

# Learning Generalizable Robot Policy with Human Demonstration Video as a Prompt

**Xiang Zhu\***  
Tsinghua University  
Shanghai Qi Zhi Institute  
zhuxiang24@mails.tsinghua.edu.cn  
\*Equal contribution. Project Co-lead.

**Yichen Liu\***  
Tsinghua University  
liu-yc22@mails.tsinghua.edu.cn  
\*Equal contribution. Project Co-lead.

**Hezhong Li**  
Tsinghua University  
lhz21@mails.tsinghua.edu.cn

**Jianyu Chen<sup>†</sup>**  
Tsinghua University  
Shanghai Qi Zhi Institute  
RobotEra TECHNOLOGY CO., LTD.  
jianyuchen@mail.tsinghua.edu.cn  
<sup>†</sup> Corresponding Author.

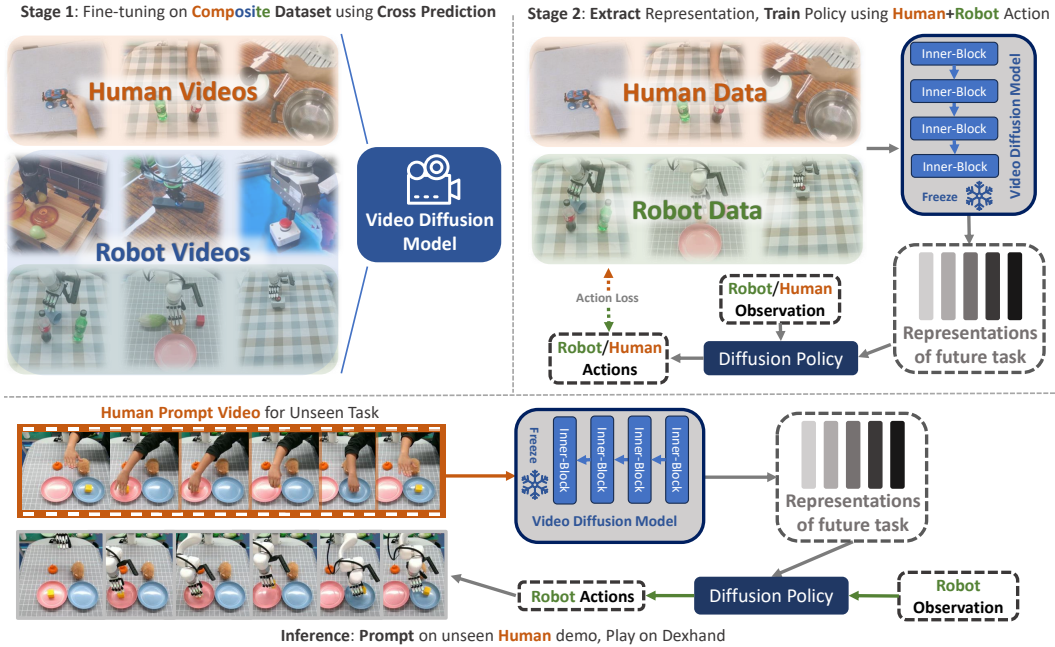


Figure 1: First, fine-tune a video diffusion model on diverse datasets to obtain informative representations. In the second stage, use a video generation model to extract information from human prompt videos for skill learning with both generation and robot data. Finally, during inference, the model prompts on unseen human demos to perform tasks based on the human input.

## Abstract:

Recent robot learning methods commonly rely on imitation learning from massive robotic dataset collected with teleoperation. When facing a new task, such methods generally require collecting a set of new teleoperation data and finetuning the policy. Furthermore, the teleoperation data collection pipeline is also tedious and expensive. Instead, human is able to efficiently learn new tasks by just watching others do. In this paper, we introduce a novel two-stage framework that utilizes human demonstrations to learn a generalizable robot policy. Such policy can directly take human demonstration video as a prompt and perform new tasks without any new teleoperation data and model finetuning at all. In the first stage, we train video generation model that captures a joint representation for both the human and robot demonstration video data using cross-prediction. In the second stage, we fuse the learned representation with a shared action space between human and robot using

a novel prototypical contrastive loss. Empirical evaluations on real-world dexterous manipulation tasks show the effectiveness and generalization capabilities of our proposed method.

**Keywords:** Learn from Human, Human Video Prompt

## 1 Introduction

Traditional robot learning methods typically train language-conditioned policies on large datasets collected through teleoperation[1, 2]. While effective for known tasks, this paradigm faces two fundamental limitations when encountering novel tasks. First, language instructions, though intuitive for humans, provide only categorical information and lack the rich spatial and temporal details crucial for physical manipulation. Second, adapting to new tasks usually requires collecting additional robot demonstrations, a process that is both time-consuming and expensive due to the complexity of teleoperation systems.

Humans, in contrast, can efficiently acquire new skills simply by observing others perform tasks. This observation suggests that visual demonstrations may offer a more natural and information-rich medium for teaching robots. Videos inherently capture not just what task to perform, but also how to perform it, including critical aspects like object relationships, motion trajectories, and timing. Moreover, human demonstration videos are significantly more scalable to acquire compared to robot data, whether captured in laboratory settings or obtained from existing online resources.

Current approaches have only begun to explore this direction. While works like EgoMimic[3] have incorporated human demonstrations, it focuses on single-task scenarios. The field still lacks a general framework that can leverage the full potential of human videos for robot learning.

We propose a two-stage human-prompted learning framework that combines robotic datasets with human demonstration data to address challenges in task learning. In the first stage, we use a video generation model that receives a prompt video of a human performing a task and an image of the robotic hand. The model generates a video of the robot performing the task, embedding embodiment transfer information through a cross-prediction strategy. This helps the model learn a representation that captures the task, context, and target modality effectively. In the second stage, we fine-tune the representation using a diffusion policy, incorporating both human and robot data. A unified action space bridges the gap between the two modalities, while a cluster-based loss enhances skill separation and multi-skill imitation performance. Experiments on real-world tasks demonstrate the framework’s effectiveness in improving human-robot interaction and versatile manipulation.

## 2 Related Works

### 2.1 Learn from Human Videos

Recent advances in robot manipulation have increasingly utilized human video data to enhance both dexterous and gripper-based manipulation. In dexterous manipulation, works like [4, 5, 6] focus on fine-grained control of multi-fingered systems, while [7] integrates affordance cues. For gripper manipulation, end-to-end video-conditioned policies such as [8, 9, 10, 11] translate visual cues into actionable policies. Approaches using paired human–robot demonstration data, including [12], [13, 14], [15], [16], [3], and [17], address the domain gap by linking human actions with robot trajectories. More recently, generative video techniques like [18], [19], [20], and [21] leverage video synthesis and textual cues to generate visuomotor policies. These studies highlight the growing trend of using human video demonstrations, paired data, and generative methods to create more adaptable and robust robot manipulation policies.

### 2.2 Video as Prompt for Robotic Learning

Recent works [22, 23] have increasingly used human demonstration videos to guide robot learning. For example, [14] addresses human-robot embodiment mismatch by converting human videos into robot-centric demonstrations using unsupervised domain adaptation and keypoint extraction. [24] enables zero-shot generalization by conditioning robot policies on pretrained video embeddings. Similarly, [8] maps human videos to robot actions through cross-attention transformers, while [25] improves sample efficiency and generalization through contrastive learning, imitation, and limited

adaptation. [11] focuses on cross-embodiment skill discovery for transferable representations. In our work, we extend these insights by using video as a prompt and integrating video generation model to more effectively bridge the human-robot gap by embodiment transfer and to achieve robust, generalized policy learning for complex manipulation tasks.

### 2.3 Diffusion Models for Robot Policy Learning

Diffusion models have shown great success in generative computer vision [26, 27], leading to their adaptation in robotic policy learning. Pioneering works [28, 29, 30] demonstrated their ability to generate denoised robot actions and capture multi-modal behavior distributions. Scaling efforts, such as [31], showcased generalization across diverse robotic platforms with transformer-based diffusion policies pretrained on the Open X-Embodiment dataset. Models like MDT [32] and RDT-1B [33] use transformer-based diffusion models, replacing the traditional U-Net. RDT-1B further unifies action representations across robots and incorporates multi-robot data for bimanual manipulation.

Diffusion-based policies can model multimodal action distributions in high-dimensional spaces [28, 34, 35]. [36] augments these models with unsupervised clustering and intrinsic rewards to maintain multiple behavior modes, while [37] adds an entropy regularizer to enhance robustness. Building on these approaches, we propose injecting SPCL’s [38, 39] prototypical contrastive loss into diffusion policy training, leveraging task category labels as supervised prototypes. Thereby sharpening policy representations along known task dimensions and enhancing multi-task imitation performance. Our method can be framed as “in-context learning”, where human demonstrations serve as the contextual basis for the policy to achieve zero-shot transfer of unseen robotic skills.

## 3 METHODS

In this section, we detail our framework by outlining its two-stage training process. Our goal is to enable the agent to learn and extract meaningful features and representations from human demonstration videos to perform specific tasks. The algorithm is divided into two stages. In the first stage, we develop a robust representation from a large collection of human data, robotic gripper data, and a small amount of dexterous hand data. This representation captures the target modality, the intended task, and the scene context. In the second phase, using this solid representation, we learn human manipulation patterns from abundant human data and a limited set of high-quality dexterous hand operation data, transferring these patterns to the dexterous hand to achieve better generalization with only a small amount of valuable dexterous hand data.

### 3.1 Augment Video Generation by Cross Prediction(VGCP)

Recent advancements in video generation models have leveraged extensive online video datasets containing prior knowledge about physical world dynamics. However, these models lack data relevant to robotic manipulation. To address this, we fine-tune an existing model with a custom dataset focused on robotic and object manipulation. Our goal is to enable agents to perform tasks based

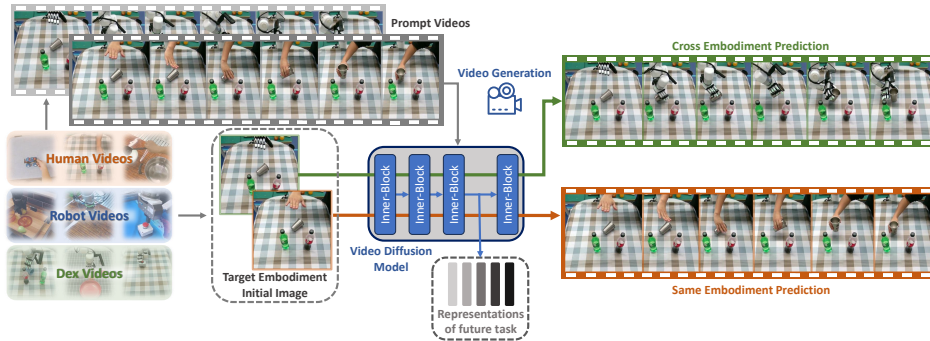


Figure 2: **Stage 1:** We fine-tune a pre-trained video generation model using cross-prediction, enabling the model to retain physical knowledge while gaining the ability for embodiment modality transfer.

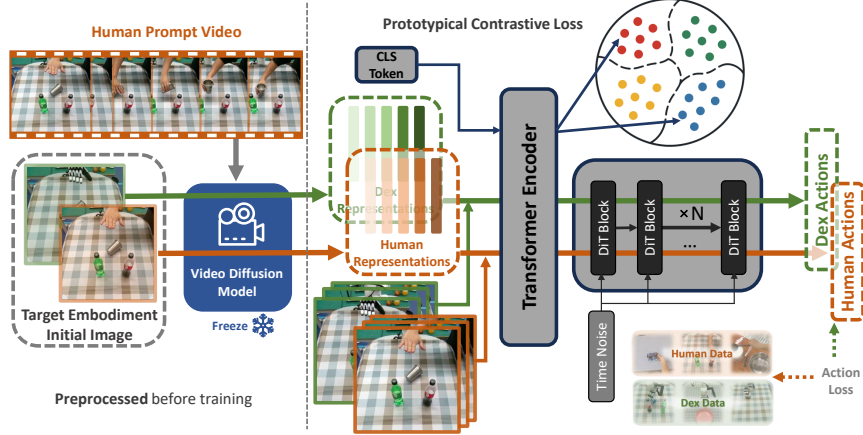


Figure 3: **Stage 2:** We use the diffusion model trained in the first stage to combine information from human prompt videos with the target embodiment, generating an informative representation. This representation, along with a shared action between human and robot and a prototypical contrastive loss, enables the diffusion policy to learn common task and skill information for both.

on human prompt videos, particularly complex dexterous hand operations. Therefore, the video generation model is trained to include details on human manipulation, object motion, scene context, and affordances. Our dataset consists mainly of human manipulation videos, supplemented by robotic gripper and self-collected dexterous hand videos. While gripper videos differ from dexterous hand operations, they provide valuable information on object movement and scene understanding, enhancing the model’s performance.

**Cross-Prediction:** To further leverage our data and enhance the quality of the learned representation, we propose a method called cross-prediction. As shown in Fig.2, our video generation model receives a video prompt showing a task performed by a **source embodiment** and an initial scene for a **target embodiment**, then generates a video of the target embodiment performing the task, effectively transferring the embodiment. For example, given a prompt video of a human grasping a cup, the model produces a video of a dexterous hand performing the same action. During training, we choose cross-prediction with probability  $P$  (using different **source** and **target** embodiments) or normal-prediction with probability  $1-P$  (using the same embodiment). This approach embeds embodiment transfer information into the video generation model, further improving the transferability of the learned representation. Our cross-prediction method randomly selects different source ( $s$ ) and target ( $t$ ) embodiments with probability  $P$ , and identical ones with probability  $1 - P$ . When the two embodiments are the same, the process mirrors typical video generation models that generate subsequent frames from an initial frame. Our goal is for the model to learn the modality transfer between human and robot, while preserving existing knowledge of their manipulation. We believe that this approach enables the model to capture a video prompt representation that encompasses the skills used by the source embodiment, the manipulated objects, and some environmental context. The fine-tuned model will then be frozen during the second stage of training.

### 3.2 Skill Learning by Human Video-Action Pair Boosting

Training manipulation policies with dexterous hand data has become a key approach in robotic manipulation. However, collecting such data through teleoperation is time-consuming and costly. To address this, we reduce reliance on robotic hand data by leveraging human hand demonstrations to enhance manipulation capabilities. The availability of human hand data is virtually unlimited, with videos easily sourced from the internet or self-collected, requiring minimal time and infrastructure. We propose incorporating human hand demonstrations into a compatible format alongside teleoperated robotic hand data, which are then jointly trained using Imitation Learning (IL). This forms the core methodology of our Stage 2 algorithm.

**Human Data Preprocessing:** The objective of human data processing is to establish correspondence between human demonstration data and robot motion data. Given that we only have third-



person view RGB videos of human demonstrations, while the robot data contains both third-person view RGB videos and joint state information. The robot’s end-effector state comprises two components: 6D Wrist pose (wrist position  $\in \mathbb{R}^3$  and wrist orientation  $\in \text{SO}(3)$ ) and finger joint  $\theta^n$ . We employ the hand-tracking method WiLoR[40] for hand localization in video frames and 3D hand mesh reconstruction, and the model outputs 21 keypoints  $J_{h_i} \in \mathbb{R}^3$  for each frame. We are able to obtain wrist 6D from these 21 keypoint, with detailed methodology provide in Appendix.

And for the wrist position, we can compute the 3D position of the human hand in the camera coordinate system using the camera intrinsic parameters. To align this with the robot’s base coordinate system, we transform all coordinate systems to the camera coordinate. To establish joint-level correspondence between human and robotic systems, hand motion retargeting is required to map the kinematic configurations from the human demonstrator’s hand to the target robot manipulator. We follow Anyteleop[41] to formulate the retargeting problem as an optimization.

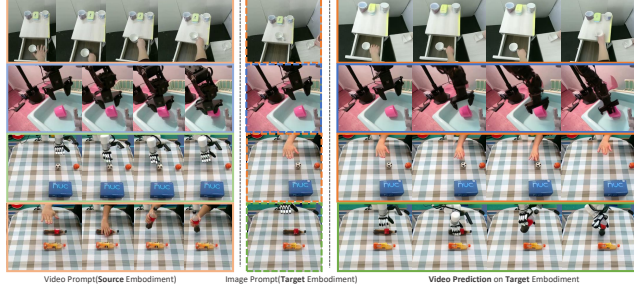


Figure 4: Examples of Cross-Prediction Video Generation.

**Representation Conditioned Diffusion Policy:** Our framework employs a diffusion policy that models the conditional action distribution  $p(a_t|s_t, z)$ , where:  $s_t = f_{resnet34}(o_t) \in \mathbb{R}^{1000}$  is the visual observation feature extracted by a pretrained ResNet-34,  $z \in \mathbb{R}^{4096}$  is the stage-one representation,  $a_t \in \mathbb{R}^{19}$  is the action vector at time  $t$ . The diffusion process operates in the action space through:

$$a_t^i = \sqrt{\alpha_i} a_0 + \sqrt{1 - \alpha_i} \epsilon, \epsilon \sim \mathcal{N}(0, \mathbf{I}) \quad (1)$$

while the reverse process learns to predict and remove the noise through a transformer-based denoising network  $\epsilon_\theta$ .

For dexterous hand control tasks, we decompose the action space into three distinct components: finger joint angles  $a_t^{finger} \in \mathbb{R}^{12}$ , wrist orientation  $a_t^{rot} \in \text{SO}(3)$ , and wrist position  $a_t^{pos} \in \mathbb{R}^3$ . This decomposition enables specialized handling of each action modality through separate prediction heads in the denoising network. The training loss combines weighted component losses:

$$\mathcal{L}_{action} = \lambda_f \mathcal{L}_{finger} + \lambda_r \mathcal{L}_{rotation} + \lambda_p \mathcal{L}_{postion} \quad (2)$$

where:

$$\mathcal{L}_* = \mathbb{E}_{i, a_0, \epsilon} [\|\epsilon_* - \epsilon_\theta(a_t^i, i, h_0)\|^2] \quad (3)$$

$h_0 = [s_t, z] W_e + b_e$  is the concatenated input token,  $W_e$  is the embedding matrix,  $b_e$  is the bias.  $\lambda_f, \lambda_r, \lambda_p$  are weighting coefficients determined through validation.

### 3.3 ProtoDiffusion Contrastive Policy(PDCP)

Siamese Prototypical Contrastive Learning (SPCL)[38, 39] first groups feature embeddings into prototypes via K-Means and then applies a Siamese-style metric loss to pull together embeddings within each prototype while pushing apart those from different prototypes, alongside a prototypical cross-entropy loss that treats prototype assignments as soft labels to sharpen each sample’s affinity to its cluster. By leveraging learned cluster structure rather than individual instances, SPCL mitigates false negatives and yields more stable positive sets, thereby enhancing semantic discriminability in self-supervised learning.

To adapt SPCL for diffusion policy training, we add a learnable clustering token to the input of our Diffusion Transformer (DiT) encoder and use the encoder’s output for that token—denoted  $h$ —as a prototype-aware latent. During training we jointly optimize three losses:

(1) the NT-Xent contrastive loss  $\mathcal{L}_{contra}$ , which treats tasks sharing the same skill as positives and tasks of different skills as negatives (analogous to data-augmented positives in vision); (2) a prototypical cross-entropy loss  $\mathcal{L}_{proto}$ , which encourages each the learned latent  $h$  to align with the

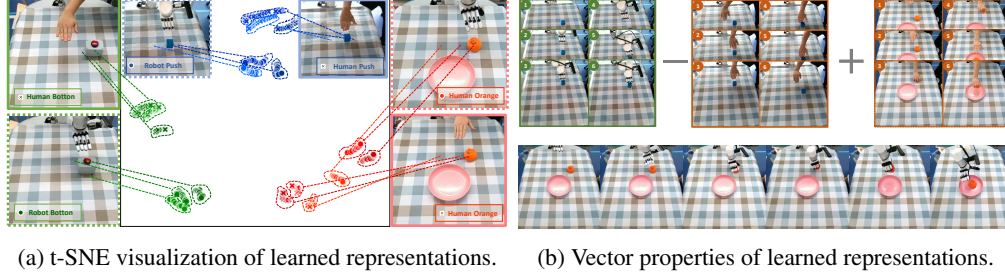


Figure 5: Evaluation for the Learned Representations

cross-entropy distribution over its K-Means prototype label, promoting tighter intra-prototype clustering and clearer inter-prototype separation. And (3) a Siamese-style metric loss  $\mathcal{L}_{\text{metric}}$  at the prototype level, which further enforces proximity among same-prototype samples in metric space while repelling different-prototype samples, thus reducing semantic confusion. Please refer to the appendix for detailed information about these three losses. The total loss for PDCP consist of the action loss and a weighted sum of the mentioned loss with coefficient  $w_p, w_p, w_m$ :

$$\mathcal{L} = \mathcal{L}_{\text{action}} + w_c \mathcal{L}_{\text{contra}} + w_p \mathcal{L}_{\text{proto}} + w_m \mathcal{L}_{\text{metric}} \quad (4)$$

Together, these objectives guide the network to learn task-discriminative, skill-aware, and modality-agnostic representations which improving cross-modal alignment between human and robot demonstrations and ultimately enhancing the robustness and generalization of our diffusion policy.

## 4 EXPERIMENTS

In the next section, we conduct extensive experiments to validate our proposed method. We deploy our algorithm on the dexterous hand Xhand, performing engaging tasks such as pick-and-place, pressing, pouring water, flipping cups, push-pull operations, and even ball flicking. Our experiments aim to answer the following questions:

1. Can cross-prediction enhance the performance, effectiveness, and representational capacity of the learned representation?
2. Is it feasible to use different video prompts in combination with a video generation model for skill extraction, and what is the resulting impact?
3. Does incorporating human action data improve the generalization ability of the diffusion policy?
4. Can PDCP boosting and outperforming the norm diffusion policy?
5. Can the final learned policy demonstrate “in-context learning” capabilities?

### 4.1 Evaluation for Cross Prediction Video Generation

Our cross-prediction mechanism extracts the task performed by the source embodiment from the prompt video and, using the scene context and target embodiment information form another given image, generates a video of the target embodiment executing that task. As shown in Fig.4, we show two examples of same-embodiment prediction alongside with two examples of cross-embodiment prediction. The generated videos faithfully capture the task details from the prompt, and in the cross-embodiment cases, the generated target embodiment motion and scene context align precisely with the original task. This demonstrates that our model has learned not only the physical dynamics of the world but also learned how to transfer task and skill across different embodiments.

### 4.2 Evaluation for the Learned Representations

As shown in previous experiments, applying cross-prediction enables the video generation model to convert between human and robot embodiments while retaining its original capabilities. However, our goal is to extract the task information performed by the source embodiment from the prompt video using the cross-prediction fine-tuned SVD, and then transfer that information to the target embodiment. So, how do we evaluate and verify the performance of the representations extracted by the video generation model? We conducted the following experiments:

Table 1: Success Rate (SR) and Score metrics on position, scene, and background generalization

Method	Position		Scene		Background	
	SR	Score	SR	Score	SR	Score
Language+R.+H.	0.58	69.5	0.56	55.6	0.36	50.9
Repr.+R.	0.68	77.9	0.67	81.1	0.55	69.1
Repr.+R.+H.	0.74	81.6	0.67	84.4	0.64	74.5
<b>Repr.+R.+H.+DPCP (Ours)</b>	<b>0.79</b>	<b>85.3</b>	<b>0.69</b>	<b>88.9</b>	<b>0.73</b>	<b>80.9</b>

#### 4.2.1 t-SNE Visualization of the Learned Representation

We conducted a t-SNE visualization experiment. We applied t-SNE to reduce the dimensionality of the skills extracted from the **human representation** (with a human video as the prompt and a human image as the initial frame) and the **robot representation** (with a human video as the prompt and a robot image as the initial frame), and then visualized them on a two-dimensional plane. The results, shown in Fig. 5a, reveal that tasks of the same category are grouped by similar color schemes, with identical hues indicating representations at different time points. Dashed lines highlight the positions of target objects in the two-dimensional space. The plot demonstrates that representations can effectively distinguish tasks, with the same task’s representations clustering closely based on embodiment, while still distinguishing target object positions. Additionally, representations from different time steps within the same video are tightly grouped yet distinguishable, reflecting the capture of continuous temporal information. These results show that the learned representations not only encode embodiment and task information but also capture object positioning and motion details, which is challenging to obtain with language prompts, showcasing the advantage of video prompts.

#### 4.2.2 Vector Properties of the Learned Representations

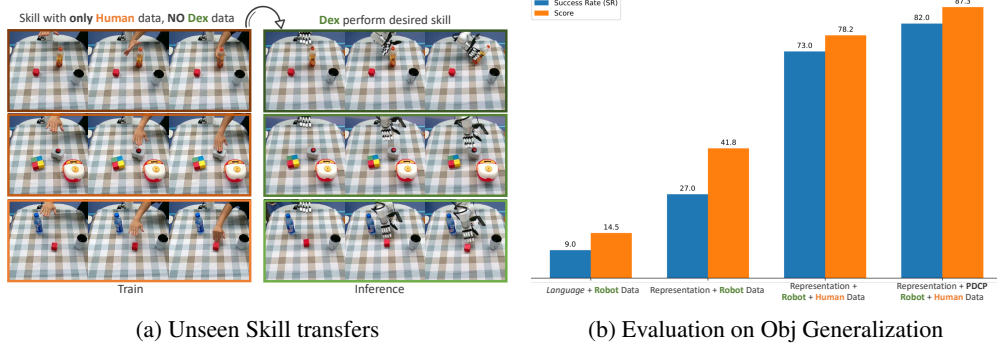
In word embeddings, a fascinating phenomenon is observed: these embeddings exhibit an intrinsic vector structure. For example, taking the representation of “king,” subtracting that of “queen,” and then adding the representation of “nurse” yields that of “doctor.” This arithmetic property is celebrated as both elegant and intriguing. Thus, we aim to explore whether our learned representations possess a similar vectorial nature. We subtracted the human representation from the corresponding robot representation for the “push blue box” task and then added the human representation for the “grasp orange” task. Using this composite representation, we generated a video that, as shown in the Fig. 5b, ultimately depicted a robot executing the “grasp orange” task. Moreover, we conducted tests on other tasks—such as “push button” and “pour water”, and observed similarly intriguing results. This demonstrates that our learned representations indeed exhibit notable vector properties.

### 4.3 Real-World Robot Experiments

Our policy relies solely on third-person RGB images from a fixed camera, and we evaluate our method on a multi-task benchmark (See Appendix). We conduct comparisons with three baselines: (1) **Representation+Robot** means Using only robotic data without human data. (2) **Language+Robot+Human** Means replacing skill representations with CLIP-encoded task labels (e.g., “pour water”) as policy inputs. (3) **Representation+Robot+Human** means using both robot data and human data

#### 4.3.1 Position, Scene, Background Generalization

To assess the generalization capabilities of our policy, we conducted comprehensive testing across three critical dimensions. **Positional Generalization** means randomly relocating target objects within the workspace. **Scene Generalization** means randomly changing or replacing inoperable objects within the workspace. **Background Generalization** means testing under different background (tablecloth, ...) To evaluate task execution performance, we introduce two metrics: Success



Rate (SR) and Task Score. The binary SR metric (1 for complete success, 0 otherwise) and a subtask metric (details provided in Appendix).

As evidenced by the comparative results in Table 1, our method demonstrates better performance across all three generalization scenarios when evaluated against the three baseline approaches. The comparative evaluation between our Stage-1 representation-guided policy and language-conditioned alternatives reveals fundamental advantages in the learned representations’ capacity to encode richer task-relevant information.

Our analysis further reveals that incorporating human demonstration data consistently enhances policy performance across all three generalization scenarios. The performance improvements observed suggest that high-quality human demonstrations provide critical behavioral patterns and manipulation strategies that are often underrepresented in robot-collected data. Table 1 clearly demonstrates the effectiveness of the proposed PDCP algorithm in mitigating mis-grasping issues when handling closely spaced objects. We posit that while the Stage-1 representation primarily focuses on spatial variations, the PDCP module provides complementary functionality by enabling task-aware feature clustering.

#### 4.3.2 Object Generalization

To investigate the complementary benefits of our representation learning framework and human demonstration data, we design object-level generalization experiments to evaluate their individual and combined effects. While our representations perform well across all generalization scenarios, these experiments isolate and quantify the contribution of human data with novel object instances.

As shown in Figure 6b, we test four distinct tasks involving unseen objects. For each task, we compare two training conditions: (1) using robot-collected data with similar, non-identical objects (excluding target objects), and (2) combining robot data with a larger set of human demonstrations that include the target objects. The results clearly show that human data significantly improves manipulation success rates for novel objects, highlighting its ability to transfer knowledge on object affordances and manipulation strategies.

Additionally, our representation outperforms direct language input, thanks to its incorporation of end-effector positional information. We argue that for tasks absent from the robot training set, our representation provides essential spatial cues and partial operational knowledge, while the human data offers complementary information on novel objects and manipulation strategies.

#### 4.4 Skill Generalization

Given that our policy can handle tasks unseen in the robot training dataset, we investigate whether it can also demonstrate novel skills absent from the original robotic data. To test this hypothesis, we trained three distinct policies on three novel tasks: grasping an orange juice bottle, pressing a button, and pushing a red block. Each policy was trained using approximately 50 robotic trajectories of diverse grasping tasks combined with extensive human demonstration data. For comparison, we implemented baseline models trained solely on either (1) the same 50 robotic trajectories or (2) language-augmented robot and human data. Crucially, these baseline models failed to generalize to new skills at the skill level, whereas our approach combining representation learning with joint human-robot data demonstrated clear potential for novel skill generalization. As illustrated in Fig-



ure 6a, the left panel displays human hand manipulation data, while the right panel shows the robot executing corresponding actions that were never present in its original training dataset. This comparative analysis reveals that incorporating human data significantly enhances the policy’s ability to acquire and transfer novel skills to the robotic domain.

## 5 CONCLUSIONS

This work fully leverages human data throughout both the representation generation and diffusion policy processing stages. By utilizing human videos as prompts—which provide substantially richer information than language inputs—our approach demonstrates superior performance in dexterous hand manipulation tasks.

## 6 LIMITATIONS

Our work faces inherent challenges in generalizing new skills from human videos. While we demonstrate successful cases where the robot replicates human manipulations to complete tasks (occurring when the policy successfully interprets and follows human data), many attempts fail due to misalignment between the policy’s execution and the observational prerequisites for skill activation. Currently, our new skill acquisition success rate ranges under 10%.

Although this remains a significant challenge, the human-demonstrated novel skills we introduced show measurable potential for training corresponding robotic capabilities. Importantly, we hypothesize that substantially increasing human demonstration data for specific tasks - similar to the scaling approach used with robot data - could help mitigate this issue.

## 7 APPENDIX

### A Appendix

#### A.1 Details on SVD formulation

More specifically, we selected the open-source stable video diffusion (SVD) model with 1.5 billion parameters as the base model. The original SVD model conditions on the initial-frame image  $I_0$ , we augment it by adding a CLIP image encoder  $\phi(\cdot)$  to extract features from the prompt video and integrate them into the SVD UNet block via cross-attention. To enable longer video outputs, reduce training costs, and improve efficiency, we set the model to generate videos with  $T = 16$  frames at a resolution of 256×256. We denote our model as  $G_\theta$ , with inputs comprising an initial image  $I_0^t$  of the target embodiment and a video  $V^s = I_{0:T}^t$  of a source embodiment performing the specified task. We fine-tune the base SVD model with a diffusion objective which generates video from noised samples  $x_t = \sqrt{\bar{a}_t}x_0 + \sqrt{1 - \bar{a}_t}\epsilon$ , where  $x_0 = V^t$  and  $D$  is our dataset.

$$\mathcal{L}_{CP} = \mathbb{E}_{x_0 \sim D, t, \epsilon} \|G_\theta(x_t, \phi(V^s), I_0^t) - x_0\|^2 \quad (5)$$

**Dataset for Training Video Generation by Cross Prediction:** The dataset for video generation model training comprises three parts: (1) videos of human object manipulation, (2) videos of robotic gripper manipulation, and (3) videos of dexterous hand manipulation that we collected ourselves. The human manipulation videos are sourced from HOI4D, RH20T, and our own collected footage, with the aim of teaching the model how humans handle objects. The robotic gripper videos, primarily from BRIGE and RH20T, enable the model to learn the physical dynamics of object motion and acquire the ability to transfer from human to gripper embodiments. Finally, a small set of high-quality dexterous hand videos is included to capture authentic dexterous manipulation and to facilitate the transfer between human and dexterous hand operations. The video generation model training process takes 5 days on 8 NVIDIA A100 GPUs. The second stage of policy learning requires approximately 6-12 hours on four NVIDIA A100 GPUs.

## A.2 Loss Details for ProtoDiffusion Contrastive Policy(PDCP)

(1) **The NT-Xent contrastive loss**, which treats tasks sharing the same skill as positives and tasks of different skills as negatives (analogous to data-augmented positives in vision);

$$\mathcal{L}_{\text{contra}} = -\frac{1}{|B|} \sum_{i \in B} \frac{1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(\text{sim}(z_i, z_p)/\tau)}{\sum_{a \in B \setminus \{i\}} \exp(\text{sim}(z_i, z_a)/\tau)} \quad (6)$$

where  $P(i) = \{j \in B | y_j = y_i, j \neq i\}$  indicate samples which have the same skill with  $i$ .  $z_i = g_c(h_i)$  where  $g_c(\cdot)$  is a projection head, and  $\text{sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$  is the similarity metric.

(2) **A prototypical cross-entropy loss**, which encourages each the learned latent  $\mathbf{h}$  to align with the cross-entropy distribution over its K-Means prototype label, promoting tighter intra-prototype clustering and clearer inter-prototype separation.

$$\mathcal{L}_{\text{proto}} = -\frac{1}{|B|} \sum_{i \in B} \log \frac{\exp(\text{sim}(h_i, \mu_{c(i)})/\tau)}{\sum_{j=1}^K \exp(\text{sim}(h_i, \mu_j)/\tau)} \quad (7)$$

where  $\mu_j$  is the  $j$ -th prototypes(cluster center);  $c_i$  is the cluster index for sample  $i$  and  $\tau$  is the temperature parameter.

(3) **A Siamese-style metric loss** at the prototype level, which further enforces proximity among same-prototype samples in metric space while repelling different-prototype samples, thus reducing semantic confusion.

$$\mathcal{L}_{\text{metric}} = \frac{1}{|\mathcal{P}|} \sum_{(i,j) \in \mathcal{P}} \text{CE}(g_m(\|h_i - h_j\|), 1) + \frac{1}{|\mathcal{N}|} \sum_{(i,j) \in \mathcal{N}} \text{CE}(g_m(\|h_i - h_j\|), 0) \quad (8)$$

where  $\mathcal{P}$  indicates all the sample pairs in the cluster, while  $\mathcal{N}$  indicates sample pairs in different cluster and  $g_m(\cdot)$  is also a projection head.

## A.3 Details of Human Data Preprocessing

To extract the 6D wrist pose (position + orientation) from 21 keypoints, we first obtain the 3D position directly from the wrist keypoint  $\mathbf{W} \in \mathbb{R}^3$ . For the 3D orientation, we construct a local coordinate frame using the wrist keypoint  $\mathbf{W}$ , the first metacarpophalangeal joint of the index finger  $\mathbf{I}$ , and the first metacarpophalangeal joint of the middle finger  $\mathbf{M}$ . The x-axis unit vector  $\mathbf{x}$  is defined as the normalized direction from  $\mathbf{W}$  to  $\mathbf{M}$ :

$$\mathbf{x} = \frac{\mathbf{M} - \mathbf{W}}{\|\mathbf{M} - \mathbf{W}\|} \quad (9)$$

The z-axis unit vector  $\mathbf{z}$  is computed as the normalized normal vector of the plane formed by  $\mathbf{W}$ ,  $\mathbf{I}$  and  $\mathbf{M}$ , oriented toward the palm:

$$\mathbf{z} = \frac{(\mathbf{I} - \mathbf{W}) \times (\mathbf{M} - \mathbf{W})}{\|(\mathbf{I} - \mathbf{W}) \times (\mathbf{M} - \mathbf{W})\|} \quad (10)$$

Finally, the y-axis  $\mathbf{y}$  is derived via the right-hand rule as the cross product of  $\mathbf{z}$  and  $\mathbf{x}$ :

$$\mathbf{y} = \mathbf{z} \times \mathbf{x} \quad (11)$$

The resulting rotation matrix  $\mathbf{R} = [\mathbf{x} \ \mathbf{y} \ \mathbf{z}]$  and wrist position  $\mathbf{W}$  form the complete 6D pose representation. This method ensures orthonormality and a right-handed coordinate system.

## A.4 Integrated Normalization Methodology for Robot and Human Datasets in Stage 2

In Stage 2, we normalize the actions of **both the human hand and the robot dexterous hand** to the range  $[0, 1]$  using a joint normalization approach. Specifically, for the 19-dimensional action space, we first compute the **global maximum** ( $\max_a$ ) and **global minimum** ( $\min_a$ ) for each dimension across the combined dataset of human and robotic actions. Each action component  $a_i$  (where  $i \in \{1, 2, \dots, 19\}$ ) is then normalized as follows:

$$a_{\text{norm},i} = \frac{a_i - \min_a}{\max_a - \min_a} \quad (12)$$

Table 2: Dataset Composition

Datasets	Human	Gripper	Dexhand	Number Trajectories	Stage 1	Stage 2
<b>RH20T</b>	✓	✓	×	87529	✓	×
<b>HOI4D</b>	✓	×	×	2972	✓	✓
<b>Bridge</b>	×	✓	×	25460	✓	×
<b>Self-Collected</b>	✓	×	✓	3679	✓	✓

## A.5 Real-World Robot Experiments Setting

### A.5.1 Hardware Settings

In the physical experiment phase, we employed an xArm7 robotic manipulator coupled with a dexterous robotic hand as our hardware platform. The robotic system was controlled using two distinct control signals: (1) 6D wrist pose data for arm movement control, and (2) 12D finger articulation data for precise hand manipulation. Notably, our observation inputs were exclusively derived from third-person RGB camera images.

### A.5.2 Experiment Settings and Details

We integrate the frameworks of Stage 1 and Stage 2 and evaluate our algorithm across various skills. The detailed dataset statistics are presented in Table 2. For each skill task, we collected paired robot and human demonstration videos (with one-to-one correspondence between robot and human videos) as Self-Collected Datasets. The table reports: (1) whether the dataset incorporate the Human/ Gripper/ Dexhand data, (2) number trajectories of dataset, (3) whether Stage1 and Stage 2 trained using the dataset.

To evaluate the performance of our tasks, we have defined two metrics: **Success Rate (SR)** and **Score**. The Success Rate (SR) is a binary indicator that signifies the completion status of a task, with a value of 1 representing successful task completion and 0 indicating failure. However, this metric alone is insufficient to fully capture the nuances of task performance. Therefore, we introduce the Score metric to assess the completion status of multi-stage tasks. For the Pick-and-Place task, approaching the object and initiating a picking action yields 30 points, successfully picking the object awards 40 points, and placing the object at the designated location results in an additional 30 points. In the Pour water task, approaching the bottle and initiating a grasping action earns 30 points, grasping the bottle gains 40 points, and accurately pouring water into the target cup adds another 30 points. For the Flip Cups task, approaching the cup and initiating a grasping action scores 30 points, performing the flipping action awards 30 points, and successfully standing the cup upright on the table grants 40 points. In the Push Box task, approaching the block and initiating a pushing action yields 60 points, while successfully moving the block forward adds 40 points. For the Press Buttons task, approaching the object and initiating a pressing action earns 30 points, and completing the task awards 70 points.

**Experiment Settings of Position, Scene and Background Generalization** In the generalization experiments on position, scene, and background in Section 4.3.1, we utilized the entire dataset and tested all skills to comprehensively demonstrate the performance improvement of our policy following the introduction of human datasets and representation learning.

**Experiment Settings of Object Generalization** In the novel object generalization experiments in Section 4.3.2, the following four tasks were showcased: picking up a capybara and placing it on a plate, picking up dice and placing it on a plate, pouring orange juice, and pushing a red box. For training the Pick-and-Place task, our robot dataset comprised approximately 300 grasping tasks (excluding picking capybara and dice), while the Human dataset included around 700 tasks (including picking capybara and dice). The Pour Orange Juice task was trained using roughly 150 Robot data entries (excluding pouring orange juice) and about 300 Human data entries (including pouring orange juice). For training the Push Red Block task, we employed approximately 120 Robot data entries related to pushing actions (excluding pushing the red block) and around 250 Human data entries related to pushing actions (including the red block).

## A.6 Skill Generalization

In the novel skill generalization experiments in Section 4.3.3, skill generalization refers to the capability of a robotic system to acquire and execute skills that are **NOT present in the original robot training dataset** but can be learned by leveraging **human demonstration data**. To validate this capability, we construct novel skills by partitioning our self-collected dataset. For instance, we introduce a previously unseen skill P (Push)—while the robot training dataset contains no examples of P, the human dataset includes demonstrations of this skill P. Through our proposed algorithm, the robot is trained on this combined data and subsequently demonstrates the ability to successfully perform P (Push) in deployment.

For our experiments in Section 4.3.3, three new skills were introduced by controlling the selection of datasets: pushing a box, holding a bottle, and pressing a button. The Push Block task was trained using 40 Robot data entries related to pick-and-place tasks and 240 Human data entries related to pushing a box. The Hold Cup task utilized 40 Robot data entries related to grasping and 300 human data entries related to holding cups. The Press Button task was trained with 40 robot data entries related to grasping and 160 human data entries related to pressing actions.

## A.7 Training Details

We list the hyperparameters for **Stage 1 of Our Algorithm** in Table 3 **Stage2 of Our Algorithm** in Table 4. Stage 2 model is trained for 20 epoches with global batch size of 128 across 1 A800 gpu.

Batch Size	64
Optimizer	adamw
Learning Rate	5e-6
Cross Prediction Prob	0.6
inference step	50
video frame	16

Table 3: Hyperparameters of Stage 1

Batch Size	128
Optimizer	adamw
Learning Rate	1e-4
Weight Decay	1e-6
observation horizon	2
action horizon	8
<b>Resnet34</b>	
Size of Input Image	[160, 160, 3]
<b>DiT</b>	
embed dim	512
# decoder layers	6
# encoder layers	4
No. of heads	8
goal drop	0
embed pdrop	0
attn pdrop	0.3
resid pdrop	0.1
mlp pdrop	0.05
sigma data	0.5
sigma min	0.001
sigma max	80.0
$\lambda_f$	12/19
$\lambda_r$	4/19
$\lambda_p$	3/19
<b>PDCP</b>	
$\omega_c$	1.0
$\omega_p$	1.0
$\omega_m$	1.0

Table 4: Hyperparameters of Stage 2

## References

- [1] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, J. Ibarz, B. Ichter, A. Irpan, T. Jackson, S. Jesmonth, N. J. Joshi, R. C. Julian, D. Kalashnikov, Y. Kuang, I. Leal, K.-H. Lee, S. Levine, Y. Lu, U. Malla, D. Manjunath,

- I. Mordatch, O. Nachum, C. Parada, J. Peralta, E. Perez, K. Pertsch, J. Quiambao, K. Rao, M. S. Ryoo, G. Salazar, P. R. Sanketi, K. Sayed, J. Singh, S. A. Sontakke, A. Stone, C. Tan, H. Tran, V. Vanhoucke, S. Vega, Q. H. Vuong, F. Xia, T. Xiao, P. Xu, S. Xu, T. Yu, and B. Zitkovich. Rt-1: Robotics transformer for real-world control at scale. *ArXiv*, abs/2212.06817, 2022. URL <https://api.semanticscholar.org/CorpusID:254591260>.
- [2] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, K. Choromanski, T. Ding, D. Driess, K. A. Dubey, C. Finn, P. R. Florence, C. Fu, M. G. Arenas, K. Gopalakrishnan, K. Han, K. Hausman, A. Herzog, J. Hsu, B. Ichter, A. Irpan, N. J. Joshi, R. C. Julian, D. Kalashnikov, Y. Kuang, I. Leal, S. Levine, H. Michalewski, I. Mordatch, K. Pertsch, K. Rao, K. Reymann, M. S. Ryoo, G. Salazar, P. R. Sanketi, P. Sermanet, J. Singh, A. Singh, R. Soricut, H. Tran, V. Vanhoucke, Q. H. Vuong, A. Wahid, S. Welker, P. Wohlhart, T. Xiao, T. Yu, and B. Zitkovich. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *ArXiv*, abs/2307.15818, 2023. URL <https://api.semanticscholar.org/CorpusID:260293142>.
- [3] S. Kareer, D. Patel, R. Punamiya, P. Mathur, S. Cheng, C. Wang, J. Hoffman, and D. Xu. Egomimic: Scaling imitation learning via egocentric video. *arXiv preprint arXiv:2410.24221*, 2024.
- [4] Z. Chen, S. Chen, E. Arlaud, I. Laptev, and C. Schmid. Vividex: Learning vision-based dexterous manipulation from human videos. *arXiv preprint arXiv:2404.15709*, 2024.
- [5] Y. Qin, Y.-H. Wu, S. Liu, H. Jiang, R. Yang, Y. Fu, and X. Wang. Dexmv: Imitation learning for dexterous manipulation from human videos. In *European Conference on Computer Vision*, pages 570–587. Springer, 2022.
- [6] C. Wang, H. Shi, W. Wang, R. Zhang, L. Fei-Fei, and C. K. Liu. Dexcap: Scalable and portable mocap data collection system for dexterous manipulation. *arXiv preprint arXiv:2403.07788*, 2024.
- [7] A. Kannan, K. Shaw, S. Bahl, P. Mannam, and D. Pathak. Deft: Dexterous fine-tuning for real-world hand policies. *arXiv preprint arXiv:2310.19797*, 2023.
- [8] V. Jain, M. Attarian, N. J. Joshi, A. Wahid, D. Driess, Q. Vuong, P. R. Sanketi, P. Sermanet, S. Welker, C. Chan, et al. Vid2robot: End-to-end video-conditioned policy learning with cross-attention transformers. *arXiv preprint arXiv:2403.12943*, 2024.
- [9] C. Yuan, C. Wen, T. Zhang, and Y. Gao. General flow as foundation affordance for scalable robot learning. *arXiv preprint arXiv:2401.11439*, 2024.
- [10] Y. Ju, K. Hu, G. Zhang, G. Zhang, M. Jiang, and H. Xu. Robo-abc: Affordance generalization beyond categories via semantic correspondence for robot manipulation. In *European Conference on Computer Vision*, pages 222–239. Springer, 2024.
- [11] M. Xu, Z. Xu, C. Chi, M. Veloso, and S. Song. Xskill: Cross embodiment skill discovery. In *Conference on robot learning*, pages 3536–3555. PMLR, 2023.
- [12] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay: Long-horizon imitation learning by watching human play. *arXiv preprint arXiv:2302.12422*, 2023.
- [13] L. Smith, N. Dhawan, M. Zhang, P. Abbeel, and S. Levine. Avid: Learning multi-stage tasks via pixel-level translation of human videos. *arXiv preprint arXiv:1912.04443*, 2019.
- [14] H. Xiong, Q. Li, Y.-C. Chen, H. Bharadhwaj, S. Sinha, and A. Garg. Learning by watching: Physical imitation of manipulation skills from human videos. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 7827–7834. IEEE, 2021.



- [15] C. Wen, X. Lin, J. So, K. Chen, Q. Dou, Y. Gao, and P. Abbeel. Any-point trajectory modeling for policy learning. *arXiv preprint arXiv:2401.00025*, 2023.
- [16] J. Gu, S. Kirmani, P. Wohlhart, Y. Lu, M. G. Arenas, K. Rao, W. Yu, C. Fu, K. Gopalakrishnan, Z. Xu, et al. Rt-trajectory: Robotic task generalization via hindsight trajectory sketches. *arXiv preprint arXiv:2311.01977*, 2023.
- [17] S. Ye, J. Jang, B. Jeon, S. Joo, J. Yang, B. Peng, A. Mandlekar, R. Tan, Y.-W. Chao, B. Y. Lin, et al. Latent action pretraining from videos. *arXiv preprint arXiv:2410.11758*, 2024.
- [18] J. Liang, R. Liu, E. Ozguroglu, S. Sudhakar, A. Dave, P. Tokmakov, S. Song, and C. Vondrick. Dreamitate: Real-world visuomotor policy learning via video generation. *arXiv preprint arXiv:2406.16862*, 2024.
- [19] Y. Du, S. Yang, B. Dai, H. Dai, O. Nachum, J. Tenenbaum, D. Schuurmans, and P. Abbeel. Learning universal policies via text-guided video generation. *Advances in neural information processing systems*, 36:9156–9172, 2023.
- [20] H. Bharadhwaj, D. Dwibedi, A. Gupta, S. Tulsiani, C. Doersch, T. Xiao, D. Shah, F. Xia, D. Sadigh, and S. Kirmani. Gen2act: Human video generation in novel scenarios enables generalizable robot manipulation. *arXiv preprint arXiv:2409.16283*, 2024.
- [21] M. Shridhar, Y. L. Lo, and S. James. Generative image as action models. *arXiv preprint arXiv:2407.07875*, 2024.
- [22] M. Chang and S. Gupta. One-shot visual imitation via attributed waypoints and demonstration augmentation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5055–5062. IEEE, 2023.
- [23] H. Xiong, H. Fu, J. Zhang, C. Bao, Q. Zhang, Y. Huang, W. Xu, A. Garg, and C. Lu. Robotube: Learning household manipulation from human videos with simulated twin environments. In *Conference on Robot Learning*, pages 1–10. PMLR, 2023.
- [24] E. Chane-Sane, C. Schmid, and I. Laptev. Learning video-conditioned policies for unseen manipulation tasks. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 909–916. IEEE, 2023.
- [25] Z. Qian, M. You, H. Zhou, X. Xu, H. Fu, J. Xue, and B. He. Contrast, imitate, adapt: Learning robotic skills from raw human videos. *IEEE Transactions on Automation Science and Engineering*, 2024.
- [26] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *ArXiv*, abs/2006.11239, 2020.
- [27] Y. Song, J. N. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. *ArXiv*, abs/2011.13456, 2020.
- [28] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion. *ArXiv*, abs/2303.04137, 2023.
- [29] T. Pearce, T. Rashid, A. Kanervisto, D. Bignell, M. Sun, R. Georgescu, S. V. Macua, S. Z. Tan, I. Momennejad, K. Hofmann, and S. Devlin. Imitating human behaviour with diffusion models. *ArXiv*, abs/2301.10677, 2023.
- [30] M. Reuss, M. X. Li, X. Jia, and R. Lioutikov. Goal-conditioned imitation learning using score-based diffusion policies. *ArXiv*, abs/2304.02532, 2023. URL <https://api.semanticscholar.org/CorpusID:257952177>.

- [31] O. M. Team, D. Ghosh, H. R. Walke, K. Pertsch, K. Black, O. Mees, S. Dasari, J. Hejna, T. Kreiman, C. Xu, J. Luo, Y. L. Tan, P. R. Sanketi, Q. Vuong, T. Xiao, D. Sadigh, C. Finn, and S. Levine. Octo: An open-source generalist robot policy. *ArXiv*, abs/2405.12213, 2024.
- [32] M. Reuss, Ö. E. Yagmurlu, F. Wenzel, and R. Lioutikov. Multimodal diffusion transformer: Learning versatile behavior from multimodal goals. *ArXiv*, abs/2407.05996, 2024.
- [33] S. Liu, L. Wu, B. Li, H. Tan, H. Chen, Z. Wang, K. Xu, H. Su, and J. Zhu. Rdt-1b: a diffusion foundation model for bimanual manipulation. *ArXiv*, abs/2410.07864, 2024. URL <https://api.semanticscholar.org/CorpusID:273233386>.
- [34] Y. Wang, L. Wang, Y. Jiang, W. Zou, T. Liu, X. Song, W. Wang, L. Xiao, J. Wu, J. Duan, et al. Diffusion actor-critic with entropy regulator. *Advances in Neural Information Processing Systems*, 37:54183–54204, 2024.
- [35] V. Jain, T. Akhoun-Sadegh, and S. Ravanbakhsh. Sampling from energy-based policies using diffusion. *arXiv preprint arXiv:2410.01312*, 2024.
- [36] S. Li, R. Krohn, T. Chen, A. Ajay, P. Agrawal, and G. Chalvatzaki. Learning multimodal behaviors from scratch with diffusion policy gradient. *Advances in Neural Information Processing Systems*, 37:38456–38479, 2024.
- [37] R. Zhang, Z. Luo, J. Sjölund, T. Schön, and P. Mattsson. Entropy-regularized diffusion policy with q-ensembles for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 37:98871–98897, 2024.
- [38] J. Li, P. Zhou, C. Xiong, and S. C. Hoi. Prototypical contrastive learning of unsupervised representations. *arXiv preprint arXiv:2005.04966*, 2020.
- [39] S. Mo, Z. Sun, and C. Li. Siamese prototypical contrastive learning. *arXiv preprint arXiv:2208.08819*, 2022.
- [40] R. A. Potamias, J. Zhang, J. Deng, and S. Zafeiriou. Wilor: End-to-end 3d hand localization and reconstruction in-the-wild. *ArXiv*, abs/2409.12259, 2024. URL <https://api.semanticscholar.org/CorpusID:272753318>.
- [41] Y. Qin, W. Yang, B. Huang, K. V. Wyk, H. Su, X. Wang, Y.-W. Chao, and D. Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. *ArXiv*, abs/2307.04577, 2023. URL <https://api.semanticscholar.org/CorpusID:259367735>.