# Human-Agent Joint Learning for Efficient Robot Manipulation Skill Acquisition

**Shengcheng Luo**[1*], **Quanquan Peng**[1*], **Jun Lv**[1*], **Kaiwen Hong**[2],
**Katherine Rose Driggs-Campbell**[2], **Cewu Lu**[1], **Yong-Lu Li**[1†]
Shanghai Jiao Tong University[1]    University of Illinois Urbana-Champaign[2]

**Abstract:** Employing a teleoperation system for gathering demonstrations offers the potential for more efficient learning of robot manipulation. However, teleoperating a robot arm equipped with a dexterous hand or gripper, via a teleoperation system poses significant challenges due to its high dimensionality, complex motions, and differences in physiological structure. In this study, we introduce a novel system for joint learning between human operators and robots, that enables human operators to share control of a robot end-effector with a learned assistive agent, facilitating simultaneous human demonstration collection and robot manipulation teaching. In this setup, as data accumulates, the assistive agent gradually learns. Consequently, less human effort and attention are required, enhancing the efficiency of the data collection process. It also allows the human operator to adjust the control ratio to achieve a trade-off between manual and automated control. We conducted experiments in both simulated environments and physical real-world settings. Through user studies and quantitative evaluations, it is evident that the proposed system could enhance data collection efficiency and reduce the need for human adaptation while ensuring the collected data is of sufficient quality for downstream tasks. Videos are available at https://mvig-rhos.com/joint_learning.

**Keywords:** Robotic Teleoperation, Robot Manipulation, Imitation Learning

## 1 Introduction

The long-term vision in the field of robot learning has been to enable robots to perform diverse tasks at a human level in the physical world. Recently, significant progress has been made toward this goal by learning robot manipulation policies from demonstrations. Previous studies have utilized teleoperation systems [1, 2, 3, 4, 5, 6] to collect human demonstrations, and learning-based policies [7, 8, 9] have been formulated using the gathered data. Despite the notable advancements, several challenges still need to be addressed. For example, in vision-based teleoperation systems, even with state-of-the-art 3D hand pose estimation algorithms [10, 11, 12, 13], errors persist that significantly affect the teleoperation. Additionally, discrepancies between the structures of human hands and robot end-effectors, along with the lack of haptic feedback during contact-rich manipulation, also pose challenges. As a result, current teleoperation systems require human operators to practice extensively to adapt to these differences and gather the necessary data. This means that many human adaptations are essential. Furthermore, to meet the data requirements of robotic systems, humans need to collect large amounts of data, making this process very burdensome.

Naturally, a question was raised: *in data-collection, how to make human adaptation less or even free while keeping the data quality?* Here, we aim to address this question and argue that human-agent joint learning can help. That said, an effective and efficient teleoperation system should be designed to preferentially capture the operator's intentions for directing a robot end effector and pose the *main frame*, while concurrently enabling an autonomous agent to help us ensure motion stability and *interpolate* the details. To this end, we propose a framework that achieves shared control between the

---

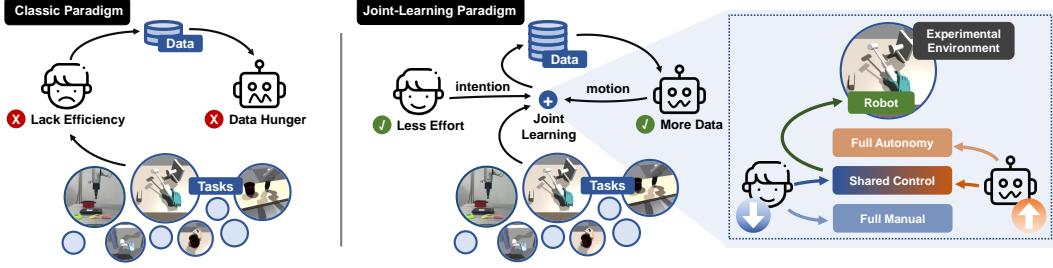* denotes equal contribution, †denotes corresponding author.

Figure 1: Traditional frameworks typically separate human and agent training, requiring operators first to learn the task environment before data collection. This often leads to inefficiencies due to delayed and insufficient data gathering. In our framework, we integrate human and agent training from the start in a joint learning model. This enables simultaneous development and adapts the agents to human operation more effectively, enhancing overall efficiency and promoting better collaboration between humans and machines allowing for human effortless adaptation data collection.

human and a learned assistive agent. As shown in Fig. 1, our *human effortless adaptation* framework seeks to balance human needs and system performance, aiming to directly enhance the efficiency and quality of data collection by reducing the time and effort in human adaptation.

Given our human-agent joint learning approach, we allow the data acquisition agent to grow and learn along with the human operator, reducing the learning burden of the human operator throughout the data collection task. Inspired by shared autonomy [14, 15, 16], we introduce a novel teleoperation system that enables collaboration between humans and learning-based agents to control a robot jointly. In particular, our proposed system provides the flexibility to adjust a "control ratio" between the human operator and a learning-based agent. A lower control ratio, in the beginning, signifies a greater emphasis on humans to teach the agent the finer-grained knowledge under the structure of human intention and principal actions. As the agent's learning improves, a higher ratio later indicates increased autonomy from the learned agent to replace the human effort to "inpaint" the whole process given only human intention and principal actions. Additionally, once a sufficient amount of data is gathered, we have the option to transition the shared autonomy agent to full autonomy by reducing the human control ratio to zero.

We implement our framework as a joint learning system, which facilitates a human operator's intention and an assistive agent's execution of teleoperation tasks. For human operators, their inputs consist of intuitive control actions based on visual feedback and past cognitive experience. For the agent, their inputs are derived from a combination of sensors, data streams, and possibly pre-processed information. This can include visual data from cameras [1], tactile data from force sensors [3], and any additional context provided by the system [6], such as object recognition or environmental mapping. The agent processes these inputs using algorithms designed to interpret the task requirements and generate actions to assist the human operator, utilizing a diffusion model [8, 17, 18, 19] as the backbone enables us to adjust the control ratio by modifying the step number of the forward and reverse processes [20], providing a customized and adaptive approach to teleoperation tasks.

We conducted experiments in six different *simulation* environments using two types of end-effectors: a dexterous hand and a gripper. Additionally, we performed experiments on three *real-world* tasks to validate our findings. Evaluation results indicate that our proposed system significantly enhances data collection efficiency, increasing the collection success rate by *30%* and nearly *doubling* the collection speed. Additionally, data collected in shared autonomy mode is as effective for downstream tasks and models as data collected directly from the teleoperation system, demonstrating comparable validity. Our main contributions are summarized as follows:

- We study how to reduce human adaptation in teleoperation data collection and propose a human-agent joint learning paradigm.

- We build a system that fosters concurrent development between the human operator and assistive agent, which not only streamlines the learning process but also expedites the robot's ability to perform robot manipulations autonomously.

- Conducting experiments to demonstrate the efficiency and effectiveness of our proposed system. Our system achieved significant performance improvements, including a *30%* increase in data collection success rate and nearly *double* the collection speed. We also deployed our system in a *real-world* environment and achieved significant results.

## 2 Related Works

**Teleoperation System.** Data has always been a crucial foundation, and robots are no exception. Teleoperation serves as a significant source for collecting robot data [7, 21, 22, 23, 24, 25]. Some works achieve teleoperation through wearable devices [1, 2, 3, 4, 26], and vision-based teleoperation systems offer a low-cost and easily developed alternative [5, 6, 27, 28]. For instance, Li et al. [28] utilizes neural networks for markerless vision-based teleoperation of dexterous robotic hands from depth images. Handa et al. [5] set up a vision-based teleoperation system to control the Allegro Hand, accomplishing various contact-rich manipulation tasks in the real world. Recently, Qin et al. [6] introduced AnyTeleop, a unified teleoperation system designed to accommodate various arms, hands, realities, and camera setups within a singular framework. In this paper, we introduce a joint learning paradigm to assist teleoperation by sharing control between the human operator and a learning-based agent, aiming to improve the efficiency of the teleoperation process.

**Human Robot Cooperation.** Collecting fine-grained human demonstration data for robotic manipulation is an effective but labor-intensive and time-consuming way to enable robots to complete a wide range of tasks [29, 30]. Previous work uses shared autonomy to assist people with disability in performing tasks by arbitrating human inputs and robot actions [31]. Many of the shared autonomy algorithms aim to estimate human intents from a set of pre-defined goals [32, 33, 34, 35] or by mapping low-dimension control input to high-dimension robot actions [31, 36]. In this work, we introduce a system that enables shared control between the human and assistive agent to facilitate the process of data collection and robot learning.

## 3 Technical Approach

The proposed system enables human operators to control the robot using a teleoperation system to gather training data (Sec. 3.1). Subsequently, utilizing the collected data, we train an agent (Sec. 3.2) to establish shared control between the human operator and the learned agent, thereby enhancing the efficiency of the data collection process (Sec. 3.3). Similar to the concept of "bootstrapping", as more data accumulates, our system raises the control ratio of the learned agent, thereby reducing the effort required from human operators. This, in turn, enables us to collect even more data and continue improving the system iteratively. Moreover, we offer the option to transition the shared control agent to full autonomy once sufficient data is acquired (Sec. 3.4).

### 3.1 Teleoperation System.

Our pipeline initially captures the raw sensory signal $\mathcal{I}$. Human hand pose $\mathcal{P}^h \in \mathbb{R}^{20 \times 3}$ can be obtained from the captured signal using off-the-shelf 3D hand pose estimation [10, 11, 13]. The pose $\mathcal{P}^h$ consists of the positions of 20 keypoints. Then, employing an inverse kinematic function $f_{IK}$, we compute the action of the robot $a \in \mathbb{R}^m$: $a = f_{IK}(\mathcal{P}^h_t, \mathcal{P}^h_{t+1})$, where it is calculated upon the change in the hand pose. Given this teleoperation system, the human operator will move the hand to produce a sequence of hand poses $\{\mathcal{P}^h_i\}_{i=0}^T$ to teleoperate the robot with an action sequence $\{a_i\}_{i=0}^T$ to achieve the task $\mathcal{T}$. The trajectory $\{(s_i, a_i)\}_{i=0}^T$ is the collected human demonstration data, where $s \in \mathbb{R}^n$ (here $n = 18$) is the robot state, could be used for downstream tasks.

### 3.2 Diffusion-Model-Based Assistive Agent.

After getting the data, we train a diffusion-model-based assistive agent to learn how to assist the human in collecting data in a shared control manner. We follow the Denoising Diffusion Probabilistic Model (DDPM) [18] training paradigm to construct the diffusion-model-based assist agent. The

*forward process* of the Diffusion Model can be regarded as adding Gaussian noise to the data $x^0$ according to a variance schedule $\beta_{1:K}$ by

$$x_k = \sqrt{\alpha_k}x_{k-1} + \sqrt{1-\alpha_k}\epsilon, \tag{1}$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \alpha_k = 1 - \beta_k$. DDPM models the output generation as a denoising process (Stochastic Langevin Dynamics). A line of works [19, 17, 8, 37] use diffusion model to generate the action for agents: given $x^K$ sampled from Gaussian noise $\mathcal{N}(\mathbf{0}, \mathbf{I})$, it utilizes a parameterized diffusion process to model how $x^K$ is denoised in order to get noise-free action $x^0$ by

$$p_\theta(x^0) = \int p(x^K) \prod_{k=1}^{K} p_\theta(x^{k-1}|x^k)\mathrm{d}x^{1:K}, \tag{2}$$

where $p_\theta(x^{k-1}|x^k) = \mathcal{N}(\mu_\theta(x^k, k), \Sigma(x^k, k))$ is usually referred as *reverse process*. Luo [38] shows that $p_\theta(x^{t-1}|x^k)$ becomes tractable when conditioned on $x_0$ and Eq. 2 can be reformulated as minimizing the error in the noise prediction. Ho et al. [18] simplify the training loss function as

$$\mathcal{L} := \mathbb{E}_{k,\boldsymbol{x}_0,\boldsymbol{\epsilon}\sim\mathcal{N}(\mathbf{0},\boldsymbol{I})}\left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\boldsymbol{x}_k(\boldsymbol{x}_0, \boldsymbol{\epsilon}), k)\|_2^2\right], \tag{3}$$

where step $k$ is sampled uniformly as $k \in [1, K]$, $\boldsymbol{\epsilon}_\theta$ is the noise prediction model. During the inference phase, we can generate $x_0$ by recursively sample $\boldsymbol{z} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$:

$$x_{k-1} = \mu_\theta(x_k, k) + \sigma_k \boldsymbol{z}. \tag{4}$$

Similar to [8, 20], with the collected trajectory $\{(s_i, a_i)\}_{i=0}^{T}$, we aim to train an agent to imitate the trajectory, accomplishing a specific task $\mathcal{T}$. Therefore, we utilize DDPM to capture the conditional distribution of $p(a|s)$ and the training loss in Eq. 3 shall be modified as

$$\mathcal{L} := \mathbb{E}_{k,(s_i,a_i),\boldsymbol{\epsilon}\sim\mathcal{N}(\mathbf{0},\boldsymbol{I})}\left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(a_i + \boldsymbol{\epsilon}, s_i, k)\|_2^2\right]. \tag{5}$$

At an abstract level, the diffusion-model-based assist agent, noted as $f(\cdot|\cdot)$, is provided with the state $s$, denoising step number $k$, and a noise action $a^k$, which could be an imperfect action gathered from the teleoperation system or sampled from a Gaussian distribution, to predict the desired action

$$a = f(a^k|s, k). \tag{6}$$

Note that during the experiment, we found that adding slight noise to $s_i$ during training will give a better result. We assume doing so will augment the dataset.

### 3.3  Data Collection with Shared Control

During data collection, the proposed system offers the option to control the robot in a shared control mode rather than directly applying the collected action $a^h$ from the teleoperation system. The classical shared autonomy method is achieved through the equation [32]:

$$a^s = \gamma a^h + (1 - \gamma)a^r, \tag{7}$$

where $a^r$ is generated by the learned agent. Considering that the agent operates as a diffusion policy (Fig. 2), we blend the action from the human with the forward and reverse processes. Given action $a^h$, a forward process diffuses the action as follows: $a^k = a^h + \epsilon^k$. Subsequently, a reverse process denoises the action $a^k$:



Figure 2: To achieve shared control between the human and agent, we blend the action from the human operator $a^h$ using the forward and reverse process. The parameter $\gamma$ governs the control ratio, where a lower $\gamma$ results in the action better aligning with the human operator's intention. In contrast, a higher $\gamma$ allows the learned agent to exert more influence over the blended action.

$$a^s = f(a^k|s, k). \tag{8}$$

By applying action $a$, the control of the robot is shared between the human and the diffusion-model-based assistive agent. We can adjust the control ratio $\gamma = k/K$ between the human operator and
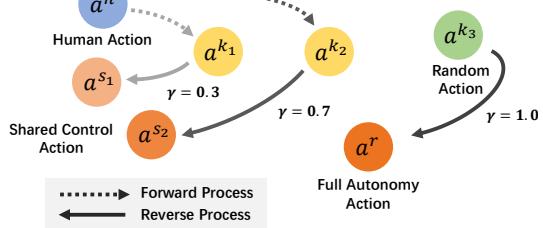
the diffusion-model-based assistive agent by varying $k$. When $\gamma = 0$, the action $a^s$ represents the teleoperation action $a^h$, which is the dexterous robot directly controlled by a human operator. As $\gamma$ approaches 1.0, the action $a^s$ transitions to full autonomy $a^r$. A higher $\gamma$ value indicates a higher level of autonomy, allowing the learning-based agent more control rights to stabilize and direct the dexterous hand. This leads to a reduced human workload during the data collection process.

### 3.4 Intergrating Data Collection and Manipulation Learning.

We outline the overall process in Algo. 1. The assistive agent is trained in four steps as follows:

*Step 1.* Initially, we collect a dataset for pre-training agent $f$ under full manual control by human operators, *i.e.*, with the control ratio $\gamma = 0$.

*Step 2.* Given the initial dataset, we train a less capable assistive agent to aid in further data collection. The training process has been formulated in Eq. 5 and Eq. 6, where a neural network $\epsilon_\theta$ is trained to predict noise $\epsilon$ out of the noisy action $a^k$.

*Step 3.* The trained agent assists in a second data collection round, aiming for higher efficiency and success. We refine the agent using data from both rounds to enhance its performance. This cycle repeats until the agent achieves full autonomy and the required data volume is collected.

---

**Algorithm 1** Overall Process

---

**Require:** The human operator $\mathcal{H}$;
**Ensure:** The collected dataset $\mathcal{D}$; the learned agent $f$; control ratio $\gamma$;
 1: Initialization: $\mathcal{D} \leftarrow \emptyset, \gamma \leftarrow 0$;
 2: **while** $|\mathcal{D}|$ is small **do**                    ▷ not enough data is collected
 3:     $\mathcal{H}$ collects data $d$ under $f$'s help;    ▷ with proper control ratio $\gamma$, details shown in Fig. 2
 4:     **if** $d$ is valid **then**
 5:         $\mathcal{D} \leftarrow \{d\} \cup \mathcal{D}$;
 6:     Finetune $f$ with $\mathcal{D}$;
 7:     Raise the agent's control ratio $\gamma$;
 8: **return** $\mathcal{D}$ and $f$;

---

Currently, the adjustment of the control ratio $\gamma$ is guided by intuitive assessments and careful monitoring of success rates. When an agent, trained with recently gathered data, exhibits a marked improvement in success rate relative to its prior performance, this enhancement is taken as a cue to fine-tune $\gamma$, thus refining control effectiveness. The assessment of the agent's performance, followed by adjustments to $\gamma$, relies on empirical data and the agent's performance in practice. This method ensures that decisions regarding the adjustment of $\gamma$ are firmly rooted in data, thus harmonizing the control strategy with the observed outcomes. It is posited that an agent reaches full autonomy when it achieves a success rate that aligns with or surpasses current state-of-the-art benchmarks.

## 4 Experiments

**Tasks.** We adopt six multi-stage manipulation tasks (Fig. 3). *Pick-and-Place* aims at picking an object on the table and placing it into a container. *Articulated-Manipulation*'s objective for the dexterous hand is to grasp and unscrew a door handle to open it, while for the gripper, it is to grab a drawer handle and pull the drawer open. *Push-cube* requires the robot to push the cube to the target position. *Tool-Use* aims at picking a hammer and using it to drive a nail into a board.

**Efficiency of Data Collection.** Our proposed system leverages shared control between human operators and learned agents to enhance the efficiency of data collection. To learn how the assistant agent could improve the data collection process, we conducted a user study.

In the user study, 10 human operators participate, collecting data under two modes: one where control is shared between the operator and the learned agent (*w/ Ours*), and the other where control
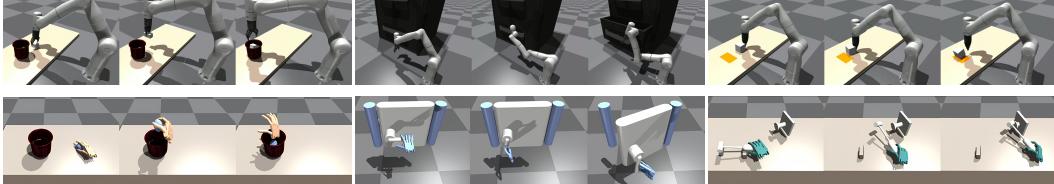
Figure 3: Overview of six task settings and their task flow for Pick-and-Place *(left)*, Articulated-Manipulation *(middle)*, Gripper-Push *(upper-right)* and Dexterous-Tool-Use *(bottom-right)*.

| | | Pick-and-Place | | | Door-Open | | | Tool-Use | | |
| | | Success Rate ↑ | Horizon Length ↓ | Collection Speed ↑ | Success Rate ↑ | Horizon Length ↓ | Collection Speed ↑ | Success Rate ↑ | Horizon Length ↓ | Collection Speed ↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| *Group1* | *w/ Ours* | **86.96** | **219.01** | **320** | **87.11** | **142.29** | **460** | **66.50** | **232.17** | **200** |
| | *w/o Ours* | 51.53 | 378.49 | 176 | 62.49 | 258.27 | 252 | 42.38 | 487.95 | 129 |
| *Group2* | *w/ Ours* | **94.06** | **214.16** | **324** | **80.29** | **134.16** | **424** | **55.55** | **275.71** | **172** |
| | *w/o Ours* | 45.42 | 471.48 | 120 | 53.45 | 317.21 | 176 | 34.47 | 511.03 | 124 |

Table 1: User studies on three dexterous hand tasks.

is directly by the operator alone (*w/o Ours*). Each participant is instructed to collect as much data as possible within 3 minutes under two different modes for 3 dexterous hand tasks. Three metrics are evaluated: **Success Rate** (Percent) indicates the percentage of attempts where data collection was successful. **Horizon Length** (Steps per Sample) measures the length of each collected trajectory, with a lower horizon length indicating smoother data collection. **Collection Speed** (Samples per Hour) refers to the number of successful trajectories that can be collected in one hour.

In Tab. 1, by sharing control between humans and learned agents, our system shows improvements in both success rate and collection speed, while the average horizon length of the collected trajectories is reduced. This suggests that our system enhances the efficiency of data collection by facilitating a process that is easier to succeed, faster, and more fluid in terms of trajectory smoothness. To ensure the fairness of the experiment wasn't compromised, we equally divided the user group into two parts, *Group 1* first collected data directly by themselves *(w/o Ours)* and then collected data with an assistive agent *(w/ Ours)*, while the *Group 2* reversed the order, first *(w/ Ours)* mode and then *(w/o Ours)* mode.

To gain deeper insight into how the learned agent assists the human operator, we visualize several keyframes from the data collection process of three dexterous hand tasks under shared control mode with $\gamma = 0.4$. From Fig. 4, it is evident that human operators are not required
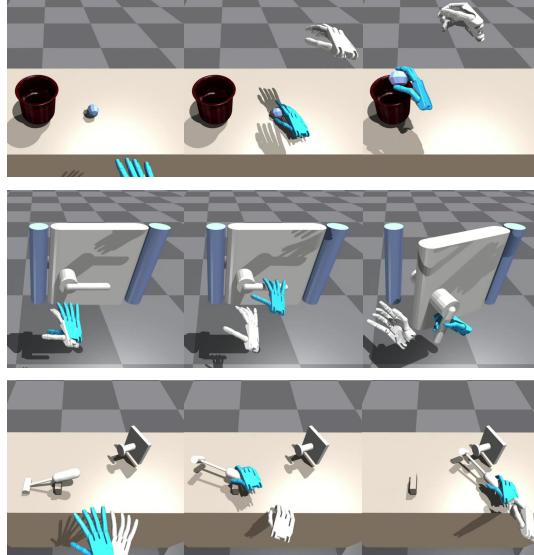


Figure 4: Shared control process overview. The white one is the hand controlled purely by the *human operator*, while the cyan one is under *shared control* between the human and the assistive agent.

to provide too precise control with the assistive agent facilitating shared control over the dexterous hand. Instead, they only need to convey high-level intentions, such as the direction of hand movement or finger grasp motions. In multi-stage tasks, like picking up a hammer and then using it to drive a nail, operators only need to provide a *trigger action* to guide the agent to transition from one sub-stage to the next. As a result, less effort and attention are required, making the data collection easier to execute successfully and speeding it up.
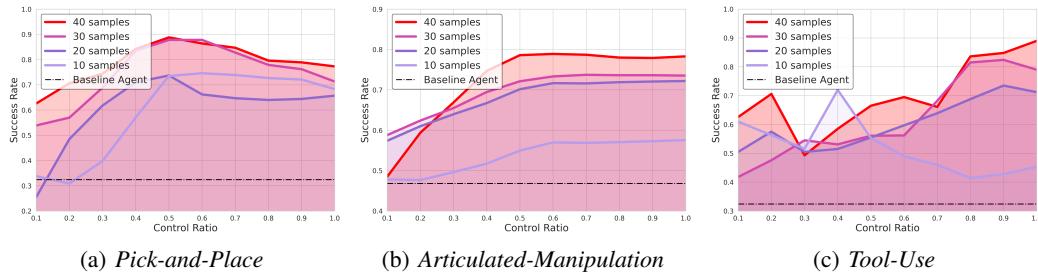
**Quantitative Evaluation.** When the learned agent shares control with users, the system effectively corrects imperfect human control signals to achieve specific tasks. Given the challenge of directly

| | | | |
|---|---|---|---|
| (a) *Pick-and-Place* | (b) *Articulated-Manipulation* | (c) *Tool-Use* |

Figure 5: Additional data collected under our framework effectively contributes to training improvements. We train a simulated operator to evaluate our system, it shows that even with limited data, the learned assist agent can improve the success rate of data collection to improve the efficiency. With the data accumulated, the performance of the learned agent keeps rising. Moreover, the learned agent could be transitioned to a full autonomy agent ($\gamma = 1.0$).

.

| Dexterous Hand | Pick-and-Place | | Articulated-Manipulation | | Tool-Use | |
|---|---|---|---|---|---|---|
| | $40\mathcal{H}$ | $10\mathcal{H} + 30\mathcal{S}$ | $40\mathcal{H}$ | $10\mathcal{H} + 30\mathcal{S}$ | $40\mathcal{H}$ | $10\mathcal{H} + 30\mathcal{S}$ |
| BC | 0.30 | **0.50** | 0.22 | **0.57** | 0.39 | **0.40** |
| BC-RNN | 0.54 | **0.67** | 0.47 | **0.50** | **0.27** | 0.25 |
| DP | 0.73 | **0.76** | 0.77 | **0.78** | 0.88 | **0.89** |

| Parallel Gripper | Pick-and-Place | | Articulated-Manipulation | | Push-cube | |
|---|---|---|---|---|---|---|
| | $40\mathcal{H}$ | $10\mathcal{H} + 30\mathcal{S}$ | $40\mathcal{H}$ | $10\mathcal{H} + 30\mathcal{S}$ | $40\mathcal{H}$ | $10\mathcal{H} + 30\mathcal{S}$ |
| BC | 0.42 | **0.44** | 0.35 | **0.37** | **0.88** | 0.85 |
| BC-RNN | **0.39** | 0.36 | 0.71 | **0.73** | 0.59 | **0.67** |
| DP | 0.51 | **0.60** | 0.42 | **0.67** | **0.83** | 0.82 |

Table 2: Data quality on downstream tasks.

| | Dexterous Tool-Use | | Gripper Push-cube | |
|---|---|---|---|---|
| | BC | DP | BC | DP |
| $10\mathcal{H}$ | 0.29 | 0.45 | 0.23 | 0.42 |
| $10\mathcal{H} + 10\mathcal{H}$ | 0.28 | 0.67 | 0.37 | 0.78 |
| $10\mathcal{H} + 20\mathcal{H}$ | 0.28 | 0.82 | 0.51 | 0.67 |
| $10\mathcal{H} + 30\mathcal{H}$ | 0.39 | 0.88 | 0.88 | 0.83 |
| $10\mathcal{H} + 10\mathcal{S}$ | 0.31 | 0.71 | 0.33 | 0.81 |
| $10\mathcal{H} + 20\mathcal{S}$ | 0.30 | 0.79 | 0.61 | 0.62 |
| $10\mathcal{H} + 30\mathcal{S}$ | 0.40 | 0.89 | 0.85 | 0.82 |

Table 3: Tool-Use and Push-cube task success rate under increasing data.

measuring the level of imperfection in user signals and the correction ability of our system, we simulate the human using a baseline agent trained with Behavior Cloning (BC). Additionally, we import noise to the agent's control signal through the diffusion policy's forward diffuse process.

In Fig. 5, the *x-axis* represents the forward ratio, where a higher forward ratio corresponds to the lower quality of action provided by the simulated operator and a higher control ratio of the learned agent on the dexterous hand. The graph illustrates that with limited data availability, the agent can assist the simulated operator more effectively. As the agent accesses and trains with more data, its ability to correct actions improves. The results show that our system gradually diminishes the demand for the operator's attention and effort, thereby enhancing the efficiency of data collection.

Furthermore, once sufficient data is collected and the assistive agent is trained, it can transition into full autonomy mode by setting $\gamma$ to 1.0 and denoising random actions from the Gaussian distribution noise. Across three different dexterous manipulation tasks, we can achieve success rates of 0.76, 0.78, and 0.89, indicating that the assistive agent can effectively transform into an automated dexterous manipulation agent.

From our experiments, we have observed that the assistive agent significantly aids human operators in managing fine control at the low level, especially in scenarios where accurate observation by humans is challenging, complicating effective action control. For instance, tasks such as grasping an egg or moving a hammer present visual challenges. It can be difficult to visually confirm whether the egg is securely grasped or if there's a risk of it being dropped. This uncertainty makes it hard for human operators to react promptly to sudden changes. However, within our proposed joint learning framework, human operators are primarily required to focus on high-level intentions and task planning during data collection, while the assistive agent manages the detailed low-level actions. This division of labor significantly reduces the burden on human operators by clearly separating strategic planning from execution tasks, streamlining the collaboration between humans and machines.

**Data Quality on Downstream Task.** In this section, we illustrate that collecting data under shared control does not compromise the quality of the data. We gather dexterous hand and gripper manipulation demonstrations via the proposed system in two modes: fully controlling the robots by a human
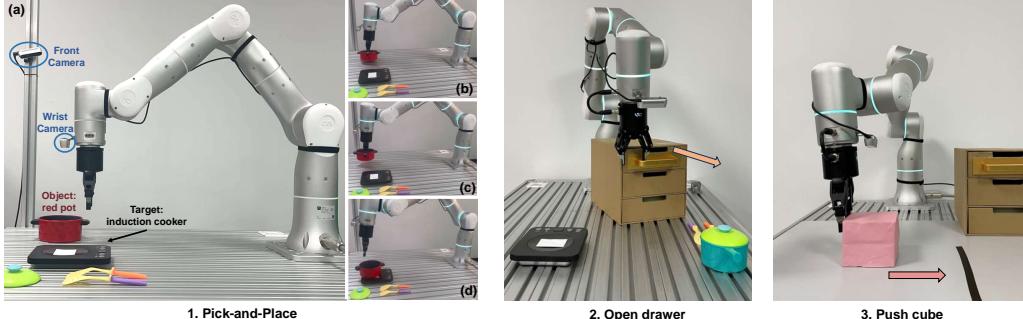
Figure 6: Real world setting. 1. *Pick-and-Place*: (a) Hardware setup. The robot gripper should pick the red pot up (b-c) and place it onto the black induction cooker (d). 2. *Articulated-Manipulation*: use a gripper to open the drawer. 3. *Push-cube*: use a gripper to push the cube across the line.

|         | Success Rate ↑ | Horizon Length ↓ | Collection Speed ↑ |
|---------|----------------|------------------|--------------------|
| w/ Ours | **0.79**       | **18.72**        | **151**            |
| w/o Ours| 0.70           | 21.54            | 121                |

Table 4: Real world parallel gripper Pick-and-Place task user study.

|     | Pick-and-Place | | Articulated-Manipulation | | Push-cube | |
|-----|--------|----------|-------|-----------|--------|------------|
|     | 40H    | 20H + 20S | 30H   | 10H + 20S | 20H    | 10H + 10S  |
| BC  | 13 / 20 | 14 / 20  | 18 / 20 | 19 / 20  | 15 / 20 | 15 / 20   |
| DP  | 11 / 20 | 12 / 20  | 16 / 20 | 12 / 20  | 15 / 20 | 13 / 20   |

Table 5: Real world parallel gripper experiments of data quality.

($\mathcal{H}$) and sharing control ($\mathcal{S}$) between the human operator and the learned assistive agent. And utilize these data to train different kinds of agents, like BC, BC-RNN [7], and Diffusion Policy (DP) [8].

In Tab. 2, compared to directly collecting human demonstrations from the expert human operator, who can achieve success rates and efficiency comparable to those with agent assistance, the data collected by sharing control between the human and the assistive agent can achieve comparable or even surprisingly better results with BC and BC-RNN. Their results are comparable with DP, possibly as DP can better fit the tasks, which is in line with [8].

In Tab. 3, we compare the effects of using different sets of data to train BC and DP. We can find that utilizing more data collected under the shared control mode leads to comparable performance on the tool-use task. This verifies that the new data contributes significantly to policy learning and can achieve a similar effect compared to the data from human experts but at a much lower cost. These results indicate that the data collected under the proposed paradigm have sufficient quality and efficiency for downstream tasks.

**Real World Experiment.** To better evaluate our system, we further conduct real-world experiments. Three tasks are adopted: Pick-and-Place, Articulated-Manipulation, and Push-cube in Fig. 8. Following the same rules as Sec. 4, four human volunteers are invited to participate in the user study to collect data under two modes: one where control is shared between the human operator and the learned agent (*w/ Ours*), and the other where control is directly by the human operator alone (*w/o Ours*). Our proposed system achieves significant improvements in success rate and collection speed by sharing control between human operators and learned agents, as demonstrated in Tab. 4. Additionally, data gathered under our proposed joint learning shared control mode yield performance on the three tasks that are comparable to those pure human datasets using BC and DP, further substantiated by the results presented in Tab. 5.

## 5 Conclusion

In this paper, we introduce a novel human-agent joint learning paradigm that enables simultaneous human demonstration collection and robot manipulation teaching. This approach allows the human operator to share control with a diffusion-model-based assistive agent within a vision-based teleoperation system to control multiple robot end-effectors such as grippers and dexterous hands. Given our paradigm, the human operator can reduce the effort spent on data collection and adjust the control ratio between the human and agent based on different scenarios. Our system offers a more efficient and flexible solution for data collection and robot manipulation learning via teleoperation.

## References

[1] S. P. Arunachalam, I. Güzey, S. Chintala, and L. Pinto. Holo-dex: Teaching dexterity with immersive mixed reality. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5962–5969. IEEE, 2023.

[2] Z. Gharaybeh, H. Chizeck, and A. Stewart. *Telerobotic control in virtual reality*. IEEE, 2019.

[3] H. Liu, X. Xie, M. Millar, M. Edmonds, F. Gao, Y. Zhu, V. J. Santos, B. Rothrock, and S.-C. Zhu. A glove-based system for studying hand-object manipulation via joint pose and force sensing. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 6617–6624. IEEE, 2017.

[4] H. Liu, Z. Zhang, X. Xie, Y. Zhu, Y. Liu, Y. Wang, and S.-C. Zhu. High-fidelity grasping in virtual reality using a glove-based system. In *2019 international conference on robotics and automation (icra)*, pages 5180–5186. IEEE, 2019.

[5] A. Handa, K. Van Wyk, W. Yang, J. Liang, Y.-W. Chao, Q. Wan, S. Birchfield, N. Ratliff, and D. Fox. Dexpilot: Vision-based teleoperation of dexterous robotic hand-arm system. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9164–9170. IEEE, 2020.

[6] Y. Qin, W. Yang, B. Huang, K. Van Wyk, H. Su, X. Wang, Y.-W. Chao, and D. Fox. Anyteleop: A general vision-based dexterous robot arm-hand teleoperation system. *arXiv preprint arXiv:2307.04577*, 2023.

[7] A. Mandlekar, D. Xu, J. Wong, S. Nasiriany, C. Wang, R. Kulkarni, L. Fei-Fei, S. Savarese, Y. Zhu, and R. Martín-Martín. What matters in learning from offline human demonstrations for robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.

[8] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy: Visuomotor policy learning via action diffusion, 2023.

[9] P. Florence, C. Lynch, A. Zeng, O. A. Ramirez, A. Wahid, L. Downs, A. Wong, J. Lee, I. Mordatch, and J. Tompson. Implicit behavioral cloning. In *Conference on Robot Learning*, pages 158–168. PMLR, 2022.

[10] J. Lv, W. Xu, L. Yang, S. Qian, C. Mao, and C. Lu. Handtailor: Towards high-precision monocular 3d hand recovery. *British Machine Vision Conference (BMVC)*, 2021.

[11] Y. Rong, T. Shiratori, and H. Joo. Frankmocap: A monocular 3d whole-body pose estimation system via regression and integration. In *IEEE International Conference on Computer Vision Workshops*, 2021.

[12] T. Schmidt, R. A. Newcombe, and D. Fox. Dart: Dense articulated real-time tracking. In *Robotics: Science and systems*, volume 2, pages 1–9. Berkeley, CA, 2014.

[13] F. Weichert, D. Bachmann, B. Rudak, and D. Fisseler. Analysis of the accuracy and robustness of the leap motion controller. *Sensors*, 13(5):6380–6393, 2013.

[14] S. Javdani, S. S. Srinivasa, and J. A. Bagnell. Shared autonomy via hindsight optimization. *Robotics science and systems: online proceedings*, 2015, 2015.

[15] S. Reddy, A. D. Dragan, and S. Levine. Shared autonomy via deep reinforcement learning. *arXiv preprint arXiv:1802.01744*, 2018.

[16] C. Schaff and M. R. Walter. Residual policy learning for shared autonomy. *arXiv preprint arXiv:2004.05097*, 2020.

[17] A. Ajay, Y. Du, A. Gupta, J. Tenenbaum, T. Jaakkola, and P. Agrawal. Is conditional generative modeling all you need for decision-making?, 2023.

[18] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

[19] M. Janner, Y. Du, J. B. Tenenbaum, and S. Levine. Planning with diffusion for flexible behavior synthesis, 2022.

[20] T. Yoneda, L. Sun, G. Yang, B. C. Stadie, and M. R. Walter. To the noise and back: Diffusion for shared autonomy. In *Robotics: Science and Systems XIX, Daegu, Republic of Korea, July 10-14, 2023*, 2023.

[21] A. Brohan, N. Brown, J. Carbajal, Y. Chebotar, J. Dabis, C. Finn, K. Gopalakrishnan, K. Hausman, A. Herzog, J. Hsu, et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.

[22] F. Ebert, Y. Yang, K. Schmeckpeper, B. Bucher, G. Georgakis, K. Daniilidis, C. Finn, and S. Levine. Bridge data: Boosting generalization of robotic skills with cross-domain datasets. *arXiv preprint arXiv:2109.13396*, 2021.

[23] H.-S. Fang, H. Fang, Z. Tang, J. Liu, J. Wang, H. Zhu, and C. Lu. Rh20t: A robotic dataset for learning diverse skills in one-shot. *arXiv preprint arXiv:2307.00595*, 2023.

[24] J. Kofman, X. Wu, T. J. Luu, and S. Verma. Teleoperation of a robot manipulator using a vision-based human-robot interface. *IEEE transactions on industrial electronics*, 52(5):1206–1219, 2005.

[25] A. Mandlekar, Y. Zhu, A. Garg, J. Booher, M. Spero, A. Tung, J. Gao, J. Emmons, A. Gupta, E. Orbay, et al. Roboturk: A crowdsourcing platform for robotic skill learning through imitation. In *Conference on Robot Learning*, pages 879–893. PMLR, 2018.

[26] J. I. Lipton, A. J. Fay, and D. Rus. Baxter's homunculus: Virtual reality spaces for teleoperation in manufacturing. *IEEE Robotics and Automation Letters*, 3(1):179–186, 2017.

[27] D. Antotsiou, G. Garcia-Hernando, and T.-K. Kim. Task-oriented hand motion retargeting for dexterous manipulation imitation. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.

[28] S. Li, X. Ma, H. Liang, M. Görner, P. Ruppel, B. Fang, F. Sun, and J. Zhang. Vision-based teleoperation of shadow dexterous hand using end-to-end deep neural network. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 416–422. IEEE, 2019.

[29] H. Liu, S. Nasiriany, L. Zhang, Z. Bao, and Y. Zhu. Robot learning on the job: Human-in-the-loop autonomy and learning during deployment. *arXiv preprint arXiv:2211.08416*, 2022.

[30] H. R. Walke, K. Black, T. Z. Zhao, Q. Vuong, C. Zheng, P. Hansen-Estruch, A. W. He, V. Myers, M. J. Kim, M. Du, et al. Bridgedata v2: A dataset for robot learning at scale. In *Conference on Robot Learning*, pages 1723–1736. PMLR, 2023.

[31] H. J. Jeon, D. P. Losey, and D. Sadigh. Shared autonomy with learned latent actions. *arXiv preprint arXiv:2005.03210*, 2020.

[32] A. D. Dragan and S. S. Srinivasa. A policy-blending formalism for shared control. *The International Journal of Robotics Research*, 32(7):790–805, 2013.

[33] S. Javdani, H. Admoni, S. Pellegrinelli, S. S. Srinivasa, and J. A. Bagnell. Shared autonomy via hindsight optimization for teleoperation and teaming. *The International Journal of Robotics Research*, 37(7):717–742, 2018.

[34] K. Muelling, A. Venkatraman, J.-S. Valois, J. E. Downey, J. Weiss, S. Javdani, M. Hebert, A. B. Schwartz, J. L. Collinger, and J. A. Bagnell. Autonomy infused teleoperation with application to brain computer interface controlled manipulation. *Autonomous Robots*, 41:1401–1422, 2017.

[35] D. Sadigh, S. S. Sastry, S. A. Seshia, and A. Dragan. Information gathering actions over human internal state. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 66–73. IEEE, 2016.

[36] D. P. Losey, H. J. Jeon, M. Li, K. Srinivasan, A. Mandlekar, A. Garg, J. Bohg, and D. Sadigh. Learning latent actions to control assistive robots. *Autonomous robots*, 46(1):115–147, 2022.

[37] M. Xu, Z. Xu, C. Chi, M. Veloso, and S. Song. Xskill: Cross embodiment skill discovery, 2023.

[38] C. Luo. Understanding diffusion models: A unified perspective, 2022.

[39] T. Wu, M. Wu, J. Zhang, Y. Gan, and H. Dong. Graspgf: Learning score-based grasping primitive for human-assisting dexterous grasping, 2023.

[40] 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.

[41] S. R. Buss and J.-S. Kim. Selectively damped least squares for inverse kinematics. *Journal of Graphics tools*, 10(3):37–49, 2005.

[42] A. N. Pechev. Inverse kinematics without matrix inversion. In *2008 IEEE International Conference on Robotics and Automation*, pages 2005–2012. IEEE, 2008.

[43] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

# 6 Appendix

## 6.1 Implementation details

Here we lay down the details of the data collection, training, and testing process. More technical details are given here to illustrate our method and implementations better.

### 6.1.1 Shadow Hand and Parallel Gripper Teleoperate System.

To adapt to Isaac Gym and our vision system, we made certain modifications to the XML file of the Shadow Hand. We followed [39, 6, 40], removed the entire arm part, and added six degrees of freedom to the base mount of the Shadow Hand. This allows it to move freely in the virtual environment without depending on a base. Similarly, to obtain the rigid body Jacobian matrices of the five fingertips of the Shadow Hand, we added a massless rigid body to the tips of all five fingers of the Shadow Hand. This facilitates direct inverse kinematics calculations for the entire finger. In inverse kinematics (IK) calculations, we employed the Damped Least Squares (DLS) method [41, 42], this approach helps to prevent instability issues when approaching singularity points. Additionally, the DLS method supports real-time applications because it can provide fast and stable solutions, which is particularly crucial for teleoperation systems. Focusing solely on the five fingertips and wrist is regarded as the most balanced approach between computational efficiency and the precision required for complex

hand movements in real-time applications. The system operates on a computer with an RTX 4070 graphics card and a monitor.

To mitigate the accumulation of errors, the process involves mapping hand motion from the real world into the virtual environment and then comparing each action with the action from the previous frame to calculate a delta action. The reason for calculating delta action is to identify and apply only the changes in movement from one frame to the next, rather than applying the absolute positions and orientations directly. This approach helps reduce the accumulation of errors that might occur due to discrepancies between the real-world movements and their representation in the simulated environment. By focusing on the changes (delta) rather than absolute values, the system can more accurately replicate the intended movements in the simulator, leading to more precise and consistent control of the shadow hand.

**6.1.2   Baselines.** In this section, we provide the implementation details for BC and BC-RNN models. In Behavior Cloning (BC), the objective is to minimize $\mathbb{E}_{(s,a)\sim\mathcal{D}}||\pi_\theta(s)-a||^2$. We use a 3-layer multi-layer perception (MLP) with a ReLU activation function. All layers are fully connected layers with 128 hidden dimensions with a learning rate of $2 \cdot 10^{-3}$. We also use the AdamW [43] to be the optimizer. The training epoch in dexterous tasks Pick-and-Place, Articulated-Manipulation, and Tool-Use is $60, 100, 100$ separately.

As for BC-RNN, we use an LSTM as the backbone network for BC-RNN [7], which we find a slight performance improvement compared to the vanilla RNN model. Following [7], during the training phase, a state-action sequence $\{(s_i, a_i), \cdots, (s_{i+T-1}, a_{i+T-1})\}$ of length $T$ is sampled from the dataset $\mathcal{D}$ and the network will predict the action sequence based on the states as its input. During the inference phase $a_t, h_{t+1} = \pi_\theta(s_t, h_t)$ where $h_t, h_{t+1}$ are the hidden states. Here we set the learning rate to be $2 \cdot 10^{-3}$, and the training epoch to be 60.

## 6.2   Experiment Setups

***Dexterous Hand Pick-and-Place*** aims at picking an object on the table and placing it into a container. The observation space is 24 dimensions, including the dexterous robot hand state (18-dim), the object's position (3-dim), and the container's position (3-dim). The dexterous robot hand state is the position of each fingertip (15-dim) and the wrist position (3-dim). The action space is 28 dimensions, including the state change of each joint (22-dim) and the wrist transformation (6-dim). The object's position is randomized for each attempt within a $10cm \times 10cm$ square on the table.

***Dexterous Hand Articulated-Manipulation*** aims at grasping and unscrewing the door handle to open the door. The observation space is 32 dimensions, including the dexterous robot hand state (18-dim), the door handle's position (3-dim) and quaternion (4-dim), and the door base's position (3-dim) and quaternion (4-dim). In contrast, the action space is 28 dimensions. The door's position is randomized for each attempt within a $40cm \times 40cm$ square on the floor.

***Dexterous Hand Tool-Use*** aims at picking a hammer and using it to drive a nail into a board. The observation space is 32 dimensions, including the dexterous robot hand state (18-dim), hammer's position (3-dim) quaternion (4-dim), and nail's position (3-dim). At the same time, the action space is 28 dimensions. The nail's position is randomized for each attempt within a $10cm \times 10cm$ square on the table.

***Parallel Gripper Pick-and-Place*** aims at picking an object on the table and placing it into a container. The observation space is 27 dimensions, including the five rigid bodies of the gripper to object distances (15-dim), the distance between left and right grippers (3-dim), the object's position (3-dim), the distance between object and target (3-dim,) and the distance between flange and target (3-dim). The action space is 8 dimensions, including the state change of each joint (7-dim) and gripper (1-dim). The object's position is randomized for each attempt within a $10cm \times 10cm$ square on the table.

***Parallel Gripper Articulated-Manipulation*** aims at picking an object on the table and placing it into a container. The observation space is 16 dimensions, including the five rigid bodies of gripper to
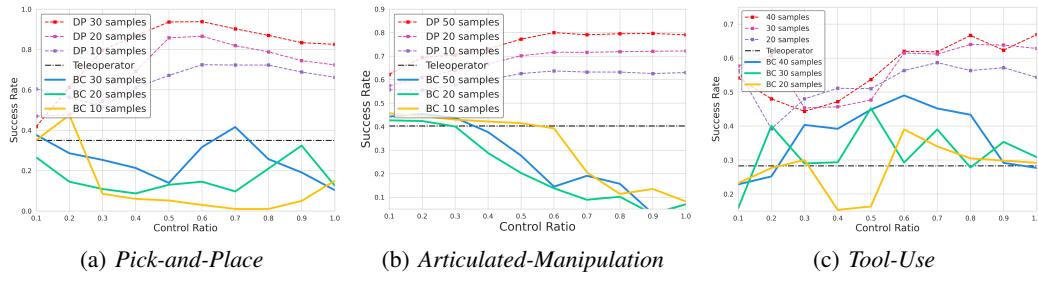
Figure 7: Ablation on different dexterous agents trained with different compositions of data.

object distances (15-dim), and the distance between object and target (1-dim). The action space is 7 dimensions, including the state change of each joint (7-dim). The object's position is randomized for each attempt within a $10cm \times 10cm$ square on the table.

***Parallel Gripper Cube-Push*** aims at pushing an object on the table to the target position. The observation space is 22 dimensions, including the three rigid bodies of the gripper to object distances (9-dim), the flange's position (7-dim), the distance between object and target (3-dim,) and the distance between flange and target (3-dim). The action space is 7 dimensions, including the state change of each joint (7-dim). The object's position is randomized for each attempt within a $5cm \times 5cm$ square on the table.

**6.2.1 Ablation study.** We implement the shared control agent with different methods like the diffusion model and BC. BC adapts a classical way for blending policy to achieve shared control [32]. We use it in the ablation study to blend BC policy with pure human action to achieve shared control in Fig.7. Compared to the classical way which explicitly averages human action $a^h$ and agent action $a^r$ to get the shared action $a^s$, we instead use the diffusion model, which is a popular implicit model, to blend two actions. It models the process as the forward and reverse process. The forward/diffuse process is about adding Gaussian noise to human action $a^h$, and the reverse process uses a neural network $f(\cdot|\cdot)$ to denoise $a^k$ to get the shared action $a^s$.

BC agent is trained using a specific sequence of data collection and fine-tuning steps to optimize performance across different levels of shared control. Initially, we collect data sets of 10, 10, and 20 episodes under various task conditions. These initial datasets are used to train a preliminary agent. Following this initial training phase, we employ the trained agent to assist in further data collection under three different control ratios represented by $\gamma$ values of 0.25, 0.5, and 0.75. The data collected with the assistance of the agent under these $\gamma$ settings are then used to fine-tune the agent.

As shown in Fig.7, experiments demonstrated that the success rate of an assistive agent based on BC is lower than that of an agent based on diffusion models, indicating a reduced capacity for assistance. In certain instances, the action even becomes worse at particular control ratios.

Table 6: Agent performance on human expert or amateur datasets.

| Dexterous Hand | Pick-and-Place | | Articulated-Manipulation | | Tool-Use | |
|---|---|---|---|---|---|---|
| | Skilled | Unskilled | Skilled | Unskilled | Skilled | Unskilled |
| BC | 0.45 | 0.02 | 0.43 | 0.18 | 0.40 | 0.05 |
| BC-RNN | 0.41 | 0.05 | 0.62 | 0.04 | 0.27 | 0.05 |
| DP | 0.71 | 0.01 | 0.68 | 0.10 | 0.81 | 0.03 |

We test the performance of training with different data compositions. For a task, we gathered two manipulation datasets from both skilled and unskilled human operators. We consider operators to be skilled workers if they can practice for more than five hours and reach a success rate and efficiency comparable to those with assistive agents. As shown in Tab. 6, the performance of agents trained

Table 7: Ablation study on DP performance between $r$.

| | Pick-and-Place | Articulated-Manipulation | Tool-Use |
|---|---|---|---|
| $r = 0.0$ | 0.565 | 0.661 | 0.012 |
| $r = 0.1$ | **0.620** | **0.681** | **0.547** |
| $r = 0.2$ | 0.575 | 0.407 | 0.115 |
| $r = 0.3$ | 0.435 | 0.216 | 0.029 |

on the dataset of unskilled operators is much lower than that on the dataset of skilled operators. Therefore, all the human operation datasets $\mathcal{H}$ we use in the main text are from skilled operators.

In our framework, $r$ represents the modification ratio of noise between the state and action. Specifically, during the training, the noise added to state $s$ satisfies $\epsilon_s = u \cdot \mathcal{N}(\mathbf{0}, \mathbf{I})$ while the noise added to action $a$ satisfies $\epsilon_s = v \cdot \mathcal{N}(\mathbf{0}, \mathbf{I})$. Then $r = \frac{u}{v}$. We test different $r$ as shown in Tab. 7, to ensure the best agent performance. We default to using $r = 0.1$ in our model.

**6.2.2 Real World Experiment.** In this section, we evaluate the real-world performance of our method. We use the setup shown in Fig.8, which includes a Flexiv Rizon4 arm equipped with a gripper and two Intel RealSense D435i RGB-D cameras. One camera is mounted on the wrist of the robotic arm, while the second is positioned on the side. One task here is to pick the red pot shown in Fig.8 and place it onto the induction cooker. For more details please refer to CoRL-17.mp4

During the real-world data collection phase, we estimate the human hand's pose using RGBD input. Considering the significant difference in morphology between the human hand and a 7-DoF robotic arm, we chose to track the end effector's position by monitoring the position of the hand's wrist. Additionally, we used the action of closing or opening the human hand as the condition for determining whether to grasp or release an object. This approach leverages the greater dexterity of the human hand to enhance the control and precision of the robotic arm. We record RGB images from two camera views, joint poses (7-dim), gripper width (1-dim), the end effector's position (3-dim), and its quaternion (4-dim). The RGB images have a size of $640 \times 480$ pixels, each episode is sampled at a frequency of 10 Hz.

In real-world experiments, the network architecture is generally similar to the simulation environment's. Our input has changed from the original hand states and object states to the position and orientation of the robot arm end effector, as well as images from the first-person and third-person perspectives. We made two main modifications: 1) For the images, we used a ResNet-18 model. We used a standard ResNet-18 (without pretraining) as the encoder with its global average pooling replaced with a spatial softmax pooling to maintain spatial information. 2) We deepened the layer of the neural network, increased its hidden layer dimension, and expanded action horizon prediction from predicting the next frame action to predicting actions for the subsequent $T$ frames, *i.e.*, $a_{t+1:t+T-1}$ (where $T = 8$).
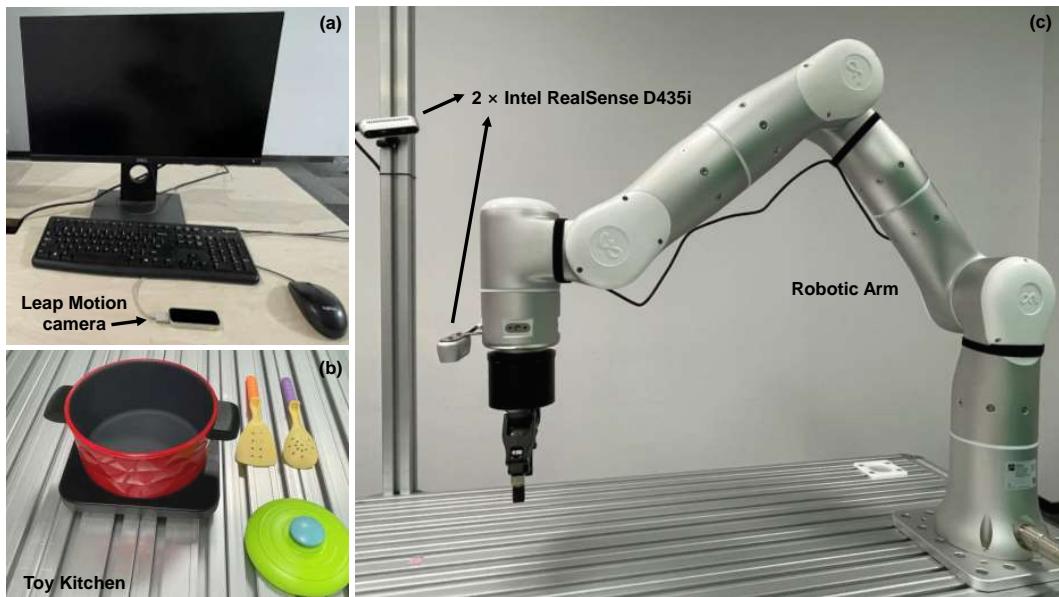
Figure 8: Realworld Pick-and-Place Experiment. The hardware setup comprises (a) a Leap Motion camera utilized for teleoperation data collection, (b) a toy kitchen environment set up for the pick-and-place task, and (c) a Flexiv Rizon4 robotic arm equipped with a gripper and two cameras. One camera is mounted on the wrist of the robotic arm, while the second one is positioned on the side.