclidean Distance between v_i and v_j :

$$\mathbf{f}_i^{\text{JD}} = \|\mathbf{v}_i - \mathbf{v}_j\|_2, \quad (12)$$

where $\mathbf{f}_i^{\text{JD}} \in \mathbb{R}^m$ and $\|\cdot\|_2$ means the Euclidean norm.

- *Body Label*: To discriminate the source of node \mathbf{v}_i , we use a one-hot vector $\mathbf{f}_i^{\text{BL}} = [0, 1]$ as one of its features. It can also help the policy capture the motion modes of different robot arms such as the symmetry.

We concatenate \mathbf{f}_i^{JC} , \mathbf{f}_i^{JD} , and \mathbf{f}_i^{BL} to form the completed feature $\mathbf{f}_i \in \mathbb{R}^{D_f}$, where D_f denotes the dimension of node features.

Spatial-Temporal Graph Learning. Given the same instruction and observation, different historical robot poses result in different predictions. Thus, we build the spatial structure graph with the temporal motion information. Specifically, we build a spatial-temporal graph $G_{\text{ST}} = \langle V_{\text{ST}}, E_{\text{ST}} \rangle$ by combining $\{G_S^i = \langle V_S^i, E_S^i \rangle\}_{i=0}^{T-1}$ from historical timesteps, where T denotes the number of history steps. In G_{ST} , the node set V_{ST} contains V_S^i from historical static spatial graph. Additionally, we add edges E' linking the same joint node joint_i^t at different timesteps, in order to establish the correlation of joint motions at continuous times. It can be formulated as follows:

$$V_{\text{ST}} = \cup \{V_S^i\}_{i=0}^{T-1}, \quad E_{\text{ST}} = \cup \{E_S^i\}_{i=0}^{T-1} \cup E', \quad (13)$$

$$E' = \{e_i^{t,t'} | t, t' \in \{0, 1, \dots, T\}, t \neq t'\}. \quad (14)$$

In this way, we obtain the whole spatial-temporal graph. A Graph Convolutional Network (GCN) is then adopted to propagate and aggregate node features across the graph. The GCN layer updates each node feature \mathbf{f}_i by aggregating the features of its neighboring nodes, thereby capturing the relational and structural information of the robotic arms. We use the node feature \mathbf{H}_G of the last encoder layer as the representation of the robot structure to condition the denoising process.

3.3.4. Kinematics Regularizer

To control the end-effector effectively, the generated pose trajectory must be processed by an Inverse Kinematics (IK) solver, which calculates the joint configurations to achieve the specified poses. However, because the predicted trajectory is generated without considering the robot kinematic constraints, it often falls outside the IK solver’s feasible range, resulting in high failure rates during execution. To address this limitation, we propose a kinematics regularizer into the end-effector pose learning objective. This regularizer aligns the predicted poses with the robot kinematic constraints, ensuring that the generated trajectory remains within the solvable space of the IK solver, thus enhancing then reliability of trajectory execution.

Differentiable Kinematics. Given a joint configuration $\Theta = [\theta_1, \theta_2, \dots, \theta_n]$, the corresponding end-effector pose $\mathbf{T} \in SE(3)$ can be computed using forward kinematics,

represented as a mapping $f(\Theta) : \mathbb{R}^n \rightarrow SE(3)$. This mapping from joint space to end-effector space is differentiable, namely Differentiable Forward Kinematics (DFK), enabling the use of gradients to optimize our control policy. Leveraging DFK, our policy learns to predict the next joint configuration $\hat{\mathbf{a}}^{\text{joint}}$, from which we compute an intermediate end-effector pose \mathbf{H}_R . By using \mathbf{H}_R as a reference, we guide a denoising process to generate a precise and executable end-effector pose.

Specifically, we combine structure features \mathbf{H}_G with the last hidden state \mathbf{H}_B , projecting to the joint space and using the DFK to obtain reference as \mathbf{H}_R follows:

$$\hat{\mathbf{a}}^{\text{joint}} = \text{Proj}([\mathbf{H}_B, \mathbf{H}_G]), \quad \mathbf{H}_R = \text{DFK}(\hat{\mathbf{a}}^{\text{joint}}). \quad (15)$$

To ensure the consistency between predicted and actual joint angles, we minimize the joint loss:

$$\mathcal{L}_{\text{joint}} = \mathbb{E}_{\mathbf{a}^{\text{joint}} \sim q(\mathbf{a}_0)} [\|\mathbf{a}^{\text{joint}} - \hat{\mathbf{a}}^{\text{joint}}\|^2]. \quad (16)$$

Conditioning Diffusion Process on Kinematics. To enforce kinematic consistency, we condition the diffusion process on the reference representation \mathbf{H}_R , an auxiliary input encoding kinematic constraints. This allows the predicted pose trajectory to stay within feasible space. Given the diffusion steps from Eq.(9) and Eq.(10), we have:

$$\mathbf{a}_{k-1}^{\text{EE}} \sim \pi_\phi(\mathbf{a}_k^{\text{EE}} | \mathbf{a}_{k-1}^{\text{EE}}, \mathbf{H}_B, \mathbf{H}_R), \quad (17)$$

$$\begin{aligned} & \pi_\phi(\mathbf{a}_{k-1}^{\text{EE}} | \mathbf{a}_k^{\text{EE}}, \mathbf{H}_B, \mathbf{H}_R) = \\ & \mathcal{N}(\mathbf{a}_{k-1}^{\text{EE}}; \sqrt{\bar{\alpha}_{k-1}} \pi_\phi(\mathbf{a}_k^{\text{EE}}, \mathbf{H}_B, \mathbf{H}_R); (1 - \bar{\alpha}_{k-1})\mathbf{I}). \end{aligned} \quad (18)$$

Incorporating DFK into the diffusion process allows gradients from the pose loss to propagate back through the kinematic function, ensuring that each denoising step maintains compliance with joint constraints, thereby optimizing the end-effector’s control accuracy and robustness.

3.4. Training and Inference

Training. We use the conditional action generation mode to train KStar Diffuser, which is formulated as a conditional denoising diffusion. The loss function is defined as the Mean Square Error (MSE) as follows:

$$\mathbf{a}_k = \sqrt{\bar{\alpha}_k} \mathbf{a}_0 + \sqrt{1 - \bar{\alpha}_k} \epsilon, \quad (19)$$

$$\mathcal{L}_{\text{EE}} = \mathbb{E}_{\mathbf{a}_0^{\text{EE}} \sim q(\mathbf{a}_0), k} [\|\mathbf{a}_0^{\text{EE}} - \pi_\phi(\mathbf{a}_k^{\text{EE}}, k, C)\|^2], \quad (20)$$

$$\mathcal{L}_{\text{joint}} = \mathbb{E}_{\mathbf{a}_0^{\text{joint}} \sim q(\mathbf{a}_0)} [\|\mathbf{a}_0^{\text{joint}} - \pi_\phi(\mathbf{H}_B, \mathbf{H}_G)\|^2], \quad (21)$$

$$\mathcal{L} = \lambda \mathcal{L}_{\text{EE}} + (1 - \lambda) \mathcal{L}_{\text{joint}}, \quad (22)$$

where \mathbf{a}_k is obtained by the forward diffusion process and C is the combination of \mathbf{H}_B , \mathbf{H}_G , \mathbf{H}_R . λ is the trade-off coefficient.

Inference. Sampling from a Gaussian noise ϵ_k , the policy π_ϕ performs K iterations to gradually denoise a random noise ϵ_k into the noise-free action \mathbf{a}_0 :

$$\mathbf{a}_{k-1}^{\text{EE}} = \sqrt{\bar{\alpha}_{k-1}} \pi_\phi(\mathbf{a}_k^{\text{EE}}, k, C) + \sqrt{1 - \bar{\alpha}_{k-1}} \epsilon_k. \quad (23)$$

We use the predicted \mathbf{a}_0^{EE} as the final action to control robot execution.

4. Experiment

4.1. Dataset and Evaluation Settings

Dataset. Bimanual manipulation tasks demand high levels of coordination, synchronization, and symmetry awareness between the two robotic arms, making them inherently more challenging than single-arm tasks. To assess the capabilities of KStar Diffuser in these areas, we conducted comprehensive experiments using the RL Bench2 benchmark [17], an extended version of RL Bench tailored for bimanual manipulation and comprising tasks closely resembling real-world scenarios. Please refer to Appendix A for more details about RL Bench2 and real-world tasks.

Evaluation Settings. To evaluate the policy performance, we employ success rate as the primary metric. Although the policy generates multiple sequential actions during execution, we primarily focus on the final goal achievement rather than intermediate steps. Each task is associated with a specific success criterion defined by its target state. To comprehensively assess the policy’s capability, we conduct experiments with varying numbers of demonstrations (20 and 100) during training. Figure 3 illustrates our experimental setup, including both the simulation environment and the Cobot Agilix ALOHA robot. Detailed descriptions of simulation tasks and real-world experimental settings are provided in the Appendix B.

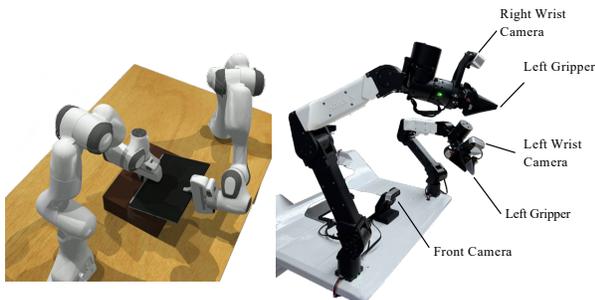


Figure 3. *The left:* the simulation environment of `pick_laptop` task. *The right:* the ALOHA device used in the real-world tasks.

4.2. Baselines

We systematically evaluate KStar Diffuser against state-of-the-art methods in two major categories:

Transformer-based methods: (1) Action Chunking with Transformers (**ACT**) [65] employs a Conditional VAE architecture, consisting of an encoder-decoder framework for joint angle sequence prediction; (2) Robotic View Transformer Leader Following (**RVT-LF**) [17] leverages RVT [15] as its backbone, incorporating a multi-view transformer for cross-view information aggregation and image re-rendering, coupled with a leader following mechanism for action prediction; (3) Perceiver-Actor Leader Following (**PerAct-LF**) [17] adopts the leader following paradigm based on PerAct [49], utilizing a perceiver Transformer to encode both instructions and voxel observations for optimal voxel action generation; (4) **PerAct2** [17] enhances PerAct by implementing a unified feature space for dual-arm actions and employing combined self-attention for synchronized bimanual action prediction.

Diffusion-based methods: (1) Joint-based Diffusion Policy (**DP-J**) [9] adopts a diffusion model to robotic manipulation within an imitation learning framework, focusing on joint angle prediction; (2) EndEffector-based Diffusion Policy (**DP-EE**) is a DP variant which we reimplement Diffusion Policy to predict end-effector poses instead of joint angles, offering an alternative control paradigm; (3) 3D Diffusion Policy (**DP3**) [61] enhances 3D perception by incorporating point clouds for joint angle prediction.

4.3. Comparison Results with SOTA Methods

Experimental Results on RL Bench2. As shown in Table 1, the KStar Diffuser significantly outperforms other state-of-the-art baselines, achieving more than 20% higher overall performance with both 20 and 100 training demonstrations. We found that:

(1) **Similar to learning a single-arm policy**, the process of learning a bimanual policy can adapt quickly and achieve a high success rate, given a relatively consistent distribution of task trajectories. In `push_box`, where the objective is for both arms to push a box toward a specified target along fixed trajectories, our KStar Diffuser and other baseline models perform well. However, as task complexity increases, success rates decrease. For example, in the `lift_ball` task, both arms must lift a large ball *simultaneously* to complete the task. Any asynchrony in movement can cause instability, leading to the ball slipping and ultimately resulting in task failure. Our KStar Diffuser achieves its robust performance on such bimanual tasks by **explicitly modeling the spatial and motion relationships between the two arms**, surpassing other methods more than 6%.

(2) **Distinct from single-arm systems**, bimanual robotic systems possess the capability for **collaborative manipulation**. Methods which are directly adapted from single-arm to bimanual manipulation exhibit high failure rates in tasks, e.g., `pick_laptop`, as they lack consideration of the spatial and motion relationships between the arms. Specifi-

Table 1. The experimental result on simulated tasks. We train the policy with the setting of different training demonstrations, *i.e.* [20, 100], to test its capability comprehensively with 100 times each task. The best results are in **bold**. Each result is reported with three seeds on average.

Method	Backbone	Push Box	Lift Ball	Handover Item (easy)	Pick Laptop	Sweep Dustpan	Overall
20 demos							
ACT [65]	Transformer	26.0 \pm 4.6	17.3 \pm 5.8	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0	8.7 \pm 1.6
RVT-LF [17]	Transformer	44.0 \pm 1.0	12.3 \pm 6.7	0.0 \pm 0.0	2.0 \pm 1.7	0.0 \pm 0.0	11.7 \pm 1.5
PerAct-LF [17]	Transformer	63.7 \pm 3.1	42.3 \pm 8.1	3.3 \pm 0.6	7.0 \pm 2.7	4.7 \pm 0.6	24.2 \pm 1.6
PerAct2 [17]	Transformer	59.7 \pm 7.6	37.7 \pm 4.7	8.0 \pm 1.0	17.3 \pm 0.6	1.7 \pm 1.5	24.9 \pm 2.1
DP-J [9]	Diffusion	30.0 \pm 3.6	19.7 \pm 2.1	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 4.0	10.0 \pm 1.0
DP-EE [9]	Diffusion	34.7 \pm 2.5	34.3 \pm 7.4	0.7 \pm 0.6	1.3 \pm 0.6	14.0 \pm 1.7	17.0 \pm 1.6
DP3 [61]	Diffusion	25.3 \pm 4.0	34.3 \pm 7.6	0.0 \pm 0.0	2.7 \pm 2.9	0.0 \pm 0.0	12.5 \pm 1.3
KStar Diffuser (Ours)	Diffusion	79.3 \pm 3.5	87.0 \pm 2.7	23.7 \pm 0.6	17.0 \pm 2.0	83.0 \pm 4.4	58.0 \pm 1.4
100 demos							
ACT [65]	Transformer	48.7 \pm 6.8	40.7 \pm 2.1	0.0 \pm 0.0	0.0 \pm 0.0	0.0 \pm 0.0	17.9 \pm 1.0
RVT-LF [17]	Transformer	76.3 \pm 1.5	28.7 \pm 1.0	0.0 \pm 0.0	1.3 \pm 1.2	2.3 \pm 3.2	21.7 \pm 0.7
PerAct-LF [17]	Transformer	66.3 \pm 3.1	68.3 \pm 6.1	7.3 \pm 2.1	14.3 \pm 3.8	8.0 \pm 2.7	32.9 \pm 2.5
PerAct2 [17]	Transformer	77.0 \pm 3.0	47.7 \pm 3.2	10.3 \pm 0.6	34.3 \pm 4.0	6.7 \pm 2.1	33.9 \pm 1.0
DP-J [9]	Diffusion	62.7 \pm 7.6	43.3 \pm 2.1	0.0 \pm 0.0	1.3 \pm 0.6	0.0 \pm 4.0	21.5 \pm 1.6
DP-EE [9]	Diffusion	61.0 \pm 2.0	59.7 \pm 3.5	2.0 \pm 1.0	10.3 \pm 3.1	69.7 \pm 2.5	40.5 \pm 1.2
DP3 [61]	Diffusion	56.0 \pm 3.6	64.0 \pm 2.7	0.0 \pm 0.0	6.3 \pm 3.1	1.7 \pm 2.1	25.6 \pm 1.4
KStar Diffuser (Ours)	Diffusion	83.0 \pm 1.7	98.7 \pm 1.5	27.0 \pm 1.7	43.7 \pm 4.5	89.0 \pm 5.2	68.2 \pm 2.1

cally, as shown in Figure 4, this task involves picking up a notebook lying flat on a cabinet surface. Given that the notebook rests fully against the tabletop, direct grasping by the robotic arm is not possible. Instead, an effective strategy is to control one arm to push the notebook outward from the cabinet by a short distance, allowing the other arm to pick it up. KStar Diffuser achieves a success rate approximately 9% higher than other methods, demonstrating its ability to **capture the coordinated motion patterns** required for collaborative object manipulation between two arms.

Real-world Experimental Results. To comprehensively evaluate the policy’s effectiveness, we build 2 tasks in real-world based on the simulation benchmark. The performances on real-world tasks are illustrated in Table 2.

Similar to the simulation results, we observe that policies which do not consider about bimanual scenarios, *i.e.*, ACT, DP, and DP3, demonstrate limited capability across all bimanual tasks, achieving around 20% success ratio on average. Although PerAct2 is designed for bimanual tasks via mapping bimanual actions into a shared learning space, it fails to capture the spatial structure of the bimanual system, leading to ineffective arm coordination during execution. Furthermore, we also found PerAct2 faces *significant inverse kinematics issues* with its predicted end-effector poses, including joint configuration conflicts and unreachable positions which are shown in Figure 4. It is likely due to PerAct2’s limited ability to capture the complex spatial

constraints and kinematic relationships within the bimanual robotic system. In contrast, KStar Diffuser achieves superior bimanual coordination, surpassing other methods by over 10%, as it successfully captures the motion patterns between dual arms and predicts feasible end-effector poses.

Table 2. The result of real-world tasks. We train all policies with 100 demonstrations and test 15 times. The best result are in **bold**.

Methods	Lift Plate	Handover	Overall
ACT [65]	37.8 \pm 8.3	0.0 \pm 0.0	18.9 \pm 11.2
DP [9]	42.1 \pm 6.5	0.0 \pm 0.0	21.1 \pm 12.5
DP3 [61]	44.0 \pm 8.3	0.0 \pm 0.0	22.0 \pm 13.4
PerAct2 [17]	51.1 \pm 6.3	8.8 \pm 3.0	29.9 \pm 10.2
KStar Diffuser (Ours)	66.7 \pm 5.3	19.7 \pm 5.3	43.1 \pm 17.8

Table 3. Ablation study of model components.

ST Graph	KR	Handover-S	Handover-R	Overall
✓	✓	27.0 \pm 1.7	19.7 \pm 5.3	23.4 \pm 5.2
✓	×	18.3 \pm 1.5	15.3 \pm 3.3	16.8 \pm 2.1
×	×	16.3 \pm 3.1	13.3 \pm 9.4	14.8 \pm 2.1

4.4. Ablation Studies

Effects of Model Components. To systematically evaluate the contribution of each component in KStar Diffuser,

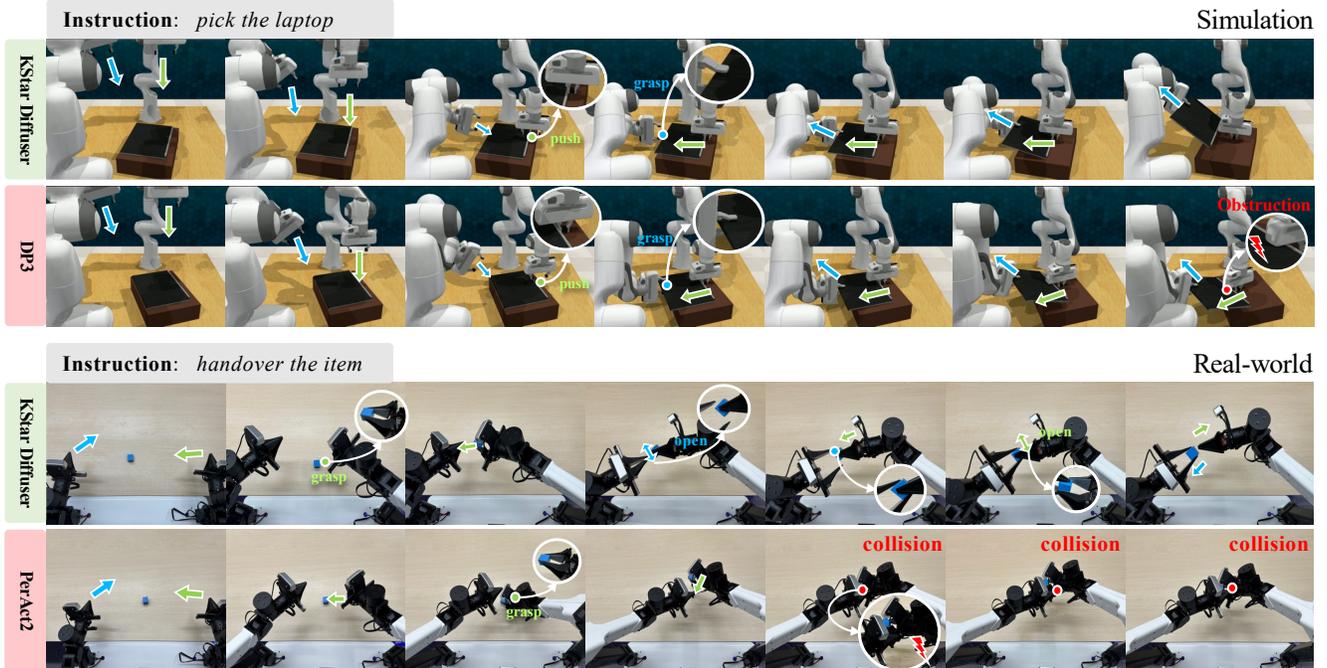


Figure 4. The visualization of bimanual manipulation on simulated RL-Bench2 and real-world tasks. The blue annotations represent the motion of the robot’s left arm, while the green annotations indicate the motion of the right arm.

we conduct ablation experiments on the `handover_item` task in both simulated and real-world environments. We design a progressive ablation process by first removing the Differential Forward Kinematics module while retaining the Spatial-Temporal Graph (ST Graph), and then completely disabling both ST Graph and Kinematics Regularizer (KR).

The experimental results in Table 3 demonstrate the crucial role of each component. The removal of KR leads to a significant decrease in success rate, particularly pronounced in real-world scenarios. This performance degradation can be attributed to the fundamental differences between simulated and real environments. While simulated environments maintain consistent and noise-free inputs, real-world scenarios introduce various perturbations, *e.g.*, sensor noise and light reflection, making the policy more susceptible to kinematic constraint violations without the regularization effect of KR. Further ablation by removing both ST Graph and KR results in a substantial performance drop across all experimental settings. This observation illustrates two key aspects: First, the ST Graph effectively captures the spatial-temporal dependencies among joints, which is essential for coordinating the relative positioning and interactions between robotic arms. Second, the graph structure’s explicit encoding of the robotic physical architecture enhances the policy’s robustness against unexpected perturbations by maintaining spatial and temporal coherence. We conduct extensive ablation studies of action chunking size,

history lengths and trade-off coefficient, please refer to Appendix C for more details.

4.5. Qualitative Analysis

We further present qualitative analysis in Figure 4. We compare the performance of KStar Diffuser with DP3 and PerAct2 in executing bimanual manipulation tasks within both simulated and real-world environments, respectively.

In the simulation task, since the laptop rests flat on the cabinet, direct lifting is not feasible. One robotic arm initiates a forward push, creating space, while the other arm concurrently grasps and elevates the laptop. KStar Diffuser effectively models this dual-arm coordination, generating a precise trajectory of synchronized actions. Conversely, DP3, adapted from a single-arm policy to a dual-arm configuration, fails to achieve effective coordination. Concretely, after executing the push motion, the right arm does not halt, obstructing the left arm’s lifting process.

In the real-world task, KStar Diffuser generates an executable item transfer trajectory between the left and right arms, with no collisions throughout the task, reflecting its strong environmental adaptability and collision avoidance capabilities. Conversely, PerAct2 encounters collisions during the handover process (marked in red), indicating less effective handling of dynamic real-world variables and a lack of kinematic awareness on the robot movements. More qualitative analysis can be found in Appendix D.

5. Conclusion

In this paper, we have proposed a novel Kinematics enhanced Spatial-Temporal gGraph Diffuser (KStar Diffuser) that explicitly incorporates both robot structures and kinematics into the bimanual motion generation process. It consists of a spatial-temporal robot graph that explicitly models the robot physical configuration to guide the generative action denoising procedure, and a kinematics regularizer that augments the NBP learning objective by introducing joint-space supervision. Extensive experiments demonstrate that KStar Diffuser outperforms the baselines by a large margin on both simulation and real-world tasks.

Limitations and Future Directions. While we explored robot structure impacts through GNN modeling and kinematic constraints, the core control logic of end-effector pose prediction and inverse kinematics remains. In the future, we aim to leverage neural networks to directly model joint movements, aligning the robot motion space with the Cartesian space of the human world.

References

- [1] Lucia Angelini, Manuela Uliano, Angela Mazzeo, Mattia Penzotti, and Marco Controzzi. Self-collision avoidance in bimanual teleoperation using collisionnik: algorithm revision and usability experiment. In *2022 IEEE-RAS 21st International Conference on Humanoid Robots (Humanoids)*, pages 112–118, 2022. 2
- [2] Andreas Aristidou, Joan Lasenby, Yiorgos Chrysanthou, and Ariel Shamir. Inverse kinematics techniques in computer graphics: A survey. In *Computer graphics forum*, pages 35–58. Wiley Online Library, 2018. 2
- [3] Weibang Bai, Ningshan Zhang, Baoru Huang, Ziwei Wang, Francesco Cursi, Ya-Yen Tsai, Bo Xiao, and Eric M Yeatman. Dual-arm coordinated manipulation for object twisting with human intelligence. In *2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 902–908. IEEE, 2021. 2
- [4] Aaron M Bestick, Samuel A Burden, Giorgia Willits, Nikhil Naikal, S Shankar Sastry, and Ruzena Bajcsy. Personalized kinematics for human-robot collaborative manipulation. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1037–1044. IEEE, 2015. 2
- [5] Botond Bócsi, Duy Nguyen-Tuong, Lehel Csató, Bernhard Schoelkopf, and Jan Peters. Learning inverse kinematics with structured prediction. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 698–703. IEEE, 2011. 3
- [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. 1
- [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
- [8] Shizhe Chen, Ricardo Garcia, Cordelia Schmid, and Ivan Laptev. Polarnet: 3d point clouds for language-guided robotic manipulation. *arXiv preprint arXiv:2309.15596*, 2023. 2
- [9] 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. 1, 2, 3, 4, 6, 7, 14
- [10] Cheng Chi, Zhenjia Xu, Chuer Pan, Eric Cousineau, Benjamin Burchfiel, Siyuan Feng, Russ Tedrake, and Shuran Song. Universal manipulation interface: In-the-wild robot teaching without in-the-wild robots. *arXiv preprint arXiv:2402.10329*, 2024. 1, 2, 3
- [11] Akos Csaszar, Jan Eilers, and Alexander Verl. On solving the inverse kinematics problem using neural networks. In *2017 24th International Conference on Mechatronics and Machine Vision in Practice (M2VIP)*, pages 1–6, 2017. 3
- [12] Aaron D’Souza, Sethu Vijayakumar, and Stefan Schaal. Learning inverse kinematics. In *Proceedings 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems. Expanding the Societal Role of Robotics in the the Next Millennium (Cat. No. 01CH37180)*, pages 298–303. IEEE, 2001. 3
- [13] Zipeng Fu, Tony Z Zhao, and Chelsea Finn. Mobile aloha: Learning bimanual mobile manipulation with low-cost whole-body teleoperation. *arXiv preprint arXiv:2401.02117*, 2024. 1, 3
- [14] Jianfeng Gao, Xiaoshu Jin, Franziska Krebs, Noémie Jaquier, and Tamim Asfour. Bi-kvil: Keypoints-based visual imitation learning of bimanual manipulation tasks. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 16850–16857. IEEE, 2024. 2
- [15] Ankit Goyal, Jie Xu, Yijie Guo, Valts Blukis, Yu-Wei Chao, and Dieter Fox. Rvt: Robotic view transformer for 3d object manipulation. In *Conference on Robot Learning*, pages 694–710. PMLR, 2023. 1, 2, 3, 6
- [16] Reinhard Grassmann, Vincent Modes, and Jessica Burgner-Kahrs. Learning the forward and inverse kinematics of a 6-dof concentric tube continuum robot in se (3). In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5125–5132. IEEE, 2018. 2
- [17] Markus Grotz, Mohit Shridhar, Tamim Asfour, and Dieter Fox. Peract2: A perceiver actor framework for bimanual manipulation tasks. *arXiv preprint arXiv:2407.00278*, 2024. 1, 2, 3, 6, 7, 13
- [18] Abdullah Aamir Hayat, Ratan OM Sadanand, and Subir K Saha. Robot manipulation through inverse kinematics. In *Proceedings of the 2015 conference on advances in robotics*, pages 1–6, 2015. 2
- [19] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [20] Binghao Huang, Yixuan Wang, Xinyi Yang, Yiyue Luo, and Yunzhu Li. 3d-vitac: Learning fine-grained manipulation

- with visuo-tactile sensing. *arXiv preprint arXiv:2410.24091*, 2024. 2
- [21] Kenneth H Hunt. Structural kinematics of in-parallel-actuated robot-arms. 1983. 1
- [22] Michael I Jordan and David E Rumelhart. Forward models: Supervised learning with a distal teacher. In *Backpropagation*, pages 189–236. Psychology Press, 2013. 3
- [23] Tsung-Wei Ke, Nikolaos Gkanatsios, and Katerina Fragkiadaki. 3d diffuser actor: Policy diffusion with 3d scene representations. *arXiv preprint arXiv:2402.10885*, 2024. 2, 3
- [24] 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. 2
- [25] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016. 2
- [26] CS George Lee. Robot arm kinematics, dynamics, and control. *Computer*, 15(12):62–80, 1982. 1
- [27] Maolin Lei, Ting Wang, Chen Yao, Huan Liu, Zhi Wang, and Yongsheng Deng. Real-time kinematics-based self-collision avoidance algorithm for dual-arm robots. *Applied Sciences*, 10(17):5893, 2020. 1
- [28] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3383–3393, 2021. 3
- [29] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3383–3393, 2021. 3
- [30] Xiaojie Li, Shaowei He, Jianlong Wu, Yue Yu, Liqiang Nie, and Min Zhang. Mask again: Masked knowledge distillation for masked video modeling. In *Proceedings of the ACM International Conference on Multimedia*, page 2221–2232. ACM, 2023. 2
- [31] Xiaojie Li, Jianlong Wu, Shaowei He, Shuo Kang, Yue Yu, Liqiang Nie, and Min Zhang. Fine-grained key-value memory enhanced predictor for video representation learning. In *Proceedings of the ACM International Conference on Multimedia*, page 2264–2274. ACM, 2023.
- [32] Xiaojie Li, Yibo Yang, Xiangtai Li, Jianlong Wu, Yue Yu, Bernard Ghanem, and Min Zhang. Genview: Enhancing view quality with pretrained generative model for self-supervised learning. In *Proceedings of the European Conference on Computer Vision*. Springer, 2024. 2
- [33] I Liu, Chun Arthur, Sicheng He, Daniel Seita, and Gaurav Sukhatme. Voxact-b: Voxel-based acting and stabilizing policy for bimanual manipulation. *arXiv preprint arXiv:2407.04152*, 2024. 2
- [34] Qi Lv, Xiang Deng, Gongwei Chen, Michael Yu Wang, and Liqiang Nie. Decision mamba: A multi-grained state space model with self-evolution regularization for offline rl. In *Advances in Neural Information Processing Systems*, pages 22827–22849. Curran Associates, Inc., 2024. 1
- [35] Xiao Ma, Sumit Patidar, Iain Haughton, and Stephen James. Hierarchical diffusion policy for kinematics-aware multi-task robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18081–18090, 2024. 1, 2
- [36] Matthew T Mason. *Mechanics of robotic manipulation*. MIT press, 2001. 2
- [37] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of thought. *Advances in Neural Information Processing Systems*, 36, 2024. 1
- [38] Yoshihiko Nakamura, Kiyoshi Nagai, and Tsuneo Yoshikawa. Dynamics and stability in coordination of multiple robotic mechanisms. *The International Journal of Robotics Research*, 8(2):44–61, 1989. 1
- [39] Valerio Ortenzi, Naresh Marturi, Michael Mistry, Jeffrey Kuo, and Rustam Stolkin. Vision-based framework to estimate robot configuration and kinematic constraints. *IEEE/ASME Transactions on Mechatronics*, 23(5):2402–2412, 2018. 2
- [40] Eric Paljug, Xiaoping Yun, and Vijay Kumar. Control of rolling contacts in multi-arm manipulation. *IEEE Transactions on Robotics and Automation*, 10(4):441–452, 1994. 1
- [41] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023. 2
- [42] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 4
- [43] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2
- [44] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 2
- [45] Allen Z Ren, Justin Lidard, Lars L Ankile, Anthony Simeonov, Pulkit Agrawal, Anirudha Majumdar, Benjamin Burchfiel, Hongkai Dai, and Max Simchowitz. Diffusion policy policy optimization. *arXiv preprint arXiv:2409.00588*, 2024. 2
- [46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2
- [47] Nilanjan Sarkar, Xiaoping Yun, and Vijay Kumar. Dynamic control of 3-d rolling contacts in two-arm manipulation. *IEEE Transactions on Robotics and Automation*, 13(3):364–376, 1997. 1

- [48] Youshik Shin and Zeungnam Bien. Collision-free trajectory planning for two robot arms. *Robotica*, 7(3):205–212, 1989. 1, 2
- [49] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Conference on Robot Learning*, pages 785–799. PMLR, 2023. 1, 2, 3, 6, 13
- [50] Christian Smith, Yiannis Karayiannidis, Lazaros Nalpan-tidis, Xavi Gratal, Peng Qi, Dimos V Dimarogonas, and Dan-ica Kragic. Dual arm manipulation—a survey. *Robotics and Autonomous systems*, 60(10):1340–1353, 2012. 1
- [51] 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. 2
- [52] Vitalis Vosylius, Younggyo Seo, Jafar Uruç, and Stephen James. Render and diffuse: Aligning image and action spaces for diffusion-based behaviour cloning. *arXiv preprint arXiv:2405.18196*, 2024. 2
- [53] Lirui Wang, Xinlei Chen, Jialiang Zhao, and Kaiming He. Scaling proprioceptive-visual learning with heterogeneous pre-trained transformers. *arXiv preprint arXiv:2409.20537*, 2024. 2
- [54] Yixiao Wang, Yifei Zhang, Mingxiao Huo, Ran Tian, Xi-ang Zhang, Yichen Xie, Chenfeng Xu, Pengliang Ji, Wei Zhan, Mingyu Ding, et al. Sparse diffusion policy: A sparse, reusable, and flexible policy for robot learning. *arXiv preprint arXiv:2407.01531*, 2024. 2
- [55] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distilla-tion. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [56] Chris Welman. Inverse kinematics and geometric constraints for articulated figure manipulation. 1993. 2
- [57] Kun Wu, Yichen Zhu, Jinming Li, Junjie Wen, Ning Liu, Zhiyuan Xu, Qinru Qiu, and Jian Tang. Discrete policy: Learning disentangled action space for multi-task robotic manipulation. *arXiv preprint arXiv:2409.18707*, 2024. 2
- [58] Wenke Xia, Dong Wang, Xincheng Pang, Zhigang Wang, Bin Zhao, Di Hu, and Xuelong Li. Kinematic-aware prompt-ing for generalizable articulated object manipulation with llms. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2073–2080. IEEE, 2024. 2
- [59] Andrea Maria Zanchettin, Paolo Rocco, Luca Bascetta, Ioanis Symeonidis, and Steffen Peldschus. Kinematic analysis and synthesis of the human arm motion during a manipula-tion task. In *2011 IEEE international conference on robotics and automation*, pages 2692–2697. IEEE, 2011. 2
- [60] Yanjie Ze, Zixuan Chen, Wenhao Wang, Tianyi Chen, Xialin He, Ying Yuan, Xue Bin Peng, and Jiajun Wu. Generalizable humanoid manipulation with improved 3d diffusion policies. *arXiv preprint arXiv:2410.10803*, 2024. 1, 2
- [61] Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, Muhan Wang, and Huazhe Xu. 3d diffusion policy. *arXiv preprint arXiv:2403.03954*, 2024. 1, 2, 6, 7
- [62] Haoyu Zhang, Meng Liu, Yuhong Li, Ming Yan, Zan Gao, Xiaojun Chang, and Liqiang Nie. Attribute-guided collabor-ative learning for partial person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12):14144–14160, 2023. 2
- [63] Haoyu Zhang, Meng Liu, Zixin Liu, Xuemeng Song, Yaowei Wang, and Liqiang Nie. Multi-factor adaptive vision selec-tion for egocentric video question answering. In *Proceedings of the 41st International Conference on Machine Learning*, pages 59310–59328. PMLR, 2024. 2
- [64] Junjie Zhang, Chenjia Bai, Haoran He, Wenke Xia, Zhigang Wang, Bin Zhao, Xiu Li, and Xuelong Li. Sam-e: Lever-aging visual foundation model with sequence imitation for embodied manipulation. *arXiv preprint arXiv:2405.19586*, 2024. 2, 3
- [65] 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, 2, 3, 6, 7, 14
- [66] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all, march 2024. URL <https://github.com/hpcaitech/Open-Sora>, 1(3):4. 2
- [67] Minjie Zhu, Yichen Zhu, Jinming Li, Junjie Wen, Zhiyuan Xu, Ning Liu, Ran Cheng, Chaomin Shen, Yaxin Peng, Feifei Feng, et al. Scaling diffusion policy in transformer to 1 billion parameters for robotic manipulation. *arXiv preprint arXiv:2409.14411*, 2024. 2

Spatial-Temporal Graph Diffusion Policy with Kinematic Modeling for Bimanual Robotic Manipulation

Supplementary Material

We provide a more comprehensive **tasks descriptions** in Section A, encompassing both simulated environments and real-world scenarios. In Section B, we elaborate on the **implementation details**, including code base of baseline methods, and hyperparameter configurations for the backbone, graph encoder modules and optimization process. Notably, to thoroughly validate efficacy of KStar Diffuser, we conduct **extensive ablation studies** analyzing the impact of demonstration quantity, action chunking size, observation history length, and the control learning objective coefficient λ . The result is presented in Section C. Finally, we present more **qualitative analysis** in Section D.

A. Task Descriptions

We conduct extensive experiments on both simulated tasks and real-world tasks to evaluate the effectiveness of our proposed KStar Diffuser. Specifically, we selected five tasks from RL Bench2, ranging from basic symmetrical tasks to advanced coordination-requiring tasks, including `push_box`, `lift_ball`, `handover_item_easy`, `sweep_dustpan`, and `pick_laptop`. For real-world evaluations, we created a similar setup with two bimanual tasks: `lift_plate` and `handover_item_easy`. We present the details about both simulated and real-world tasks in Table A. For each simulated task, we evaluate the model 100 times, whereas for each real-world task, we conduct 15 evaluations.

Table A. Tasks Details.

Task	Duration	# Keyframes	Instruction
Simulated Tasks			
<code>push_box</code>	4.33s	2.1	"Push the box to the red area."
<code>lift_ball</code>	4.40s	4.0	"Lift the ball."
<code>handover_item_easy</code>	7.17s	7.5	"Handover the item."
<code>sweep_dustpan</code>	4.93s	7.3	"Sweep the dust to the pan."
<code>pick_laptop</code>	3.97s	7.2	"Pick up the notebook."
Real-world Tasks			
<code>lift_plate</code>	6.37s	3.4	"Lift the plate."
<code>handover_item_easy</code>	9.52s	8.6	"Handover the item."

A.1. Simulated Tasks

push_box. As illustrated in Figure A(a), the task requires the robot to utilize both arms to push a heavy box, weighing 50 kg, and transport it to a designated target area through a fix trajectory. The completion of the task is defined as successfully moving the box to the specified location. The scenario involves two key elements: a large box and a target area, with the primary challenge being the considerable weight of the box, which exceeds the capacity of

a single arm to manage effectively. Notably, this task necessitates the use of both arms simultaneously, as a single robot is incapable of accomplishing it independently.

lift_ball. As depicted in Figure A(b), the task entails the robot using both arms to lift a large ball, achieving a minimum height of 0.95 meters to meet the success criteria. The object in this task is the large ball, presenting a significant coordination challenge. Due to the ball's size and the inability of the gripper to securely grasp it, the operation relies on coordinated non-prehensile manipulation, demanding precise synchronization of the arms during the lifting process. This task cannot be accomplished by a single robot because of the object's dimensions.

handover_item_easy. The task requires the robot to handover a red item by utilizing one arm to securely grasp and lift the item to a height of 80 cm while ensuring the other arm remains idle and unengaged, as shown in Figure A(c). The object involved is a single red block, and the key challenge lies in coordinating the handover process effectively. Successful completion is determined when the item is accurately identified, grasped, and positioned at the required height with no actions performed by the idle arm.

sweep_dustpan. It is shown in Figure A(d) that the task involves the robot using a broom to sweep dust into a dust pan, requiring precise coordination of the sweeping motion to effectively collect the dust. The objects involved include a broom, a dust pan, supporting objects, and the dust itself. Successful completion is defined as all the dust being gathered inside the dust pan. The primary challenge lies in the accuracy and control of the sweeping motion to ensure that the dust is properly directed into the pan.

pick_laptop. As illustrated in Figure A(e), the task requires the robot to pick up a notebook placed on top of a block by first manipulating it into a position suitable for grasping. This involves performing non-prehensile actions, such as pushing or sliding, to adjust the notebook's orientation before securely grasping and lifting it off the block. The objects involved are a notebook and a block. Successful completion is defined as the robot lifting the notebook off the block. Although the task can be performed with a single robotic arm, precise coordination is essential for effective manipulation.

A.2. Real-world Tasks

lift_plate. As shown in Figure A(f), the task involves the robot using both arms to lift a plate, maintaining

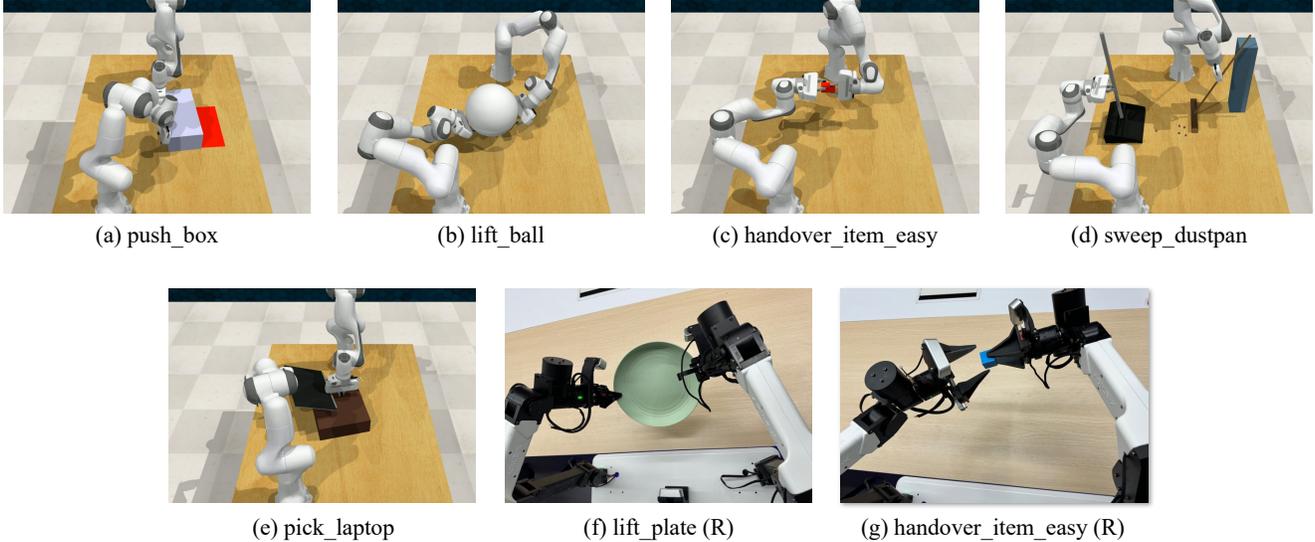


Figure A. The visualization of simulated tasks and real-world tasks. The task with “(R)” means the real-world tasks.

an elevated position for over 3 seconds to meet the success criteria. The target object is a plate requiring coordinated manipulation. Due to the plate’s width and the need for stable control, the operation demands precise bimanual manipulation, requiring synchronized lifting motions from both arms. This task cannot be accomplished by a single robot arm given the plate’s dimensions and stability requirements.

handover_item_easy. Similar to the corresponding simulated task, this task requires the robot to handover a blue item by utilizing one arm to securely grasp and lift the item as shown in Figure A(g). However, the robot do not need to lift the item to a height of 80 cm, but 30cm instead considering the security, while ensuring the other arm remains idle and unengaged, The object involved is a blue block, and the key challenge is as the same as that of in simulation. Successful completion is determined when the item is accurately identified, grasped, and positioned at the required height with no actions performed by the idle arm.

B. Implementation Details

B.1. Details of KStar Diffuser

B.1.1. Backbone

Vision Branch. Following the existing work [17, 49], we employ multiview RGB-D observation as visual input, where the resolution of RGB images is $3 \times 256 \times 256$ and the depth data is processed into point clouds (65536×3) using the camera’s intrinsic and extrinsic parameters. The camera views contain front, overhead, right_over_shoulder, left_over_shoulder, right_wrist, left_wrist. We uniformly use a Vision Transformer which is trained from scratch as the vision encoder.

Language Branch. For the language instruction, we encode the instruction with CLIP’s language encoder¹. Specifically, the input sentence is preprocessed by the CLIP’s tokenizer and then encoded to a sequence of dimensions $\mathbb{R}^{77 \times 512}$. We use its hidden state of “[CLS]” as the textual feature. It is worth noting that we extracted all text features in advance. Therefore, throughout the training phase, the weights of the text encoder do not participate in the gradient calculation of backpropagation.

Fusion Module. The fusion module begins with upsampling the visual features, followed by feature-wise modulation (FiLM) to integrate textual features, projecting them into a high-dimensional semantic space. This hierarchical fusion process is performed iteratively l_{FiLM} times, where we empirically set l_{FiLM} to 3.

The relevant hyperparameters of the backbone are presented in Table B.

Table B. The hyperparameters of Backbone.

	Vision	Text
patch size	16	N/A
hidden size	64	512
# layers	6	12
# heads	8	8
intermediate size	64	2048
dropout ratio	0.1	0.0
activation	lrelu	quick gelu
trainable	✓	×

¹<https://huggingface.co/openai/clip-vit-base-patch32>

B.1.2. Spatial-Temporal Graph

We construct the spatial graph based on the URDF file of robotic arms. In the simulation setup, we use two Franka Panda arms, each with 7 joints, resulting in a total of 14 joints. In the real-world setup, we employ a bimanual ALOHA device, which has 12 joints, with 6 joints per arm. Therefore, the number of nodes and edges in the spatial graph is set to 14 and 12 for the simulation, respectively, and 12 and 10 for the real-world setup.

For the dynamic spatio-temporal graph, we combine the spatial graphs of three consecutive timesteps and add inter-timestep edges connecting the same joint node across different timesteps. In simulated tasks, this results in a graph with 42 nodes and 33 edges. In real-world tasks, the graph has 36 nodes and 28 edges.

We use the Graph Convolutional Graph (GCN) encoder to obtain the spatial-temporal graph representation. The detailed hyperparameters are presented in Table C.

Table C. The hyperparameters of GCN Encoder.

	Simulation	Real-world
# nodes	42	36
# edges	36	28
node dimension	19	19
hidden size	128	128
intermediate size	128	128
# layers	4	4

B.1.3. Optimization Details

During training, we follow the setup of Diffusion Policy, using the DDPM scheduler to forward and denoise, where the step of forward process and reverse process is set to 100 and 1, respectively. Table D shows the training hyperparameters. Table D. The training hyperparameters.

	Values
batch size	64
learning rate	$2e - 4$
warmup step	5 k
weight decay	$1e - 6$
lr scheduler	cosine
training step	150 k
Optimizer	AdamW

B.2. Code Base

The code bases employed for our evaluations are detailed as follows:

- ACT: https://github.com/markusgrotz/peract_bimanual
- RVT-LF: https://github.com/markusgrotz/peract_bimanual

- PerAct-LF: https://github.com/markusgrotz/peract_bimanual
- PerAct2: https://github.com/markusgrotz/peract_bimanual
- DP-J: https://github.com/real-stanford/diffusion_policy
- DP3: <https://github.com/YanjieZe/3D-Diffusion-Policy>

C. Extensive Ablation Studies

To evaluate KStar Diffuser more comprehensively, we conduct following extensive experiments.

(1) The Effects of Demonstration Quantity. Given the critical role of demonstration quantity in imitation learning, we conduct an ablation study by training the KStar diffuser with 50 demonstrations, with results reported in Table E. Additionally, we provide a comparative analysis of policy performance across varying numbers of demonstrations (20, 50, and 100), as shown in Figure 5. The results demonstrate **a clear positive correlation between demonstration quantity and policy performance**. With 20 demonstrations, the policy achieves basic task completion capabilities. Increasing to 50 demonstrations yields significant performance improvements, with success rates rising by around 4.5% across the task suite. The upward trend continues as we scale to 100 demonstrations, indicating that the policy benefits from larger demonstration sets. These findings suggest that expanding the demonstration dataset consistently enhances the policy’s ability to learn and generalize manipulation tasks.

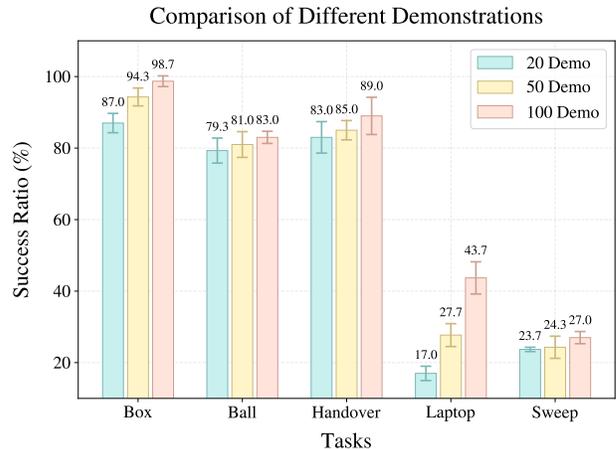


Figure B. The comparison of different number of demonstrations.

(2) The Effects of Action Chunking Size. As mentioned in previous work [9, 65], action chunking prediction serves as an effective approach to address the multimodality challenges inherent in robotic manipulation tasks. To empirically verify this mechanism within our framework, we conducted an evaluation of the KStar Diffuser under various action chunking size configurations. Our experimental setup explores three distinct action chunking strategies: 1) single-step prediction where only the next optimal pose is gen-

Table E. The experimental result on simulated tasks. We train the policy with the setting of different training demonstrations, *i.e.* [20, 100], to test its capability comprehensively. The best results are in **bold**. Each result is reported with three seeds on average.

	Push Box	Lift Ball	Handover Item (easy)	Pick Laptop	Sweep Dustpan	Overall
Demonstration Quantity						
KStar Diffuser (<i>num demos</i> =20)	79.3 \pm 3.5	87.0 \pm 2.7	23.7 \pm 0.6	17.0 \pm 2.0	83.0 \pm 4.4	58.0 \pm 1.4
KStar Diffuser (<i>num demos</i> =50)	81.0 \pm 3.6	94.3 \pm 2.5	24.3 \pm 3.1	27.7 \pm 3.2	85.0 \pm 2.7	62.5 \pm 1.4
KStar Diffuser (<i>num demos</i> =100)	83.0 \pm 1.7	98.7 \pm 1.5	27.0 \pm 1.7	43.7 \pm 4.5	89.0 \pm 5.2	68.2 \pm 2.1
Action Chunking Length						
KStar Diffuser (<i>chunking</i> =1)	92.0 \pm 1.7	98.7 \pm 0.6	23.7 \pm 5.9	16.0 \pm 5.2	12.3 \pm 4.2	48.5 \pm 0.6
KStar Diffuser (<i>chunking</i> =2)	83.0 \pm 1.7	98.7 \pm 1.5	27.0 \pm 1.7	43.7 \pm 4.5	89.0 \pm 5.2	68.2 \pm 2.1
KStar Diffuser (<i>chunking</i> =5)	81.3 \pm 1.5	97.7 \pm 2.5	1.7 \pm 1.2	14.3 \pm 3.2	99.3 \pm 1.2	58.9 \pm 1.4
Historical Observation Length						
KStar Diffuser (<i>history</i> =0)	86.3 \pm 3.5	99.7 \pm 0.6	9.7 \pm 2.1	0.0 \pm 0.0	0.0 \pm 0.0	39.1 \pm 1.1
KStar Diffuser (<i>history</i> =1)	87.7 \pm 4.0	28.3 \pm 2.5	10.3 \pm 7.4	36.3 \pm 8.5	92.7 \pm 1.5	51.1 \pm 1.5
KStar Diffuser (<i>history</i> =2)	83.0 \pm 1.7	98.7 \pm 1.5	27.0 \pm 1.7	43.7 \pm 4.5	89.0 \pm 5.2	68.2 \pm 2.1
Coefficient λ						
KStar Diffuser (λ =0.1)	84.3 \pm 4.9	99.3 \pm 1.2	12.7 \pm 2.3	19.7 \pm 4.7	89.7 \pm 2.5	61.1 \pm 1.1
KStar Diffuser (λ =0.5)	83.7 \pm 11.6	98.3 \pm 1.5	3.3 \pm 1.2	36.3 \pm 6.8	93.7 \pm 3.8	63.1 \pm 4.5
KStar Diffuser (λ =0.9)	83.0 \pm 1.7	98.7 \pm 1.5	27.0 \pm 1.7	43.7 \pm 4.5	89.0 \pm 5.2	68.2 \pm 2.1

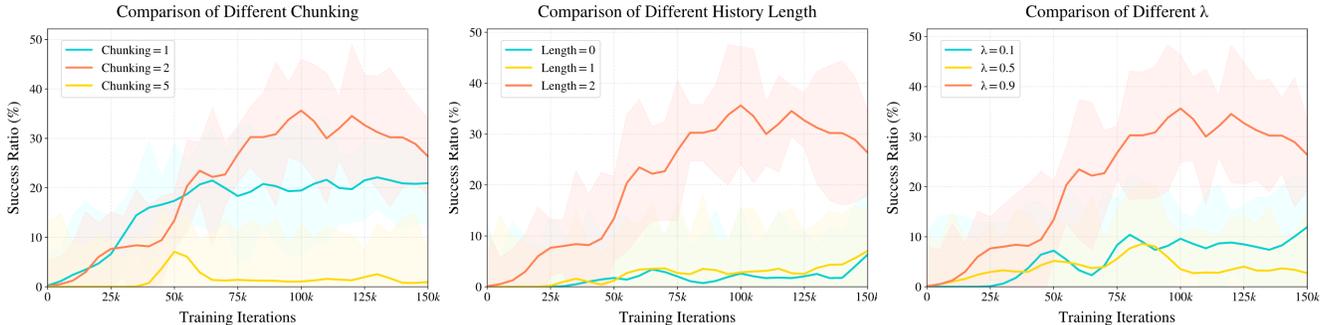


Figure C. *The left*: The result of different action chunking size. *The middle*: The result of history lengths. *The right*: The result of different coefficient λ .

erated (action chunk = 1), 2) two-step prediction of consecutive optimal poses (action chunk = 2), and 3) five-step prediction of sequential optimal poses (action chunk = 5).

The results presented in Table E reveal an **interesting trade-off between prediction horizon and model performance**. While single-step prediction (action chunk = 1) provides basic capabilities, it usually meets the multimodal problem. The two-step prediction strategy (action chunk = 2) emerges as the optimal configuration, demonstrating superior success rates and motion quality across all tasks. Notably, attempting to predict longer sequences (action chunk = 5) leads to decreased performance, with success rates dropping by around 10% compared to the two-step configuration. This performance degradation sug-

gests that when predicting next best pose, extended prediction horizons introduce excessive complexity into the learning problem, making it challenging for the policy to capture and generate accurate action sequences. As shown in Figure C, we present the variation in the success ratio for the `handover_item_easy` task, under different action chunking size configurations as the number of training steps increases.

(3) The Effects of Historical Observation Length. During our experiments, we found that the historical information plays an important role in action prediction quality. We conduct an ablation study across three history lengths: 0 (current observation only), 1, and 2 steps. The results

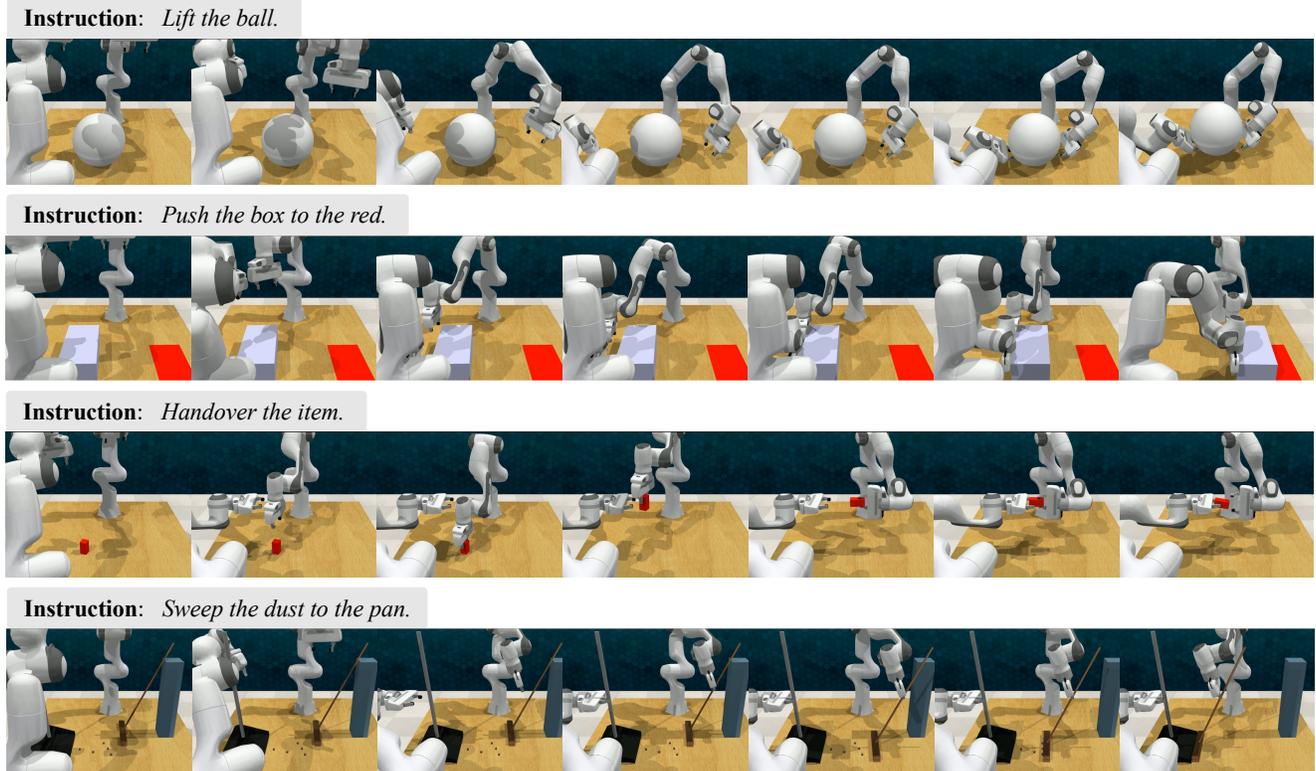


Figure D. The visualization of simulated tasks, including `push_box`, `lift_ball`, `handover_item_easy`, `sweep_dustpan`.

are shown in Table E. With no historical information (0-step), the policy fails to learn effective policies, as it cannot capture the temporal dependencies crucial for manipulation tasks. Adding one historical step enables basic learning capabilities, with the model achieving preliminary success in simpler tasks. Further extending to two historical steps yields optimal performance, showing an approximately 17% improvement over the single-step configuration and demonstrating enhanced stability across all tasks. While longer history lengths might provide more temporal context, they risk introducing redundant information that could potentially obscure relevant features, as observed in our preliminary experiments. These findings suggest that while temporal context is essential for understanding the current state and predicting actions, **an appropriate history window**, e.g., 2 steps, provides an optimal balance between capturing necessary temporal dependencies and maintaining computational efficiency. As shown in Figure C, we present the variation in the success ratio for the `handover_item_easy` task, under different history length configurations as the number of training steps increases.

(4) The Effects of Coefficient λ . In policy learning, a kinematic regularization term is incorporated into the next best pose learning objective to balance task effectiveness

and motion constraints. The regularization strength is controlled by the coefficient λ , which determines the extent of kinematic constraints imposed on the learning process. The choice of λ significantly influences policy performance. Larger values of λ (approaching 1.0) correspond to weaker kinematic constraints, enabling more flexible motion patterns. Conversely, smaller values of λ (e.g., 0.1) impose stricter kinematic constraints, yielding more conservative policies that emphasize motion smoothness over task efficiency. Empirical results, as shown in Table E, demonstrate that policy performance reaches its optimum at $\lambda = 0.9$. For smaller values of λ , the learned policies exhibit overly cautious behavior, manifesting in trajectories that prioritize smoothness at the expense of task efficiency and optimal path planning. Conversely, as λ approaches 1.0, the diminished kinematic constraints result in less regulated motion patterns, potentially compromising trajectory naturalness and precision. These findings indicate that $\lambda = 0.9$ achieves a better trade-off between preserving natural motion characteristics and ensuring efficient task execution. This configuration effectively minimizes the adverse effects of both excessive motion constraints and insufficient regulation, ultimately yielding superior performance across our task suite. As shown in Figure C, we present the variation in the success ratio for the `handover_item_easy` task,

under different coefficient λ configurations as the number of training steps increases.

D. Qualitative Analysis

We show more qualitative result in Figure D. Through the novel approach of encoding robotic arm structural information as graph representations and explicitly incorporating kinematic constraints, our model demonstrates exceptional performance in motion symmetry, synchronization, and coordination across dual-arm manipulation tasks. In the `lift_ball` experiment, the model achieves precise bilateral symmetry in spatial positioning through learned structured representations. The dual arms maintain stable symmetric configurations while preventing object instability through synchronized force application patterns, highlighting the efficacy of our structure-aware control paradigm. Furthermore, in the `push_box` task, the model exhibits remarkable geometric symmetry in motion planning and execution. By leveraging embedded kinematic information, the robotic arms consistently maintain equidistant positioning relative to the target object's center of mass while executing synchronized trajectories along parallel paths. This precise symmetrical control not only ensures operational stability during object manipulation but also establishes a robust framework for dual-arm cooperative control in complex manipulation scenarios.