# MimicPlay: Long-Horizon Imitation Learning by Watching Human Play

Chen Wang[1], Linxi Fan[2], Jiankai Sun[1], Ruohan Zhang[1],
Li Fei-Fei[1], Danfei Xu[2][3], Yuke Zhu[2][4][†], Anima Anandkumar[2][5][†]
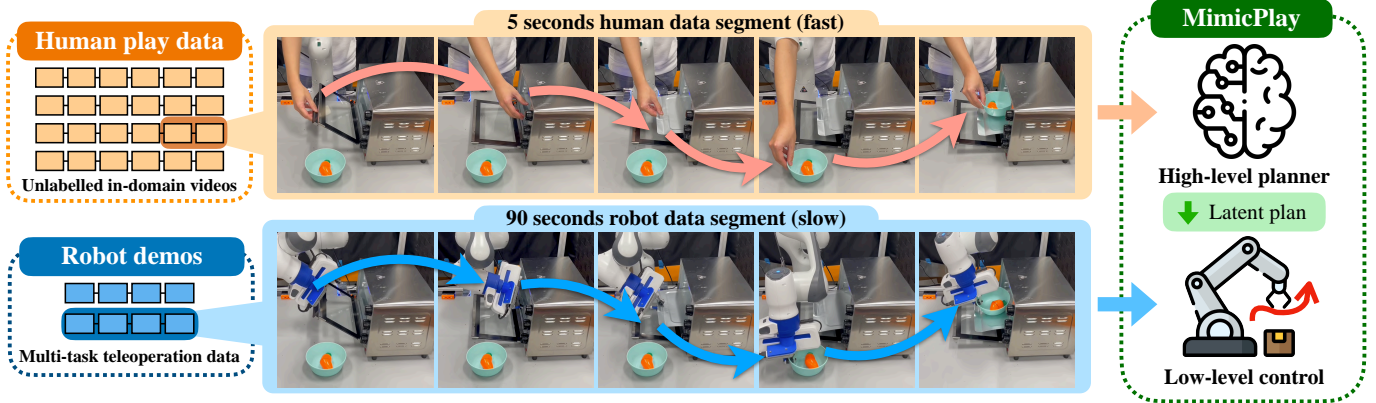[1]Stanford, [2]NVIDIA, [3]Georgia Tech, [4]UT Austin, [5]Caltech, [†]Equal Advising

Fig. 1: Human is able to complete a long-horizon task much faster than a teleoperated robot. This observation inspires us to develop `MimicPlay`, a hierarchical imitation learning algorithm that learns a high-level planner from cheap human play data and a low-level control policy from a small amount of multi-task teleoperated robot demonstrations. We show that `MimicPlay` significantly improves sample efficiency and robustness of imitation learning algorithms in long-horizon tasks.

*Abstract*—**Imitation Learning from human demonstrations is a promising paradigm to teach robots manipulation skills in the real world, but learning complex long-horizon tasks often requires an unattainable amount of demonstrations. To reduce the high data requirement, we resort to *human play data*— video sequences of people freely interacting with the environment using their hands. We hypothesize that even with different morphologies, human play data contain rich and salient information about physical interactions that can readily facilitate robot policy learning. Motivated by this, we introduce a hierarchical learning framework named MimicPlay that learns latent plans from human play data to guide low-level visuomotor control trained on a small number of teleoperated demonstrations. With systematic evaluations of 14 long-horizon manipulation tasks in the real world, we show that MimicPlay dramatically outperforms state-of-the-art imitation learning methods in task success rate, generalization ability, and robustness to disturbances. More details and video results could be found at https://mimic-play.github.io.**

## I. INTRODUCTION

Efficiently teaching robots to perform general-purpose manipulation tasks is a long-standing challenge. Imitation Learning (IL) has recently made considerable strides towards this goal, especially through supervised training using either human teleoperated demonstrations or trajectories of expert policies [37, 61]. Despite this promise, IL methods have been confined to learning short-horizon primitives, such as opening a door or picking up a specific object. The time cost and labor intensity of collecting long-horizon demonstrations, especially for complex real-world tasks with wide initial and goal condition distributions, remains a key barrier to the widespread adoption of IL methods.

To scale up imitation learning to complex long-horizon tasks, two connected directions have emerged in recent literature: *hierarchical imitation learning* and *learning from play data*. The former aims to increase learning sample efficiency by decoupling end-to-end deep imitation learning into the learning of high-level planners and low-level visuomotor controllers [32, 49], and the latter leverages an alternative form of robot training data, named *play data* [29]. Play data is typically collected through human-teleoperated robots interacting with their environment without specific task goals or guidance. Prior works show that data collected this way covers more diverse behaviors and situations compared to typical task-oriented demonstrations [11, 29]. The methods for learning from such play data often seek to uncover such diverse behaviors by training a hierarchical policy [29], where the high-level planner captures the distribution of intent and the low-level policies learn goal-directed control. However, collecting such play data in the real world can be very costly. For example, C-BeT [11] requires 4.5 hours of play data to learn manipulation skills in a specific scene, and TACO-RL [38] needs 6 hours of play data for one 3D tabletop environment.

In this work, we argue that the data required for learning high-level plan and low-level control can come in different forms, and doing so could substantially reduce the cost of

imitation learning for complex long-horizon tasks. Based on this argument, we introduce a new learning paradigm, in which robots learn high-level plans from *human play data*, where humans use their hands to interact with the environment freely. Human play data is much faster and easier to collect than robot teleoperation data (Fig. 1). It allows us to collect data *at scale* and cover a wide variety of situations and behaviors. We show that such scalability plays a key role in strong policy generalization. The robot then learns low-level manipulation policies from a small amount of *demonstration data*, which is collected by humans teleoperating with the robots. Unlike human play data, demonstration data is expensive to collect but does not lead to issues due to the mismatch between human and robot embodiments.

To scale imitation learning to long-horizon manipulation tasks, we present `MimicPlay`, a new imitation learning algorithm that leverages the complementary strengths of two data sources mentioned above: human play data and demonstration data. The robot first learns a goal-conditioned latent planner from *human play data* by predicting the future 3D human hand trajectories conditioned on the goal images. Such latent plans provide coarse guidance at each time step, tackling the challenging long-horizon manipulation problem by converting it into a guided motion generation process. Conditioned on these latent plans, the low-level controller incorporates state information that is essential for fined-grained manipulation to generate the final actions. We evaluate our method on 14 long-horizon manipulation tasks with a Franka Emika robot arm in six environments. Our results demonstrate significant improvement over state-of-the-art imitation learning methods in terms of sample efficiency and generalization abilities.

To summarize, the main contributions of our work are as follows:

1) We introduce a new paradigm for learning latent plans from easy-to-collect human play data.
2) We develop a plan-guided multi-task imitation learning algorithm based on the pre-trained latent planner for solving long-horizon manipulation tasks efficiently.
3) In 14 long-horizon evaluation tasks across six scenes, our approach outperforms the state-of-the-art baselines for over 50% in task success rate. It also exhibits stronger generalization in previously unseen tasks, with more than 40% improvements in task success rate.

## II. RELATED WORK

### A. Imitation Learning from Demonstrations

Imitation Learning (IL) has enabled robots to successfully perform various manipulation tasks [1, 4, 9, 14, 15, 22, 26, 40]. Traditional IL algorithms such as DMP and PrMP [25, 35, 36, 41] enjoy high learning sample efficiency but are limited in their ability to handle high-dimensional observations and settings that require closed-loop control. In contrast, recent IL methods built upon deep neural networks can learn reactive policies from raw demonstration data [16, 17, 32, 33, 56, 61, 62]. While these methods offer greater flexibility, they

require a large number of human demonstrations to learn even simple pick-and-place tasks. Collecting high-quality human demonstrations on robots remains labor- and resource-intensive [6, 7]. For example, RT-1 [7] reported using 130k teleoperated demonstrations collected over the course of 17 months across 13 robots. Our work instead proposes to leverage human play data, which does not require robot hardware and can be collected efficiently, to reduce the need for on-robot demonstration data dramatically.

Our idea of learning a hierarchical policy from offline demonstrations is also related to prior works [29, 31, 32, 49, 59]. However, all previous methods focus on learning both planning and control with a single type of data—teleoperated robot demonstrations, which is expensive to collect. Our approach uses cheap human play data for learning high-level planning and a small number of robot demonstrations for learning low-level control, which significantly strengthens the model's planning capability while keeping a low demand on demonstration data.

Noticeably, there are other approaches used for improving sample efficiency for imitation learning, such as grounding action on discrete observation space [23, 44, 50, 51]. In this work, we focus on analyzing the cheap human play data for improving sample efficiency, which is orthogonal to the selection of action space and could be combined with existing works for better performance.

### B. Learning from Human Videos

A plethora of recent research has explored leveraging large-scale human video data to improve robot policy learning [10, 12, 13, 42, 43, 46, 47, 48, 53, 58, 60]. Closely related to our work are R3M [34] and MVP [57], which use an Internet-scale human video dataset Ego4D [18] to pretrain visual representations for subsequent imitation learning. However, due to diversity in the data source and large domain gaps, transferring the pre-trained representation to a specific manipulation task might be difficult. Notably, Hansen et al. [20] found simple data augmentation techniques could have similar effects as these pre-trained representations. To reduce the domain gap, another thread of work [2, 27, 28, 47, 52, 53, 58] utilizes in-domain human videos, where human directly interacts with the robot task environment with their own hands. Such type of data has a smaller gap between human and robot domains, which allows sample-efficient reward shaping for training RL agents [2, 27, 52, 53] and imitation learning [28, 47, 58]. However, these works focus on learning either task rewards or features from human videos, which doesn't directly help the low-level robot action generation. In this work, we extract meaningful trajectory-level task plans from human play data, which provides high-level guidance for the low-level controller for solving challenging long-horizon manipulation tasks.

### C. Learning from Play Data

Our idea of leveraging human play data is heavily inspired by learning from play [3, 11, 29, 39], an alternative imitation learning paradigm that focuses on learning from play data, a
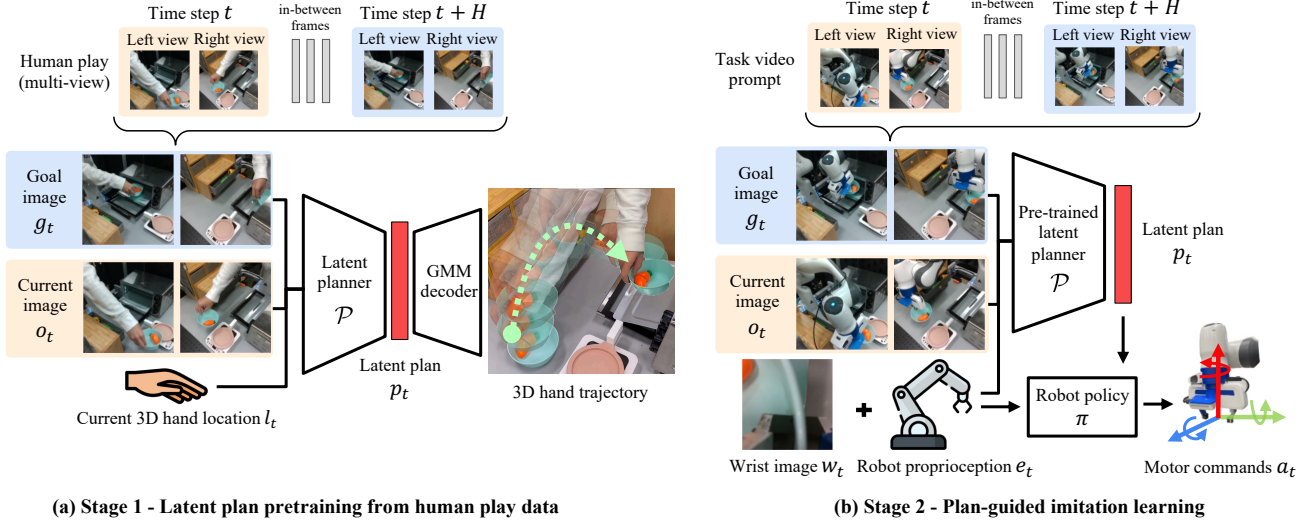
(a) Stage 1 - Latent plan pretraining from human play data

(b) Stage 2 - Plan-guided imitation learning

Fig. 2: `MimicPlay` is a two-stage framework for learning long-horizon manipulation tasks. **Stage 1**: a goal-conditioned trajectory generation model is trained to extract latent plans from *human play data*. The goal is to learn a latent plan space that contains coarse guidance for diverse task goals. **Stage 2**: the pre-trained planner is leveraged to generate latent plans conditioned on the goal image sampled from the task video prompt. Based on generated latent plans, the robot policy is trained with a plan-guided imitation learning algorithm.

form of on-robot human teleoperation demonstration provided without a specific goal. Compared to common goal-directed demonstrations, play data exhibit greater diversity in behavior and coverage of the environment state space [29].

However, collecting play data still requires the laborious teleoperation process [11, 38]. In this work, instead of using teleoperated play data, we learn from human play data, where humans freely interact with the scene with their hands. This form of data is significantly more efficient to collect than teleoperated play data.

We empirically show that we can overcome the domain gap by incorporating a small amount of additional robot demonstration data, and this combination can substantially improve the model performance and generalization in unseen situations.

## III. MIMICPLAY

Training a robot to perform a long-horizon task is challenging since it could require learning both high-level plans and low-level control. `MimicPlay` is based on the key insight that high-level planning is best learned from human play data that are fast to collect, while low-level control skills are best acquired from teleoperated demonstration data that do not have any embodiment gap. In particular, due to the difference in the embodiments of humans and robots, it is crucial to find an intermediate representation that can bridge the gap between the two data sources. `MimicPlay` addresses this challenge by learning a latent planning space to extract diverse plans from cost-effective human play data.

More specifically, `MimicPlay` (Fig. 2) is a two-stage framework for learning long-horizon manipulation by leveraging cheap and easy-to-collect human play data. In the first stage, we use a goal-conditioned 3D hand trajectory generation task as a pretraining objective to form the latent plan space from human play data. The goal is to learn a latent plan space

that encodes the high-dimensional visual inputs into coarse guidance for diverse task goals. In the second stage, the robot is trained to learn a low-level control policy given guidance via goal-conditioned behavior cloning. Here we first introduce the algorithm that trains the latent planner from human play data, then describe the details of the plan-guided imitation learning algorithm.

### A. Learning latent plans from human play data

We consider the problem of learning a goal-conditioned latent planner $\mathcal{P}$ that outputs a latent plan $p_t$ given input $I_t = \{o_t, g_t, l_t\}$ consisting of a visual observation $o_t$, a future goal image $g_t$, and current hand location $l_t$. During training, the goal image $g_t$ is specified as the future frame after $H$ steps of actions. During testing, a task video prompt $\mathcal{V}$ is used as the source of goal images that are sent to $\mathcal{P}$. The goal of the latent planner is to generate high-level guidance for the low-level controller, which means it needs to be capable of handling diverse task goals and situations, e.g., in the setting of multi-task learning. As a result, training such a planner requires a large dataset. We address this issue by leveraging a cheap and easy-to-collect data source—*human play data*.

**Collecting human play data.** A human operator is asked to interact with interesting objects in an environment and engage in behaviors that satisfy their own curiosity. For instance, in a kitchen environment, the operator might open the oven then pull out the tray or pick up a pan and place it on the stove.

This type of data contains rich state transitions and implicit human knowledge of object affordances situated in an environment. More importantly, collecting human play data is extremely cheap and efficient since it does not require task definition or labeling and only takes a small fraction of time to collect—a human operator can finish a 90-second robot teleoperation task within just five seconds (Fig. 1). In our experiments, we collect 10 minutes of human play data for

each scene as the dataset for training the latent planner, which contains the number of image-wise state transitions equivalent to a 3-hour robot teleoperation data.

**Tracking 3D human hand trajectories.** Common human video datasets are collected with single-view observations. Our task environments have rich 3D structures hence human hand motion is in 3D space. The latent plans generated based on 2D observation are ambiguous along the depth axes. We instead use a calibrated dual-camera system for tracking 3D hand trajectories from human play data. From each viewpoint, we use an off-the-shelf hand detector [45] to localize the hand location in the 2D image space. Based on multi-view detection results, we can reconstruct the full 3D hand trajectory $\tau$ in the world coordinate frame from the human play data (Fig. 2(a)).

**Learning latent plans.** With the collected human play data and corresponding 3D hand trajectory $\tau$, we formalize the latent plan learning process as a goal-conditioned 3D trajectory generation task. More specifically, as shown in Fig. 2(a), two convolutional networks (CNNs) separately process the visual observations $o_t$ and goal image $g_t$ into low-dimensional features, which are further concatenated together with the hand location $l_t$ and processed by an MLP-based encoder network into a latent plan feature $p_t$. From the latent plan $p_t$, an MLP-based decoder network generates the prediction of the 3D hand trajectory. However, simple regression of the trajectory cannot fully cover the rich multimodal distribution of human motions. In fact, even for the same human operator, one task goal might be achieved in different ways. To address this challenge, we use an MLP-based Gaussian Mixture Model (GMM) [5] to model the trajectory distribution from the latent plan $p_t$.

For a GMM as Eq. 1 shows:

$$p(\boldsymbol{\tau}|\boldsymbol{\theta}) = \sum_{\boldsymbol{z}} p(\boldsymbol{\tau}|\boldsymbol{\theta}, \boldsymbol{z}) p(\boldsymbol{z}|\boldsymbol{\theta}), \tag{1}$$

where $\boldsymbol{\theta}$ is the parameters of the GMM and $p(\boldsymbol{\tau}|\boldsymbol{\theta}, \boldsymbol{z}_k)$ is a Gaussian distribution $\mathcal{N}(\boldsymbol{\tau}|\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k)$ with $\boldsymbol{z}$ consistsing of $K$ components. A specific weight $\eta_k$ represents the probability of the $k$-th component.

GMMs are more expressive than simple Gaussians, because they are designed to capture the multi-modality which is inherited in the human play data. The final learning objective of our GMM model is to minimize the negative log-likelihood of the detected 3D human hand trajectory $\tau$ as Eq. 2

$$\mathcal{L}_{\text{GMM}}(\boldsymbol{\theta}) = -\mathbb{E}_{\boldsymbol{\tau}} \log \left( \sum_{k=1}^{K} \eta_k \mathcal{N}(\boldsymbol{\tau}|\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k) \right),$$

$$\text{where } 0 \leq \eta_k \leq 1, \sum_{k=1}^{K} \eta_k = 1 \text{ and } \boldsymbol{\theta} = \{\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k, \eta_k\}_{k=1}^{K}$$

$$\tag{2}$$

**Handling visual gap between human and robot domains.** In this work, we consider the setup where the human and the robot interact in the same environment. However, different visual appearance (for example, top vs. bottom row in Fig. 1) between human and robot domains poses a challenge in transferring the learned latent planner to the robot downstream

task. We handle this challenge by introducing another learning objective that leverages a small set of robot demonstration data (introduced in Sec. III-B). To reduce the visual gap between the human domain $\mathcal{Q}_H$ and the robot domain $\mathcal{Q}_R$, we minimize the output distribution of the visual encoders with a Kullback–Leibler (KL) divergence loss, as Eq. 3 shows:

$$\mathcal{L}_{\text{KL}} = D_{\text{KL}}(\mathcal{Q}_R||\mathcal{Q}_H) = \mathcal{Q}_R \log \frac{\mathcal{Q}_R}{\mathcal{Q}_H}, \tag{3}$$

Another input to the pre-trained latent planner is the hand location $l_t$ in 3D space. When processing the data from the robot domain $Q_R$, we use the robot's end-effector 3D location as the $l_t$. The final loss function for training the latent planner is:

$$\mathcal{L} = \mathcal{L}_{\text{GMM}} + \lambda \cdot \mathcal{L}_{\text{KL}}, \tag{4}$$

where $\lambda$ is a hyperparameter that controls the weights between two losses.

### B. Plan-guided multi-task imitation learning

The robot low-level policy $\pi$ is trained with a guided behavior cloning algorithm with robot demonstration data collected through teleoperation. Prior works often learn visuomotor policies end-to-end from scratch [33] or rely on pre-trained visual representations [34, 57]. However, since visual inputs do not provide task-level information, the policy model needs to learn both planning and low-level control, which reduces sample efficiency. In this work, we leverage our pre-trained latent planner $\mathcal{P}$ to generate latent plan $p_t$ as guidance for the low-level robot action generation $a_t = \pi((w_t, e_t)|p_t)$ (Fig. 2(b)), where $w_t$ is the wrist camera observation and $e_t$ is the end-effector 6D proprioception of the robot. In the following, we introduce how to generate the latent plan $p_t$ with the pre-trained latent planner $\mathcal{P}$ and the details of training the low-level controller $\pi$.

**Video-conditioned latent plan generation.** Prompting a robot for a long-horizon task is challenging since it requires a clear specification of the entire task, which might contain multiple subgoals. In this work, we use teleoperated robot task video $\mathcal{V}$ to prompt the pre-trained latent planner to generate corresponding plans $p_t = \mathcal{P}(o_t, g_t, l_t), g_t \in \mathcal{V}$. During training, we specify the goal image $g_t$ as the frame $H$ steps after the input observation $o_t$ in the training demonstration. $H$ is a uniformly sampled integer number within the range of $[200, 600]$, which equals 10-30 seconds of robot motion in wall-clock time. $l_t$ here is the 3D location of the robot's end-effector. During inference, we assume access to a single-shot robot video which is used as a source of goal images. The goal image will start at the 200 frame of the task video and move to the next $i$ frame after each step. We use $i = 1$ in all our experiments. Based on the inputs, the latent planner generates a latent plan feature embedding $p_t$ of shape $\mathbb{R}^{1 \times d}$, which is used as guidance for the low-level robot policy.

**Transformer-based plan-guided imitation.** The latent plan $p_t$ output by the pre-trained planner contains high-level guidance of what the robot should do at the current time step, which allows the policy $\pi$ to focus on learning low-level

**(a) Kitchen**

**(b) Study desk**

**(c) Flower**
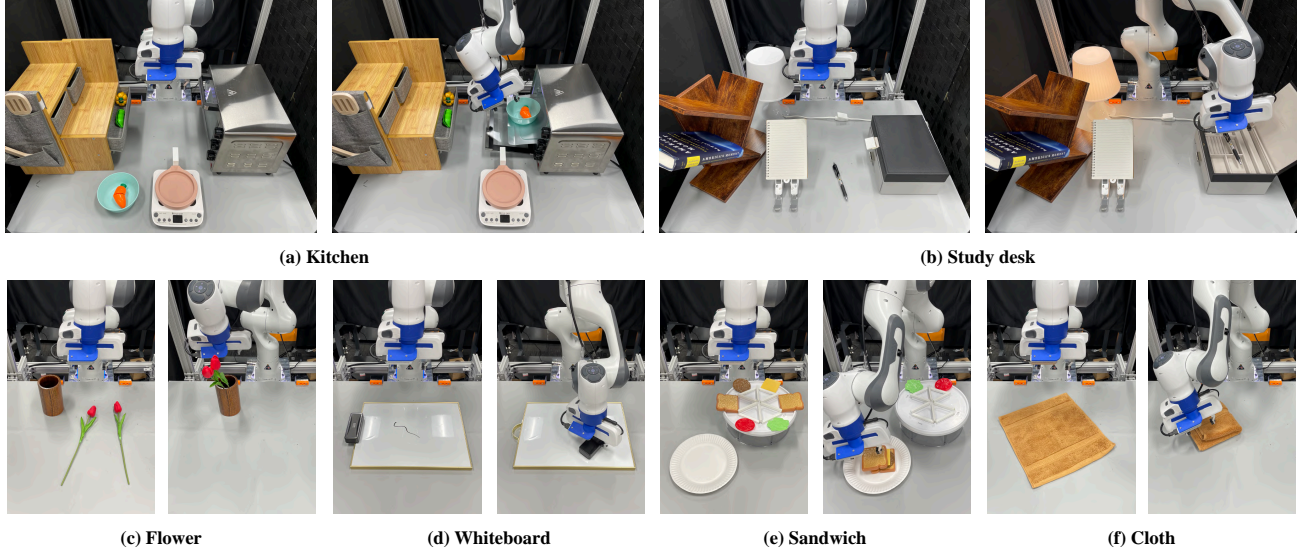
**(d) Whiteboard**

**(e) Sandwich**

**(f) Cloth**

Fig. 3: **Evaluation Tasks.** We design 6 environments with a Franka Emika robot arm. The left images for each environment are the initial states, and the right images are the goal states. **(a) Kitchen (long-horizon)**: open the oven, pull out the tray, pick up the bowl, and place the bowl on the tray. **(b) Study Desk (long-horizon)**: turn on the light, open the box, pick up the pen, and place the pen in the box. **(c) Flower**: pick up each flower and insert them into the vase. The flowers are randomly placed in the center area of the table. The vase is placed on either the left or the right side. **(d) Whiteboard**: pick up the eraser and erase randomly drawn curve lines on the whiteboard. **(e) Sandwich**: select the ingredients and place them on a side plate to make a cheeseburger or a sandwich. **(f) Cloth**: fold the towel twice. The initial location of the towel is randomly sampled.

actions by following the guidance. We introduce a plan-guided imitation algorithm based on the pre-trained latent plans. The policy model $\pi$ first processes the robot wrist camera observation $w_t$ and the proprioception $e_t$ information into low-dimensional feature vectors of shape $\mathbb{R}^{1 \times d}$. Then, we concatenate these features with the generated latent plan $p_t$ to get a one-step token embedding $s_t = [w_t, e_t, p_t]$. The embedding sequence of $T$ time steps is then represented as $s_{t:t+T} = [w_t, e_t, p_t, \cdots, w_{t+T}, e_{t+T}, p_{t+T}]$, which passes through a transformer architecture [55]. The transformer model $f_{\text{trans}}$ processes the input embeddings using its $N$ layers of self-attention and feed-forward neural networks. Given an embedding sequence of $T - 1$ time steps, $f_{\text{trans}}$ generates the embedding of trajectory prediction in an autoregressive way, as Eq. 5 shows:

$$x_T = f_{\text{trans}}(w_{1:T-1}, e_{1:T-1}, p_{1:T-1}). \quad (5)$$

where $x_T$ is the predicted action embedding at time step $T$.

The transformer architecture uses the multi-head self-attention mechanism to gather context and dependencies from the entire history trajectory at each step. The final robot control commands $a_t$ are computed by processing the action feature $x_t$ with a two-layer fully-connected network. To handle the multimodal distribution of robot actions, we also use a MLP-based GMM model [5] for the action generation.

**Multi-task learning.** One reason for learning latent plans from human play data is to extract diverse task-level plans, which is suitable for multi-task learning. In our experiment setting, each environment contains multiple tasks. For all tasks in the same environment, we train them within a single model

architecture $\{\mathcal{P}, \pi\}$ conditioning on different video prompts. For each training sample, the prompting video is uniformly sampled from the training videos of the same task category.

*C. Implementation details*

Here we lay down the details of the data collection, training, and testing process.

**Learning from human play data.** The human play data is collected by letting a human operator directly interact with the scene with a single hand for 10 minutes for each scene. The entire trajectory $\tau$ is recorded at the speed of 60 frames per second and is used without cutting or labeling. The 3D hand trajectory is detected with an off-the-shelf multi-view human hand tracker [45]. The total number of video frames within 10 minutes of human play video is around 36k. We train one latent planner for each environment with the collected human play data. For the multi-environment setup (for the experiments in Tab. III), we merge the human play data from each scene to train a single latent planner. The latent planner contains two ResNet-18 [21] networks for image processing and MLP-based encoder-decoder networks together with a GMM model, which has $K = 5$ distribution components. We train 100k iterations for the latent planner which takes a single GPU machine for 12 hours.

**Learning from robot demonstrations.** The robot teleoperation data is collected with an IMU-based phone teleoperation system RoboTurk [30]. The control frequency of the robot arm is 17-20Hz and the gripper is controlled at 2Hz. For each task, we collect 20 demonstrations. In the experiments, we also have a 40 demonstration dataset for testing the sample efficiency of
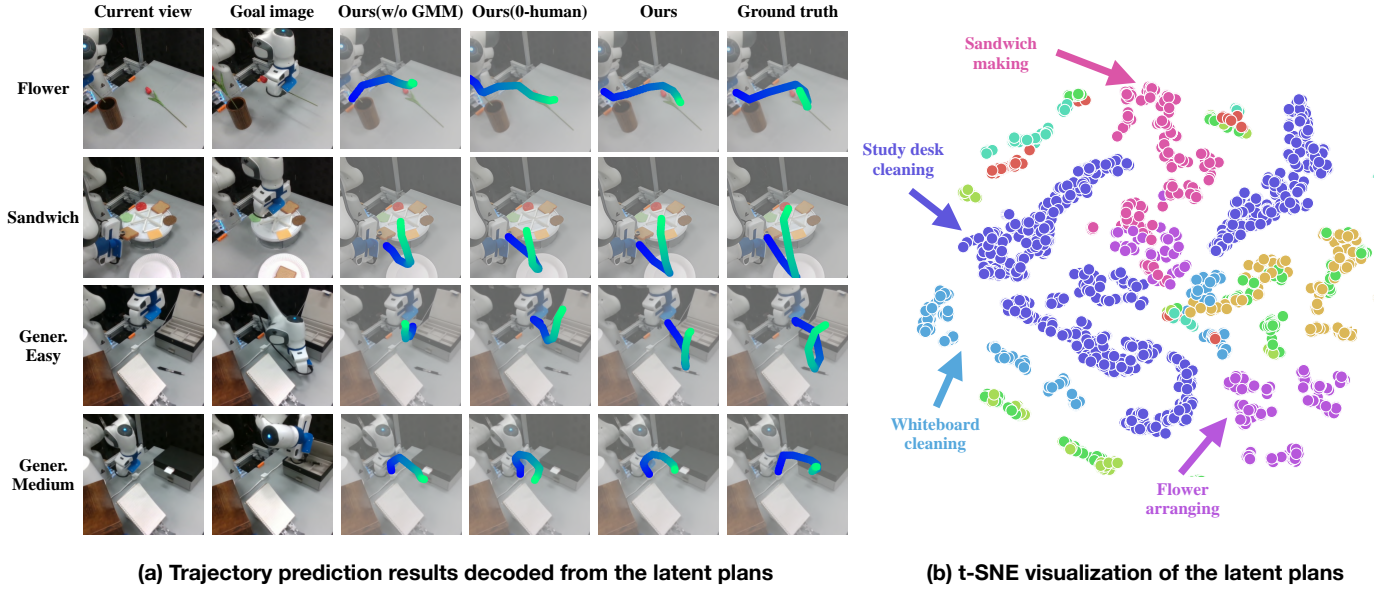
| Current view | Goal image | Ours(w/o GMM) | Ours(0-human) | Ours | Ground truth |

**(a) Trajectory prediction results decoded from the latent plans**

**(b) t-SNE visualization of the latent plans**

Fig. 4: Qualitative visualization of the learned latent plan. (**a**) Visualization of the trajectory prediction results decoded from the latent plans learned by different methods. The fading color of the trajectory from blue to green indicates the time step from 1 to 10. (**b**) t-SNE visualization of latent plans, the latent plans of the same task tend to cluster in the latent space.

TABLE I: Quantitative evaluation results in the Kitchen environment in terms of the number of successful runs out of all trials.

| | | Subgoal (first subgoal) | | | | Long horizon ($\geq$ 3 subgoals) | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Task-1 | Task-2 | Task-3 | ALL | Task-1 | Task-2 | Task-3 | ALL |
| 20 demos | GC-BC (BC-RNN) [33] | $1/10$ | $0/10$ | $1/10$ | $2/30$ | $0/10$ | $0/10$ | $0/10$ | $0/30$ |
| | GC-BC (BC-trans) [8] | $2/10$ | $0/10$ | $0/10$ | $2/30$ | $0/10$ | $0/10$ | $0/10$ | $0/30$ |
| | C-BeT [11] | $5/10$ | $6/10$ | $0/10$ | $11/30$ | $0/10$ | $0/10$ | $0/10$ | $0/30$ |
| | LMP [29] | $3/10$ | $1/10$ | $2/10$ | $6/30$ | $1/10$ | $0/10$ | $1/10$ | $2/30$ |
| | R3M-BC [34] | $9/10$ | $0/10$ | $0/10$ | $9/30$ | $0/10$ | $0/10$ | $0/10$ | $0/30$ |
| | Ours (0-human) | $\mathbf{10/10}$ | $5/10$ | $3/10$ | $18/30$ | $3/10$ | $1/10$ | $3/10$ | $7/30$ |
| | Ours | $\mathbf{10/10}$ | $8/10$ | $7/10$ | $25/30$ | $7/10$ | $3/10$ | $4/10$ | $14/30$ |
| 40 demos | GC-BC (BC-RNN) [33] | $1/10$ | $2/10$ | $2/10$ | $5/30$ | $0/10$ | $0/10$ | $1/10$ | $1/30$ |
| | GC-BC (BC-trans) [8] | $3/10$ | $7/10$ | $6/10$ | $16/30$ | $0/10$ | $0/10$ | $1/10$ | $1/30$ |
| | C-BeT [11] | $4/10$ | $\mathbf{10/10}$ | $0/10$ | $14/30$ | $0/10$ | $0/10$ | $0/10$ | $0/30$ |
| | LMP [29] | $6/10$ | $3/10$ | $2/10$ | $11/30$ | $3/10$ | $1/10$ | $0/10$ | $4/30$ |
| | R3M-BC [34] | $5/10$ | $4/10$ | $0/10$ | $9/30$ | $5/10$ | $0/10$ | $0/10$ | $5/30$ |
| | Ours (0-human) | $\mathbf{10/10}$ | $5/10$ | $5/10$ | $20/30$ | $4/10$ | $3/10$ | $5/10$ | $12/30$ |
| | Ours | $\mathbf{10/10}$ | $9/10$ | $8/10$ | $27/30$ | $7/10$ | $6/10$ | $8/10$ | $21/30$ |

different approaches. The robot policy model is a GPT-style transformer [8], which consists of four multi-head layers with four heads. We train 100k iterations for the policy with a single GPU machine in 12 hours. For a fair comparison with our method, the baseline approaches trained without human play data have five more demonstrations during training the latent planner $\mathcal{P}$ and the low-level policy $\pi$.

**Testing.** We perform real-time inference on a Franka Emika robot arm with a control frequency of 17Hz—directly from raw image inputs to 6-DoF robot end-effector and gripper control commands with our trained models. The robot is controlled with the Operational Space Control (OSC) [24].

## IV. EXPERIMENTS

### A. Experiment setups

**Environments.** We design six environments with a total of 14 tasks for a Franka Emika robot arm, as illustrated in Fig. 3. These environments feature several manipulation challenges, such as contact-rich tool manipulation (cleaning the whiteboard), articulated-object manipulation (opening the oven and the box on the study desk), high-precision tasks (inserting flowers and turning on the lamp by pressing the button), and deformable object manipulation (folding cloth). See Appendix for more environment design details.

**Tasks.** We design three tasks in the Kitchen environment

TABLE II: Quantitative evaluation results in the Study desk environment

|  |  | Trained tasks | | | | | Generalization | | | |
|---|---|---|---|---|---|---|---|---|---|---|
|  |  | Task-1 | Task-2 | Task-3 | Task-4 | ALL | Easy | Medium | Hard | ALL |
| 20 demos | GC-BC (BC-trans) [8] | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{40}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{30}$ |
|  | LMP [29] | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{40}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{30}$ |
|  | Ours (0-human) | $\frac{2}{10}$ | $\frac{3}{10}$ | $\frac{1}{10}$ | $\frac{2}{10}$ | $\frac{8}{40}$ | $\frac{2}{10}$ | $\frac{1}{10}$ | $\frac{0}{10}$ | $\frac{3}{30}$ |
|  | Ours (50-human) | $\frac{3}{10}$ | $\frac{4}{10}$ | $\frac{1}{10}$ | $\frac{4}{10}$ | $\frac{12}{40}$ | $\frac{4}{10}$ | $\frac{3}{10}$ | $\frac{1}{10}$ | $\frac{8}{30}$ |
|  | Ours (w/o KL) | $\frac{3}{10}$ | $\frac{7}{10}$ | $\frac{3}{10}$ | $\frac{2}{10}$ | $\frac{15}{40}$ | $\frac{4}{10}$ | $\frac{2}{10}$ | $\frac{0}{10}$ | $\frac{6}{30}$ |
|  | Ours (w/o GMM) | $\frac{4}{10}$ | $\frac{2}{10}$ | $\frac{2}{10}$ | $\frac{3}{10}$ | $\frac{11}{40}$ | $\frac{2}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{2}{30}$ |
|  | Ours | $\frac{6}{10}$ | $\frac{7}{10}$ | $\frac{4}{10}$ | $\frac{5}{10}$ | $\frac{22}{40}$ | $\frac{7}{10}$ | $\frac{5}{10}$ | $\frac{2}{10}$ | $\frac{14}{30}$ |

TABLE III: Quantitative evaluation results of multi-task learning

|  | Spatial generalization | | Extreme long horizon | Deformable | |
|---|---|---|---|---|---|
|  | Flower | Whiteboard | Sandwich | Cloth | ALL |
| LMP-single | $\frac{1}{10}$ | $\frac{0}{10}$ | $\frac{1}{10}$ | $\frac{3}{10}$ | $\frac{5}{40}$ |
| LMP [29] | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{0}{10}$ | $\frac{2}{10}$ | $\frac{2}{40}$ |
| R3M-single | $\frac{2}{10}$ | $\frac{1}{10}$ | $\frac{3}{10}$ | $\frac{4}{10}$ | $\frac{10}{40}$ |
| R3M [34] | $\frac{1}{10}$ | $\frac{1}{10}$ | $\frac{2}{10}$ | $\frac{2}{10}$ | $\frac{6}{40}$ |
| Ours-single | $\frac{5}{10}$ | $\frac{5}{10}$ | $\frac{6}{10}$ | $\frac{7}{10}$ | $\frac{23}{40}$ |
| Ours | $\frac{4}{10}$ | $\frac{2}{10}$ | $\frac{8}{10}$ | $\frac{8}{10}$ | $\frac{22}{40}$ |

and four tasks in the Study desk environment. All these tasks have different goals. In this work, we focus on long-horizon tasks that require the robot to complete several sub-goals. To better analyze the performance of each method, we define the *Subgoal* task category that only counts whether the first subgoal of the task has been achieved and the *Long horizon* task category which is the full task. In the Study desk environment, we design three tasks for testing the compositional generalization ability of the models to novel task goal sequences, which are not included in the training dataset. These three tasks are classified as *Easy*, *Medium*, and *Hard* depending on their difference compared to the training tasks. The *Easy* task is a simple concatenation of two trained tasks and their subgoals. The *Medium* task contains an unseen composition of a pair of subgoals that is not covered by any trained tasks, i.e., the transition from subgoal $A$ to subgoal $B$ is new. The model needs to generate novel motions to reach these subgoals. The *Hard* task contains two such unseen transitions. For the rest four environment, each scene has one task goal and features different types of challenges in manipulation, *e.g.*, generalization to new spatial configuration, extremely long horizon, and deformable object manipulation. The details of each task are included in Appendix.

**Baselines.** We compare with five prior approaches:

- GC-BC (BC-RNN) [33]: Goal-conditioned behavior cloning algorithm [29] implemented with recurrent neural networks (RNN) [19].
- GC-BC (BC-trans) [8]: Another goal-conditioned behavior cloning algorithm implemented with GPT-like transformer architecture.
- C-BeT [11]: Goal-conditioned learning from teleoperated robot play data algorithm implemented with Behavior Transformer (BeT) [44].

- LMP [29]: A learning from teleoperated robot play data algorithm designed to handle variability in the play data by learning an embedding space. LMP (single) is a variant by training each task with a separate model.
- R3M-BC [34]: A goal-conditioned imitation learning framework that leverages R3M visual representation pre-trained with internet-scale human video dataset Ego4D [18]. R3M-BC (single) is a variant by training each task with a separate model.

**Ablations** We compare four variants of our model to showcase the effectiveness of our architecture design:

- Ours: `MimicPlay` with full collection (10 min) of human play data. Ours (single) is a variant by training each task with a separate model.
- Ours (0-human): variant of our model without using human play data. The pre-trained latent plan space is trained only with the teleoperated robot demonstrations.
- Ours (50-human): variant of our model where the latent planner is trained with 50% of human play data (5 min).
- Ours (w/o GMM): variant without using the GMM model for learning the latent plan space from human play data.
- Ours (w/o KL): Our approach without using KL loss for addressing the visual gap between human and robot data when pre-training the latent planner.

### B. Results

**Human play data significantly improves the performance.** Our method outperforms Ours (0-human) by more than 23% in long-horizon task settings over all trained tasks, as shown in both Tab. I(ALL) and II(Trained tasks ALL). This result showcases that learning a latent plan space doesn't need to fully rely on teleoperated robot demonstration data. A 10-minutes of cheap and unlabelled human play data brings large improvements in the task success rate and sample efficiency.

**The two-stage framework is important for learning long-horizon tasks.** Ours (0-human) trained with our two-stage framework outperform prior end-to-end learning methods in the long-horizon task settings by more than 15%, as is shown in Tab. I(ALL) and II (Trained tasks ALL). This result show-cases that end-to-end learning for both planning and control is less effective than learning to act based on pre-trained latent plans for long-horizon tasks.

**Latent plan pre-training helps multi-task learning.** In Tab. III, we test how each method performs when training each

task with a separate model. For end-to-end learning approaches (e.g., LMP and R3M-BC), training task-specific models will lead to better performance (LMP-single vs. LMP; R3M-single vs. R3M). These results showcase the difficulty of learning multiple tasks with a single model. However, our approach has the smallest performance drop in multi-task training (Ours-single vs. Ours). These findings highlight the advantage of learning plan-guided low-level robot motions based on the pre-trained latent task space. We use t-SNE [54] to visualize the generated latent plans conditioned on different tasks, as shown in Fig. 4(b). We find that the latent plans of the same task tend to cluster in the latent space, which shows the effectiveness of our approach in distinguishing different tasks. However, we do observe an uneven performance drop with our method (the success rate of the whiteboard task drops from $^5/_{10}$ to $^2/_{10}$). The reason might be that the length of the demonstration for the whiteboard task is shorter than the other tasks, which leads to an imbalanced training dataset.

**Transformer architecture helps multi-task learning.** In Tab. I, GC-BC (BC-trans) with the GPT transformer architecture outperforms GC-BC (BC-RNN) by more than 30% in a 40-demos Subgoal setting. However, the performance of GC-BC (BC-trans) quickly drops to the same level as GC-BC (BC-RNN) in 20-demos settings. The result showcases that training vision-based transformer policy end-to-end requires more data.

**GMM is crucial for learning latent plans from human play data.** In Tab. II, our full with GMM model largely outperforms Ours (w/o GMM). Although being trained with full human play data, Ours (w/o GMM) even fails to match the performance of Ours (0-human) in the generalization task settings. In addition, in Fig. 4, the quality of the trajectory generated by Ours (w/o GMM) is also the worst. This result highlights the importance of using the GMM model to handle the multimodal distribution of the hand trajectory when learning the latent plan space.

**KL loss helps minimize the visual gap between human and robot data.** In Tab. II, although Ours (w/o KL) baseline outperforms most baselines in trained tasks, its success rate is 17% lower than Ours. In the generalization setting, Ours (w/o KL) fails to match the performance of Ours (50-human). These results showcase that the visual gap between the human play data and robot data exists, and KL loss helps close the gap when training the latent planner. More analysis of the distribution shifts between human play data and robot data can be found in the Appendix.

**The scale of the human play data matters.** In Tab. II, we compared the model variants with 50% human play data (Ours (50-human)) and found it fails to match the performance of Ours, which has access to 100% human play data. Most critically, in the unseen task settings, using more human play data to cover unseen cases in the training set significantly benefits generalization (Ours vs Ours (50-human)).

**Human play data improves generalization to new subgoal compositions.** In Tab. II generalization setting, Ours surpasses all baselines by more than 35%. This result high-
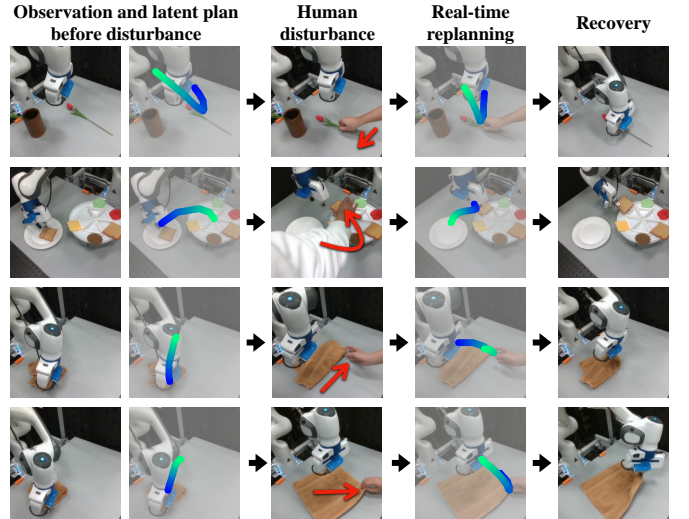


Fig. 5: Qualitative visualization of the latent plans before the disturbance and re-planning. **Column 1**: third-person view. **Column 2**: visualization of the latent plan before disturbance. **Column 3**: human disturbance; the red arrow indicts the direction of disturbance. **Column 4**: visualization of real-time re-planning capabilities, which show robustness against disturbance. **Column 5**: robot recovers with the updated task plan.

lights that our approach indeed extracts novel latent plans from human play data and guides the robot's low-level policy to generalize to new compositions of subgoals. We visualize the 3D trajectory decoded from the latent plan by projecting it onto the 2D image in Fig. 4. In the last two rows, we showcase the results of two unseen subgoal transitions. The trajectory generated by our model is most similar to the ground truth trajectory, while Ours (0-human) is overfitted to the subgoal transitions in the training set and generates the wrong latent plan. For instance, in the training data, the robot only learns to open the box after turning off the lamp, meanwhile in the *Easy* setting of generalization tasks, the robot is prompted to pick up the pen after turning off the lamp. Ours (0-human) variant still outputs a latent plan to open the box, which causes the task to fail since the box is already open.

**Real-time planning capability is robust against disturbance.** In Fig. 5, we showcase how our model reacts to unexpected human disturbance. None of these disturbances appears in the teleoperated robot demonstration data. For instance, in the cloth folding task (Fig. 5 third row), the human unfolds the towel after the robot has folded it, and the robot replans and folds the towel again. Since our whole system (including the vision-based latent planner, low-level guided policy, and robot control) is running at a speed of 17Hz, our model is able to achieve real-time re-planning capability against disturbance and manipulation errors. For more details and results, please see our demo videos.

## V. CONCLUSION

We introduce `MimicPlay`, a scalable imitation learning algorithm that exploits the complementary strengths of two data sources: cost-effective human play data and small-scale teleoperated robot demonstration data. Using human play data,

the high-level controller learns goal-conditioned latent plans by predicting future 3D human hand trajectories given the goal image. Using robot demonstration data, the low-level controller then generates the robot actions from the latent plans. With this hierarchical design, `MimicPlay` outperforms prior arts by over 50% in 14 challenging long-horizon manipulation tasks. `MimicPlay` paves the path for future research to scale up robot imitation learning with affordable human costs.

## REFERENCES

[1] Brenna D Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 57 (5):469–483, 2009.

[2] Shikhar Bahl, Abhinav Gupta, and Deepak Pathak. Human-to-robot imitation in the wild. *arXiv preprint arXiv:2207.09450*, 2022.

[3] Suneel Belkhale and Dorsa Sadigh. Plato: Predicting latent affordances through object-centric play. *arXiv preprint arXiv:2203.05630*, 2022.

[4] Aude Billard, Sylvain Calinon, Ruediger Dillmann, and Stefan Schaal. Robot programming by demonstration. In *Springer handbook of robotics*, pages 1371–1394. Springer, 2008.

[5] Christopher M Bishop. Mixture density networks. 1994.

[6] brian ichter, Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, Dmitry Kalashnikov, Sergey Levine, Yao Lu, Carolina Parada, Kanishka Rao, Pierre Sermanet, Alexander T Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Mengyuan Yan, Noah Brown, Michael Ahn, Omar Cortes, Nicolas Sievers, Clayton Tan, Sichun Xu, Diego Reyes, Jarek Rettinghouse, Jornell Quiambao, Peter Pastor, Linda Luu, Kuang-Huei Lee, Yuheng Kuang, Sally Jesmonth, Kyle Jeffrey, Rosario Jauregui Ruano, Jasmine Hsu, Keerthana Gopalakrishnan, Byron David, Andy Zeng, and Chuyuan Kelly Fu. Do as i can, not as i say: Grounding language in robotic affordances. In *6th Annual Conference on Robot Learning*, 2022.

[7] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Joseph Dabis, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Tomas Jackson, Sally Jesmonth, Nikhil Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Isabel Leal, Kuang-Huei Lee, Sergey Levine, Yao Lu, Utsav Malla, Deeksha Manjunath, Igor Mordatch, Ofir Nachum, Carolina Parada, Jodilyn Peralta, Emily Perez, Karl Pertsch, Jornell Quiambao, Kanishka Rao, Michael Ryoo, Grecia Salazar, Pannag Sanketi, Kevin Sayed, Jaspiar Singh, Sumedh Sontakke, Austin Stone, Clayton Tan, Huong Tran, Vincent Vanhoucke, Steve Vega, Quan Vuong, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Tianhe Yu, and Brianna Zitkovich. Rt-1: Robotics transformer for real-world control at scale. In *arXiv preprint arXiv:2212.06817*, 2022.

[8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[9] Sylvain Calinon, Florent D'halluin, Eric L Sauser, Darwin G Caldwell, and Aude G Billard. Learning and reproduction of gestures by imitation. *IEEE Robotics & Automation Magazine*, 17(2):44–54, 2010.

[10] Annie S Chen, Suraj Nair, and Chelsea Finn. Learning generalizable robotic reward functions from" in-the-wild" human videos. *Robotics: Science and Systems (RSS)*, 2021.

[11] Zichen Jeff Cui, Yibin Wang, Nur Muhammad, Lerrel Pinto, et al. From play to policy: Conditional behavior generation from uncurated robot data. *arXiv preprint arXiv:2210.10047*, 2022.

[12] Neha Das, Sarah Bechtle, Todor Davchev, Dinesh Jayaraman, Akshara Rai, and Franziska Meier. Model-based inverse reinforcement learning from visual demonstrations. In *Conference on Robot Learning*, pages 1930–1942. PMLR, 2021.

[13] Ashley D Edwards and Charles L Isbell. Perceptual values from observation. *arXiv preprint arXiv:1905.07861*, 2019.

[14] Peter Englert and Marc Toussaint. Learning manipulation skills from a single demonstration. *The International Journal of Robotics Research*, 37(1):137–154, 2018.

[15] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. In *Conference on robot learning*, pages 357–368. PMLR, 2017.

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

[17] Peter Florence, Lucas Manuelli, and Russ Tedrake. Self-supervised correspondence in visuomotor policy learning. *IEEE Robotics and Automation Letters*, 5(2):492–499, 2019.

[18] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jack-

son Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, et al. Ego4d: Around the world in 3,000 hours of egocentric video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18995–19012, 2022.

[19] Alex Graves and Alex Graves. Long short-term memory. *Supervised sequence labelling with recurrent neural networks*, pages 37–45, 2012.

[20] Nicklas Hansen, Zhecheng Yuan, Yanjie Ze, Tongzhou Mu, Aravind Rajeswaran, Hao Su, Huazhe Xu, and Xiaolong Wang. On pre-training for visuo-motor control: Revisiting a learning-from-scratch baseline. *arXiv preprint arXiv:2212.05749*, 2022.

[21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[22] A.J. Ijspeert, J. Nakanishi, and S. Schaal. Movement imitation with nonlinear dynamical systems in humanoid robots. In *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No.02CH37292)*, volume 2, pages 1398–1403 vol.2, 2002. doi: 10.1109/ROBOT.2002.1014739.

[23] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. Vima: General robot manipulation with multimodal prompts. *arXiv preprint arXiv: Arxiv-2210.03094*, 2022.

[24] O. Khatib. A unified approach for motion and force control of robot manipulators: The operational space formulation. *IEEE Journal on Robotics and Automation*, 3(1):43–53, 1987. doi: 10.1109/JRA.1987.1087068.

[25] Jens Kober and Jan Peters. Learning motor primitives for robotics. In *2009 IEEE International Conference on Robotics and Automation*, pages 2112–2118. IEEE, 2009.

[26] Jens Kober and Jan Peters. Imitation and reinforcement learning. *IEEE Robotics & Automation Magazine*, 17(2):55–62, 2010.

[27] Sateesh Kumar, Jonathan Zamora, Nicklas Hansen, Rishabh Jangir, and Xiaolong Wang. Graph inverse reinforcement learning from diverse videos. *Conference on Robot Learning (CoRL)*, 2022.

[28] YuXuan Liu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. Imitation from observation: Learning to imitate behaviors from raw video via context translation. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1118–1125. IEEE, 2018.

[29] Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre Sermanet. Learning latent plans from play. In *Conference on robot learning*, pages 1113–1132. PMLR, 2020.

[30] Ajay Mandlekar, Yuke Zhu, Animesh Garg, Jonathan Booher, Max Spero, Albert Tung, Julian Gao, John Emmons, Anchit Gupta, Emre Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In *Conference on Robot Learning*, pages 879–893. PMLR, 2018.

[31] Ajay Mandlekar, Fabio Ramos, Byron Boots, Silvio Savarese, Li Fei-Fei, Animesh Garg, and Dieter Fox. Iris: Implicit reinforcement without interaction at scale for learning control from offline robot manipulation data. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4414–4420. IEEE, 2020.

[32] Ajay Mandlekar, Danfei Xu, Roberto Martín-Martín, Silvio Savarese, and Li Fei-Fei. Learning to generalize across long-horizon tasks from human demonstrations. *arXiv preprint arXiv:2003.06085*, 2020.

[33] Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. In *5th Annual Conference on Robot Learning*, 2021. URL https://openreview.net/forum?id=JrsfBJtDFdI.

[34] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. In *6th Annual Conference on Robot Learning*, 2022. URL https://openreview.net/forum?id=tGbpgz6yOrI.

[35] Alexandros Paraschos, Christian Daniel, Jan R Peters, and Gerhard Neumann. Probabilistic movement primitives. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL https://proceedings.neurips.cc/paper/2013/file/e53a0a2978c28872a4505bdb51db06dc-Paper.pdf.

[36] Alexandros Paraschos, Christian Daniel, Jan Peters, and Gerhard Neumann. Using probabilistic movement primitives in robotics. *Autonomous Robots*, 42(3):529–551, 2018.

[37] Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. *Advances in neural information processing systems*, 1, 1988.

[38] Erick Rosete-Beas, Oier Mees, Gabriel Kalweit, Joschka Boedecker, and Wolfram Burgard. Latent plans for task-agnostic offline reinforcement learning. *arXiv preprint arXiv:2209.08959*, 2022.

[39] Erick Rosete-Beas, Oier Mees, Gabriel Kalweit, Joschka Boedecker, and Wolfram Burgard. Latent plans for task agnostic offline reinforcement learning. 2022.

[40] Stefan Schaal. Is imitation learning the route to humanoid robots? *Trends in cognitive sciences*, 3(6):233–242, 1999.

[41] Stefan Schaal. Dynamic movement primitives-a framework for motor control in humans and humanoid robotics. In *Adaptive motion of animals and machines*, pages 261–280. Springer, 2006.

[42] Karl Schmeckpeper, Annie Xie, Oleh Rybkin, Stephen Tian, Kostas Daniilidis, Sergey Levine, and Chelsea Finn. Learning predictive models from observation and interaction. In *European Conference on Computer Vision*, pages 708–725. Springer, 2020.

[43] Karl Schmeckpeper, Oleh Rybkin, Kostas Daniilidis, Sergey Levine, and Chelsea Finn. Reinforcement learning with videos: Combining offline observations with interaction. In Jens Kober, Fabio Ramos, and Claire Tomlin, editors, *Proceedings of the 2020 Conference on Robot Learning*, volume 155 of *Proceedings of Machine Learning Research*, pages 339–354. PMLR, 16–18 Nov 2021. URL https://proceedings.mlr.press/v155/schmeckpeper21a.html.

[44] Nur Muhammad Mahi Shafiullah, Zichen Jeff Cui, Ariuntuya Altanzaya, and Lerrel Pinto. Behavior transformers: Cloning $k$ modes with one stone. *arXiv preprint arXiv:2206.11251*, 2022.

[45] Dandan Shan, Jiaqi Geng, Michelle Shu, and David Fouhey. Understanding human hands in contact at internet scale. In *CVPR*, 2020.

[46] Lin Shao, Toki Migimatsu, Qiang Zhang, Karen Yang, and Jeannette Bohg. Concept2robot: Learning manipulation concepts from instructions and human demonstrations. *The International Journal of Robotics Research*, 40(12-14):1419–1434, 2021.

[47] Pratyusha Sharma, Deepak Pathak, and Abhinav Gupta. Third-person visual imitation learning via decoupled hierarchical controller. *Advances in Neural Information Processing Systems*, 32, 2019.

[48] Kenneth Shaw, Shikhar Bahl, and Deepak Pathak. Videodex: Learning dexterity from internet videos. *CoRL*, 2022.

[49] Kyriacos Shiarlis, Markus Wulfmeier, Sasha Salter, Shimon Whiteson, and Ingmar Posner. Taco: Learning task decomposition via temporal alignment for control. In *International Conference on Machine Learning*, pages 4654–4663. PMLR, 2018.

[50] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipulation. In *Proceedings of the 5th Conference on Robot Learning (CoRL)*, 2021.

[51] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Perceiver-actor: A multi-task transformer for robotic manipulation. In *Proceedings of the 6th Conference on Robot Learning (CoRL)*, 2022.

[52] Maximilian Sieb, Zhou Xian, Audrey Huang, Oliver Kroemer, and Katerina Fragkiadaki. Graph-structured visual imitation. In *Conference on Robot Learning*, pages 979–989. PMLR, 2020.

[53] Laura Smith, Nikita Dhawan, Marvin Zhang, Pieter Abbeel, and Sergey Levine. Avid: Learning multi-stage tasks via pixel-level translation of human videos. *arXiv preprint arXiv:1912.04443*, 2019.

[54] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008. URL http://www.jmlr.org/papers/v9/vandermaaten08a.html.

[55] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[56] Chen Wang, Rui Wang, Ajay Mandlekar, Li Fei-Fei, Silvio Savarese, and Danfei Xu. Generalization through hand-eye coordination: An action space for learning spatially-invariant visuomotor control. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8913–8920. IEEE, 2021.

[57] Tete Xiao, Ilija Radosavovic, Trevor Darrell, and Jitendra Malik. Masked visual pre-training for motor control. *arXiv preprint arXiv:2203.06173*, 2022.

[58] Haoyu Xiong, Quanzhou Li, Yun-Chun Chen, Homanga Bharadhwaj, Samarth Sinha, and Animesh 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.

[59] Danfei Xu, Suraj Nair, Yuke Zhu, Julian Gao, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Neural task programming: Learning to generalize across hierarchical tasks. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3795–3802. IEEE, 2018.

[60] Kevin Zakka, Andy Zeng, Pete Florence, Jonathan Tompson, Jeannette Bohg, and Debidatta Dwibedi. Xirl: Cross-embodiment inverse reinforcement learning. In *Conference on Robot Learning*, pages 537–546. PMLR, 2022.

[61] Tianhao Zhang, Zoe McCarthy, Owen Jow, Dennis Lee, Xi Chen, Ken Goldberg, and Pieter Abbeel. Deep imitation learning for complex manipulation tasks from virtual reality teleoperation. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5628–5635. IEEE, 2018.

[62] Yifeng Zhu, Abhishek Joshi, Peter Stone, and Yuke Zhu. VIOLA: Object-centric imitation learning for vision-based robot manipulation. In *6th Annual Conference on Robot Learning*, 2022. URL https://openreview.net/forum?id=L8hCfhPbFho.

## VI. APPENDIX

### A. Details of system setups

We illustrate the system designs for the data collection in Fig. 6. The human play data is collected by having a human operator directly interact with the environment with one of its hands (Fig. 6(a)). The left and right cameras record the video at the speed of 100 frames per second. During the collection process of human play data, no specific task goal is given and the human operator freely interacts with the scene for interesting behaviors based on its curiosity. For each scene in our experiments, we collect 10 minutes of human play data.

The robot teleportation demonstration is collected with a phone teleoperation system RoboTurk [30] (Fig. 6(b)). The left, right, and end-effector wrist cameras record the video at the speed of 20 frames per second, which is aligned with the control speed of the robot arm (20Hz). Each sequence of robot demonstration has a pre-defined task goal. During the data collection, the human demonstrator completes the assigned sub-goals one by one and finally solves the whole task. For each training task in our experiments, we collect 20 demonstrations. In the Kitchen environment, we collect 40 demonstrations for each task to figure out which approach is more sample inefficiency.
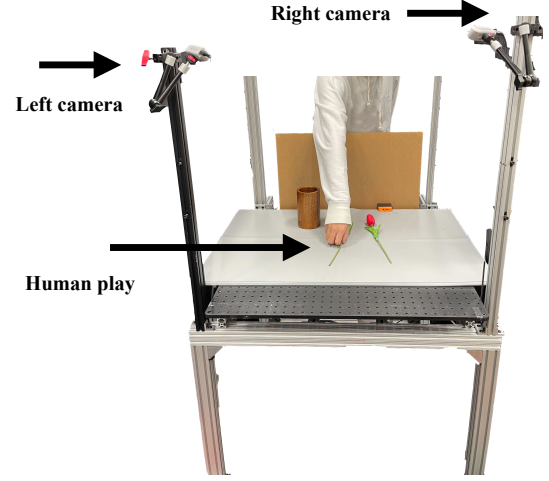
We visualize the collected human play data and robot demonstration data in Tab. 8. For the human play data, we use an off-the-shelf hand detector [45] to localize the hand's 2D location on the left and right image frame, which are visualized as red bounding boxes in Tab 8. For the robot demonstration data, we directly project the 3D location of the robot end-effector to the left and right image frames, which are visualized as blue bounding boxes in Tab 8.
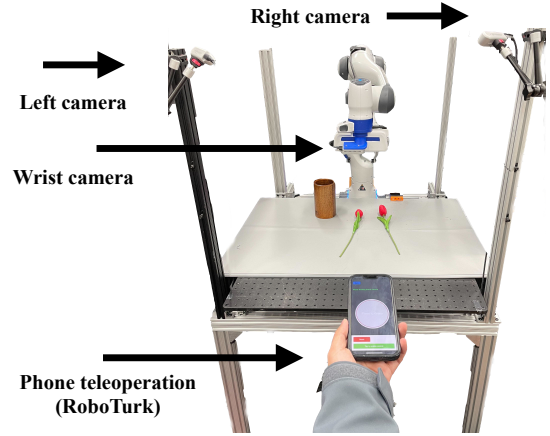
### B. Details of the task designs

The definition of our long-horizon tasks is listed in Tab. VII. For each task, the initial state and subgoals are pre-defined. The whole task is completed if and only if all subgoals are completed in the correct order.

### C. Analysis of the visual gap between human and robot data

As is introduced in the method Sec. III-A, to minimize the visual gap between human play data and robot demonstration data, we use a KL divergence loss over the feature embeddings outputted by the visual encoders. In Fig. 7, we use t-SNE to process and visualize the learned feature embeddings generated by Ours and the model variant Ours (w/o KL) on the 2D distribution plots. To better visualize the distribution overlap, we use slashes to highlight the overlap area in both plots. We observe that our approach with KL loss has a 23% larger overlap between the human data and the robot data compared to Ours (w/o KL). This result showcases the effectiveness of our KL divergence loss and supports the result in Tab. II (Ours (w/o KL) is inferior to Ours in task success rate).



**(a) Human play data collection**



**(b) Robot demonstration data collection**

Fig. 6: System setups for the data collection. (**a**) Human play data collection. A human operator directly interacts with the scene with one of its hand and perform interesting behaviors based on its curiosity without a specific task goal. (**b**) Robot demonstration data collection. A human demonstrator uses a phone teleoperation system to control the 6 DoF robot end-effector. The gripper of the robot is controlled by pressing a button on the phone interface.

### D. Training hyperparameters

We list the hyperparameters for training the models in Tab. IV for the latent planner $\mathcal{P}$ and Tab. V for the robot policy $\pi$. The hyperparameters that are named starting with GMM are related to the MLP-based GMM model. The hyperparameters that are named starting with GPT are related to the transformer architecture. We also list the hyperparameters for the baseline GC-BC (BC-trans) in Tab. VI.

**(a) Distribution overlap of Ours (w/o KL)**



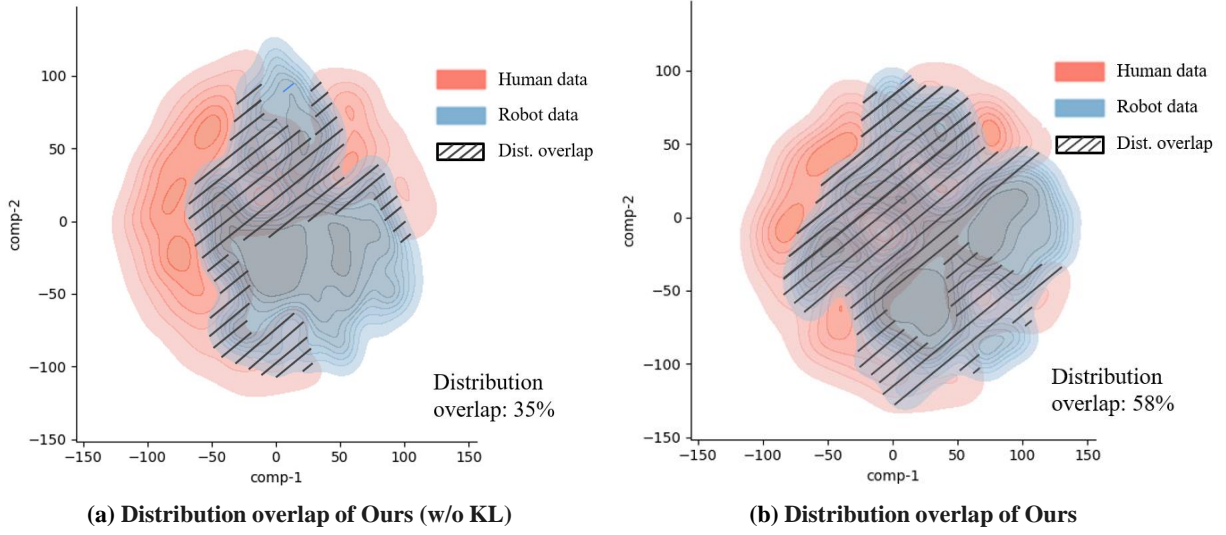**(b) Distribution overlap of Ours**

Fig. 7: t-SNE visualization of the generated feature embeddings by taking human data and robot data as inputs. The slashes refer to the overlap region of two data distributions. (**a**) Feature visualization results of our method without using KL divergence loss. (**b**) Feature visualization results of our method with KL divergence loss. Our approach covers 23% more area than the baseline.

| Hyperparameter | Default |
|---|---|
| Batch Size | 16 |
| Learning Rate (LR) | 1e-4 |
| Num Epoch | 1000 |
| LR Decay | None |
| KL Weights $\lambda$ | 1000 |
| MLP Dims | [400, 400] |
| Image Encoder - Left View | ResNet-18 |
| Image Encoder - Right View | ResNet-18 |
| Image Feature Dim | 64 |
| GMM Num Modes | 5 |
| GMM Min Std | 0.0001 |
| GMM Std Activation | Softplus |

TABLE IV: Hyperparameters - Ours (Latent Planner $\mathcal{P}$)

| Hyperparameter | Default |
|---|---|
| Batch Size | 16 |
| Learning rate (LR) | 1e-4 |
| Num Epoch | 1000 |
| Train Seq Length | 10 |
| LR Decay Factor | 0.1 |
| LR Decay Epoch | [300, 600] |
| MLP Dims | [400, 400] |
| Image Encoder - Wrist View | ResNet-18 |
| Image Feature Dim | 64 |
| GMM Num Modes | 5 |
| GMM Min Std | 0.01 |
| GMM Std Activation | Softplus |
| GPT Block Size | 10 |
| GPT Num Head | 4 |
| GPT Num Layer | 4 |
| GPT Embed Size | 656 |
| GPT Dropout Rate | 0.1 |
| GPT MLP Dims | [656, 128] |

TABLE V: Hyperparameters - Ours (Robot Policy $\pi$)

| Hyperparameter | Default |
|---|---|
| Batch Size | 16 |
| Learning rate (LR) | 1e-4 |
| Num Epoch | 1000 |
| Train Seq Length | 10 |
| LR Decay Factor | 0.1 |
| LR Decay Epoch | [300, 600] |
| MLP Dims | [400, 400] |
| Image Encoder - Wrist View | ResNet-18 |
| Image Encoder - Left View | ResNet-18 |
| Image Encoder - Right View | ResNet-18 |
| Image Feature Dim | 64 |
| GMM Num Modes | 5 |
| GMM Min Std | 0.01 |
| GMM Std Activation | Softplus |
| GPT Block Size | 10 |
| GPT Num Head | 4 |
| GPT Num Layer | 4 |
| GPT Embed Size | 656 |
| GPT Dropout Rate | 0.1 |
| GPT MLP Dims | [656, 128] |

TABLE VI: Hyperparameters - GC-BC (BC-trans)

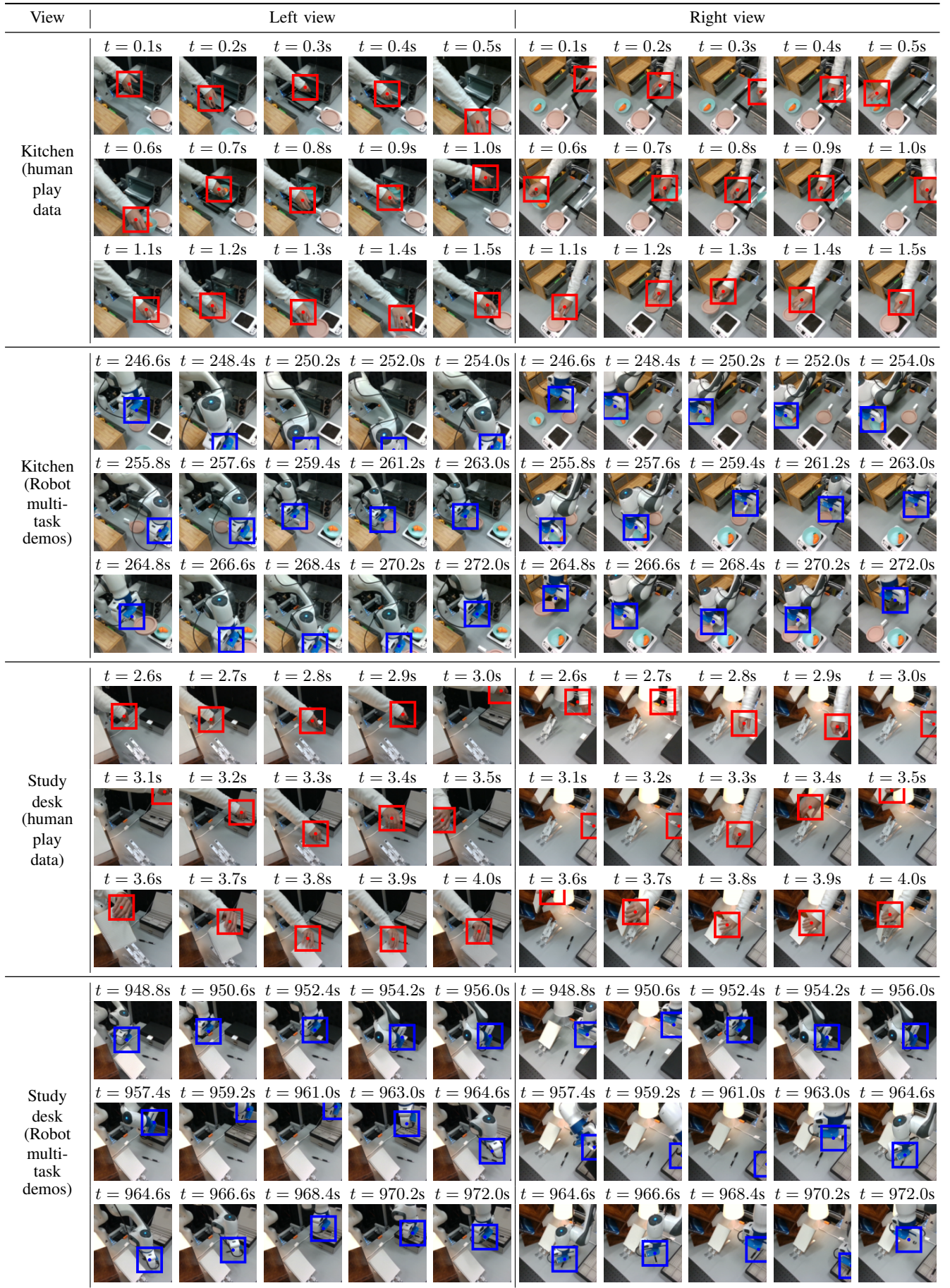| | | |
|---|---|---|
| Kitchen | Task-1 | • Initial state: A drawer is placed on the left side of the table. The drawer is not fully open and contains pumpkin and lettuce. A closed microwave oven is placed on the right side of the desktop. A bowl and a stove are placed on the lower edge of the tabletop. There is a carrot inside the bowl. A pan is placed on top of the stove.<br>• Subgoals: a) Open the microwave oven door. b) Pull out the microwave oven tray. c) Pick up the bowl. d) Place the bowl on the microwave tray. |
| | Task-2 | • Initial state: same as Kitchen Task-1.<br>• Subgoals: a) Open the drawer. b) Pick up the carrot. c) Put the carrots in the drawer. |
| | Task-3 | • Initial state: same as Kitchen Task-1.<br>• Subgoals: a) Pick up the pan. b) Place the pan on the table. c) Pick up the bowl. d) Place the bowl on the stove. |
| Study desk | Task-1 | • Initial state: The book is on the rack. The lamp is on. The box is opened and closed in a random state. The pen is located either in the center of the table or in the box.<br>• Subgoals: a) Turn off the lamp. b) Pick up the book. c) Place the book on the shelf position. |
| | Task-2 | • Initial state: The location of the book is either on the shelf or on the rack. The lamp is off. The box is closed. The pen is in the center of the table.<br>• Subgoals: a) Turn on the lamp. b) Open the box. c) Pick up the pen. d) Put it in the box. |
| | Task-3 | • Initial state: The book is on the rack. The state of the lamp is random. The box is closed. The pen is in the center of the table.<br>• Subgoal a) Open the box. b) Pick up the pen. c) Place the pen in the box. d) Pick up the book. e) Place the book on the shelf. |
| | Task-4 | • Initial state: The location of the book is either on the shelf or on the rack. The lamp is on. The box is closed. The pen is located either in the center of the table or in the box.<br>• Subgoals: a) Open the box. b) Turn off the lamp. |
| | Easy | • Initial state: The location of the book is either on the shelf or on the rack. The lamp is off. The box is closed. The pen is located either in the center of the table or in the box.<br>• Subgoals: a) Turn on the lamp. b) Open the box. c) Turn off the lamp. |
| | Medium | • Initial state: The location of the book is either on the shelf or on the rack. The lamp is on. The box is closed. The pen is in the center of the table.<br>• Subgoals: a) Open the box. b) Turn off the lamp. c) Pick up the pen. d) Place the pen in the box. |
| | Hard | • Initial state: The book is on the shelf. The lamp is on. The box is closed. The pen is located either in the center of the table or in the box.<br>• Subgoals: a) Turn off the lamp. b) Open the box. c) Pick up the book. d) Place the book on the shelf. |
| Flower | - | • Initial state: Two flowers and a vase are placed on the table. The vase will randomly be placed on the top left or top right corner of the table.<br>• Subgoals: a) Picking up a flower. b) Insert the flower into the vase. c) Pick up the other flower. d) Insert the flower into the vase. |
| Whiteboard | - | • Initial state: A whiteboard and board eraser are placed on the table. The board eraser is placed on the left side of the whiteboard.<br>• Subgoals: a) Pick up the board eraser. b) Moves over the curve line. c) Erase the curve line. d) Return the eraser to the original location. |
| Sandwich | - | • Initial state: A circular ingredient selector is placed in the upper right corner of the table. Half of the circle holds ingredients for a sandwich (bread, lettuce, sliced tomato) and half holds ingredients for a cheeseburger (bread, cheese, burger patty). A white plate is placed in the lower left corner of the table.<br>• Subgoals for a sandwich: a) Rotate the ingredient selector to the right position. Pick up a piece of bread from it and place it on the plate. b) Rotate the ingredient selector to the correct position. Pick up the lettuce and place it on top of the bread. c) Rotate the ingredient selector to the right position. Pick up the sliced tomato and place it on top of the lettuce. d) Rotate the ingredient selector to the right position. Pick up another piece of bread and place it on top of the tomato. |
| Cloth | - | • Initial state: An unfolded brown cloth is randomly placed on the table.<br>• Subgoals: a) The robot folds the cloth in half once to become 1/2 of its original size. b) The robot folds the cloth once more to become 1/4 of its original size. |

TABLE VII: Task definition.

Fig. 8: Dataset visualization.