

SWBT: Similarity Weighted Behavior Transformer with the Imperfect Demonstration for Robotic Manipulation

Kun Wu¹, Ning Liu², Zhen Zhao², Di Qiu³, Jinming Li⁴,
Zhengping Che², Zhiyuan Xu², Qinru Qiu¹, and Jian Tang^{2,†}

Abstract—Imitation learning (IL), aiming to learn optimal control policies from expert demonstrations, has been an effective method for robot manipulation tasks. However, previous IL methods either only use expensive expert demonstrations and omit imperfect demonstrations or rely on interacting with the environment and learning from online experiences. In the context of robotic manipulation, we aim to conquer the above two challenges and propose a novel framework named Similarity Weighted Behavior Transformer (SWBT). SWBT effectively learn from both expert and imperfect demonstrations without interaction with environments. We reveal that the easy-to-get imperfect demonstrations, such as forward and inverse dynamics, significantly enhance the network by learning fruitful information. To the best of our knowledge, we are the first to attempt to integrate imperfect demonstrations into the offline imitation learning setting for robot manipulation tasks. Extensive experiments on the ManiSkill2 benchmark built on the high-fidelity Sapien simulator and real-world robotic manipulation tasks demonstrated that the proposed method can extract better features and improve the success rates for all tasks. Our code will be released upon acceptance of the paper.

I. INTRODUCTION

Reinforcement learning (RL) and Imitation Learning (IL) techniques have been extensively researched and applied in the context of robotic manipulation tasks [1], [2], [3], [4], [5], [6] in recent years. RL requires an explicit reward function elaborately pre-defined to guide the agent’s action. Nonetheless, many real-world tasks are intrinsically complicated, making the design of an effective reward function challenging [7], [8]. Imitation learning (IL) aims to directly learn from the demonstrations without designing an explicit reward function and thus is more practical to be applied to real-world applications [9], [10], [11]. Many works of IL require online interactions, which are costly or dangerous, such as robotic manipulation and autonomous driving. On the other hand, offline IL enables data collection in advance, thus independent of online training. Meanwhile, early works of IL normally assume the demonstration is optimal, such as behavioral cloning [12]. The assumption of ideal demonstrations is infeasible or extensively costly since collecting high-quality demonstrations by human labor sometimes would make mistakes. For instance, SayCan [13] collected 276k episodes of data and only retained 12k successful episodes

after applying stringent filtering criteria, resulting in more than 90% of the data being wasted.

To address this challenge, recent research in offline IL [14], [15], [16] has emerged, aiming to utilize imperfect demonstrations for agent training. For instance, 2IWIL [14] introduced confidence scores that assess the likelihood of a trajectory being optimal. However, these scores still require manual labeling by human experts for each transition in the imperfect demonstration. DemoDICE [15] formulates an offline IL objective to utilize imperfect demonstrations and mitigate distribution shift from the expert demonstration distribution. However, in the context of robotic manipulation, where data dimensionality is significantly higher, DemoDICE may struggle to accurately align with the data distribution. Another approach, DWBC [17], introduces an additional discriminator to distinguish expert and non-expert demonstrations to mitigate the reward learning as well as leverage imperfect demonstrations. Nonetheless, training the discriminator necessitates costly hyperparameter tuning and struggle with high-dimensional data. Consequently, there is a need for effective offline IL methods capable of leveraging imperfect demonstrations without manual labeling in the context of robotic manipulation with high-dimensional input.

Moreover, recent works [5], [18], [19] have demonstrated that transformers are a suitable option for handling sequential decision-making problems. Recent work Decision Transformer [20] seeks to unify ideas in language modeling and RL by employing the GPT architecture to autoregressively model trajectories and output actions directly. In the context of robotic manipulation settings, we integrate the potent transformer architecture to extract the features of sequential visual-based frames within our proposed framework.

To this end, we propose a novel offline IL-based framework for robotic manipulation, which leveraging imperfect demonstrations, named Similarity Weighted Behavior Transformer (SWBT). Specifically, SWBT contains three main steps. In step 1, we pre-train a multi-modality transformer network using the mixed demonstrations in a self-supervised manner to extract representative features. Since the imperfect demonstrations are readily collected, the mixed demonstrations contain a large amount of fruitful imperfect demonstrations, enabling the transformer to learn large action space knowledge. Inspired by pretraining in Natural Language Processing such as BERT [21], we propose three self-supervised approaches to enhance the pre-trained transformer, including 1) Masked Transition Prediction, 2) Transition Reconstruction, and 3) Action Autoregression.

¹Syracuse University, NY, USA {kwu102, qiuiu}@syr.edu

²Midea Group, China {liuning22, zhaozhen8, chezp, xuzhy70, tangjian22}@midea.com

³Peking University, China qiudi@stu.pku.edu.cn

⁴Shanghai University, China ljm2022@shu.edu.cn

This work was done during Di Qiu’s and Jinming Li’s internship at Midea Group. [†]Corresponding author: Jian Tang.

These approaches enable the transformer to take on different roles like a forward dynamic model, inverse dynamic model, and behavior cloning model and thus enhance its ability to extract effective features. In step 2, we aim to identify the useful imperfect demonstrations for training the transformer. Thus, we leverage the pre-trained transformer in step 1 to extract the features and then generate the quality scores by calculating the similarity between the extracted features of the imperfect and expert demonstrations for each state-action pair in the imperfect demonstrations. By doing so, no extra modules, such as discriminator in [17] are introduced, and the quality scores indicate the quality of the imperfect demonstrations. The simplicity also makes the performance more robust since it does not require much hyper-tuning. In step 3, we perform end-to-end fine-tuning process on the transformer using both expert and imperfect demonstrations. To differentiate the quality of imperfect demonstrations and ensure their beneficial contribution to the transformer, we propose the quality scores calculated in step 2 as weights for the behavior cloning loss. Our contributions as follows.

- We introduce the Similarity Weighted Behavior Transformer (SWBT). To the best of our knowledge, SWBT is the first work to simultaneously utilize both expert and imperfect demonstrations for training robot manipulation policies within offline imitation learning settings.
- The robust transformer architecture is leveraged to pre-train an effective feature extractor and introduced quality scores to identify beneficial imperfect demonstrations and eventually improve the success rates of robotic manipulation tasks.
- Extensive experiments on the ManiSkill2 benchmark demonstrated that SWBT can effectively utilize imperfect demonstrations to improve performance. In addition, our real-world experiments verified that SWBT can be successfully applied to real-world applications.

II. RELATED WORK

A. Offline Imitation Learning with Imperfect Demonstrations

Offline IL has enormous potential to meet critical safety and cost desiderata for real-world applications. Earlier IL method Behavior Cloning (BC) [12] can be used in the offline setting seamlessly and gain comparable performance compared to GAIL methods [22]. IBC [23] and Diffusion Policy [6] further enhance BC by introducing the energy-based models and diffusion models [24]. Another category method, offline inverse reinforcement learning (Offline IRL), focuses on training a discriminator in an adversarial manner [25], [26], [27], [28] or learning a reward function [29] to match the distribution of the expert demonstrations. However, most of these works are based on the assumption that all demonstrations are optimal, thus dropping their performance with imperfect demonstrations.

To mitigate the problem induced by imperfect demonstrations, many imitation learning from imperfect demonstrations methods [14], [30], [31], [32], [29], [33], [34], [35] are proposed. 2IWIL [14] leverages the confidence score to weight

the imperfect demonstrations but requires extra human expert annotations. More recently, several works [15], [16], [17] have been proposed to solve the more challenging problem of offline imitation learning with imperfect demonstrations. DWBC [17] combines the above two algorithms by training a discriminator to learn the weights for BC objectives.

B. Transformers for Robotic Manipulation

Transformers have become the prevalent architecture in various domains. Starting from NLP [36], [37], transformers architectures recently successfully enter the field of vision [38], [39], [40] and achieve impressive results in robot learning [10], [41], [42], [43], [44], [45]. [20], [43], [44], [45] adopt the characteristics of the structure of the transformer to use multi-modality data, such as images and robot states, for policy improvement. Therefore, we leverage the sequential modeling ability of the transformer to improve the performance of the robotic manipulation tasks.

III. METHODOLOGY

A. Preliminary on Imitation Learning

We consider a standard fully observed Markov decision process (MDP) [46] to model the environment. The MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \rho_0, \gamma)$ contains the state space \mathcal{S} , action space \mathcal{A} , reward function $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$, state transition function $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$, the initial state distribution $\rho_0(s)$, and discount factor $\gamma \in (0, 1)$. In our offline IL setting with imperfect demonstrations, the reward function r is unavailable. Instead, we have a static expert demonstration dataset $\mathcal{D}_e = \{(s, a, s')\}$ and a static imperfect demonstration dataset $\mathcal{D}_i = \{(s, a, s')\}$. Our goal is to maximize the success rates for the tasks using both the expert dataset and imperfect dataset $\mathcal{D}_u = \mathcal{D}_e \cup \mathcal{D}_i$ without any online interaction with the environments.

B. Similarity Weighted Behavior Transformer

The overview of the proposed method Similarity Weighted Behavior Transformer (SWBT) is depicted in Figure 1. Specifically, SWBT comprises three main steps. Section III-C details step 1: pre-training a multi-modality transformer network using mixed demonstrations in a self-supervised fashion. Section III-D covers step 2: calculating quality scores based on the similarity between the output features of imperfect and expert demonstrations. Lastly, Section III-E presents step 3: the end-to-end fine-tuning process of the transformer network using weighted behavior cloning.

C. Transformer Pretraining via Self-Supervised Learning

We leverage the transformer as our base model since it is suitable for modeling long dependencies. To pre-train the model, we propose three self-supervised approaches to pre-train the transformer using the union dataset $\mathcal{D}_u = \mathcal{D}_e \cup \mathcal{D}_i$. As shown in Figure 2, the three tasks are 1) Masked Transition Prediction, 2) Transition Reconstruction, and 3) Action Autoregression. For all tasks, the multi-modality transformer network $T(\cdot)$ takes a trajectory segment $\tau_{in} = (o_i, m_i, a_i), \dots, (o_{i+l}, m_{i+l}, a_{i+l})$

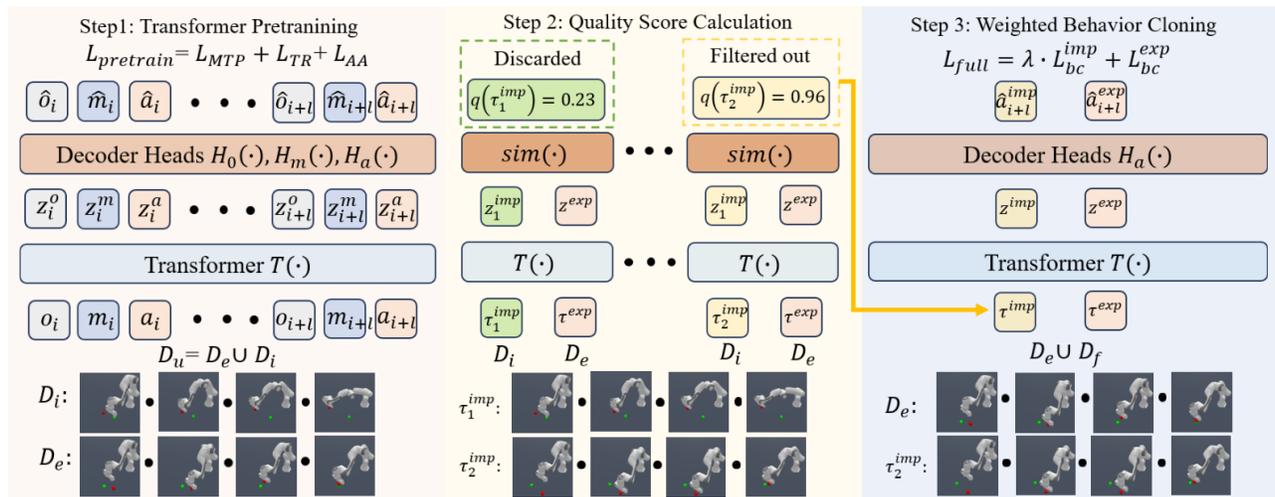


Fig. 1: Overview of SWBT. SWBT contains three steps: 1) Transformer Pretraining via Self-Supervised Learning, 2) Calculation of Quality Score by Similarity Metric, and 3) End-to-End Weighted Behavior Cloning Training. We provide an example of the StackCube task, which aims to stack the red cube on the green one. In step 1, we use both expert and imperfect demonstrations to pretrain the transformer. In step 2, we calculate the quality scores for all imperfect segments (e.g., segment τ_1^{imp} and τ_2^{imp} have quality scores of 0.23 and 0.96 respectively). Then, in step 3, the imperfect segment τ_1^{imp} is discarded and the imperfect segment τ_2^{imp} is used for weighted behavior cloning.

as input and outputs the corresponding features $\tau_{out} = (z_i^o, z_i^m, z_i^a), \dots, (z_{i+l}^o, z_{i+l}^m, z_{i+l}^a)$ for all input elements, where i is the start index of the time step, l is the segment length, o_t is the RGBD images given from the depth cameras, m_t is the robot proprioceptive states and necessary information like goal position, a_t is the action. Then we have three separate decoder heads $H_o(\cdot), H_m(\cdot), H_a(\cdot)$ for the above three kinds of modalities, which takes the corresponding features z_i^o, z_i^m, z_i^a as input and outputs the reconstruction or prediction results $\hat{o}_i, \hat{m}_i, \hat{a}_i$.

For Masked Transition Prediction (MTP) tasks, we randomly mask a part of the input trajectory segment with a pre-defined probability (randomly chosen from 0.4, 0.3, 0.2, 0.1 in our implementation), which is denoted as $M(\tau)$. We denote the optimized parameters in the transformer network $T(\cdot)$ and the decoder heads $H_o(\cdot), H_m(\cdot), H_a(\cdot)$ as θ . And the goal is given by:

$$\mathcal{L}_{MTP} = \max_{\theta} \mathbb{E}_{\tau} \sum_{t=i}^{i+l} \mathbb{I}(e_t) \log P_{\theta}(e_t | M(\tau)), \quad (1)$$

where $e_t \in \{o_t, m_t, a_t\}$ is the input element of the three modalities, $\mathbb{I}(e_t)$ is an indicator function showing the input element is masked (i.e., 1) or not (i.e., 0). By randomly masking input elements and predicting them, we encourage the transformer network to learn various roles, including forward dynamic model, inverse dynamic model, and data generating policy that boost the ability of feature extraction.

For Transition Reconstruction (TR) tasks, we randomly mask a part of the input trajectory segment and reconstruct the unmasked elements. The goal is given by:

$$\mathcal{L}_{TR} = \max_{\theta} \mathbb{E}_{\tau} \sum_{t=i}^{i+l} (1 - \mathbb{I}(e_t)) \log P_{\theta}(e_t | M(\tau)). \quad (2)$$

By reconstructing the unmasked input elements, the transformer network would learn to compress the key information

and extract more representative latent features.

To be consistent with the final objective of the robotic manipulation, which is to output action, we propose the Action Autoregression (AA) task that forces the transformer network to predict the next action based on the history transitions. The goal is given by:

$$\mathcal{L}_{AA} = \max_{\theta} \mathbb{E}_{\tau} \sum_{t=i}^{i+l} \log P_{\theta}(a_t | his(a_t)). \quad (3)$$

where $his(a_t) = (o_i, m_i, a_i), \dots, (o_{i+t}, m_{i+t})$ is the history transitions (i.e., we use causal mask in our implementations).

By combining the above three self-supervised tasks, the final pre-training step objective is as follows:

$$\mathcal{L}_{pretrain} = \mathcal{L}_{MTP} + \mathcal{L}_{TR} + \mathcal{L}_{AA}. \quad (4)$$

D. Calculation of Quality Score by Similarity Metric

To better leverage the high-quality part in the imperfect demonstrations, how to filter them out is a crucial problem. Since the expert demonstrations have high quality, we believe that the imperfect demonstrations that is more similar to expert demonstrations is of higher quality. After pre-training the transformer network in a self-supervised manner, the transformer network now has an enhanced ability to extract expressive and representative features z_i^o, z_i^m, z_i^a for the input elements. Thus we choose to use the output features of the expert demonstrations and the imperfect demonstrations to calculate the similarities and final quality scores. More specifically, for each trajectory segment τ^{imp} in the imperfect demonstrations, we extract the output features $z_{i+l}^{o,imp}, z_{i+l}^{m,imp}, z_{i+l}^{a,imp}$ at the last time step. We also extract the output features $z_{i+l}^{o,exp}, z_{i+l}^{m,exp}, z_{i+l}^{a,exp}$ of each expert trajectory segment τ^{exp} in the same way. The similarity between the imperfect features and the expert features are

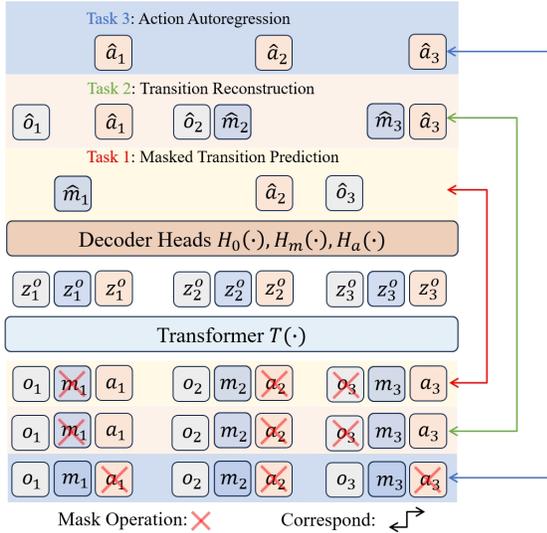


Fig. 2: The transformer pre-training process includes three tasks: 1) Masked Transition Prediction (MTP), 2) Transition Reconstruction (TR), and 3) Action Autoregression (AA). Here is an example of the pre-training process with three time-step inputs. Colored lines link the input and output.

defined as the negative L_2 distance:

$$\begin{aligned} \text{sim}(\tau^{imp}, \tau^{exp}) = & -\|z_{i+l}^{o,imp} - z_{i+l}^{o,exp}\|_2 \\ & -\|z_{i+l}^{m,imp} - z_{i+l}^{m,exp}\|_2 - \|z_{i+l}^{a,imp} - z_{i+l}^{a,exp}\|_2. \end{aligned} \quad (5)$$

The reason why we choose the features of the last time step z_{i+l} of the trajectory segment is that generally the last time step is the most important step for predicting the next action, and the attention mechanism in the transformer has extracted the historical information. For each imperfect trajectory segment, we need to calculate the similarities with all other expert trajectory segments and the similarity result $w(\tau^{imp})$ is the highest one, which is as follows:

$$w(\tau^{imp}) = \max \text{sim}(\tau^{imp}, \tau^{exp}), \quad \forall \tau^{exp} \in \mathcal{D}_e. \quad (6)$$

Once we get similarities of all imperfect segments, we obtain the quality scores by normalizing them:

$$q(\tau^{imp}) = \text{norm}(w(\tau^{imp})), \quad q(\tau^{imp}) \in [0, 1] \quad (7)$$

E. End-to-End Weighted Behavior Cloning Training

After calculating the quality scores for all imperfect trajectory segments τ^{imp} , we can rank them and filter the high-quality part out for learning the optimal policy. In order to prevent low-quality imperfect demonstrations from interfering with policy learning, we define a threshold β and choose the imperfect segments τ^{imp} as follows:

$$\mathcal{D}_f = \{\tau^{imp} | q(\tau^{imp}) > \beta\} \quad (8)$$

For the reserved high-quality segments \mathcal{D}_f , we use their quality scores as the weights for behavior cloning. Combined with the expert demonstrations \mathcal{D}_e , the final objective for fine-tuning is given by:

$$\begin{aligned} \mathcal{L}_{\text{full}} = & \max_{\theta} \mathbb{E}_{\tau \sim \mathcal{D}_e} [\log P_{\theta}(a_{i+l} | \tau)] \\ & + \lambda \mathbb{E}_{\tau \sim \mathcal{D}_f} [q(\tau) \cdot \log P_{\theta}(a_{i+l} | \tau)], \end{aligned} \quad (9)$$

where λ is a hyperparameter to balance the training of the expert demonstrations and imperfect demonstrations. We can count the distribution of quality scores $q(\tau)$ offline and then adjust the value of β accordingly. We set $l = 6$, $\lambda = 0.1$, $\beta = 0.9$ in our implementation.

IV. EXPERIMENTS

We systematically evaluate the Similarity Weighted Behavior Transformer (SWBT) in both simulated and real-world environments. For simulated environments, we built experiments on 5 different tasks on the ManiSkill2 benchmark [47] varying from manipulating rigid and fluid objects, and provided a comprehensive analysis for each component of SWBT with an extensive ablation study. To verify the effectiveness of our algorithm in real-world applications, we built a digital twin system [48] and followed the same setting in the ManiSkill2 benchmark for real-world experiments.

A. Simulation Environments and Datasets

Our simulation experiments are based on the ManiSkill2 benchmark [47], which is built on the high-fidelity Sapien simulator [49] and allows agents to be trained using static demonstrations. The tasks include three rigid object manipulation tasks and two soft-body tasks. For data collection, we trained different level behavior agents using the online DAPG+PPO [50], [51] method provided in the ManiSkill2.

PickCube-v0 is to pick up a red cube and move it to a specified endpoint. We collected 300 expert trajectories and 900 imperfect trajectories consisting of three groups of 300 trajectories, which are collected by three level agents with success rates of 0.0, 0.46, and 0.91, respectively.

StackCube-v0 is to pick up a red cube and place it onto the green one. We collected 300 expert trajectories and 900 imperfect trajectories consisting of three groups of 300 trajectories, which are generated by three level agents with success rates of 0.0, 0.60, and 0.86, respectively.

PickYCB-v0 is to pick up a YCB object [52] and move it to a specified endpoint. We selected 10 objects from all objects that are relatively easy to pick up. We collected 300 expert trajectories and 300 other imperfect trajectories from an imperfect agent with success rates of 0.41.

Fill-v0 is to fill the target beaker with clay in a bucket. 400 expert trajectories and 400 imperfect trajectories are collected from an imperfect agent with success rates of 0.35.

Hang-v0 is to hang a noodle on a specified stick. We collected 400 expert trajectories and 400 imperfect trajectories from an agent with success rates of 0.14.

For the PickYCB, Fill and Hang tasks, we only collected one level of imperfect demonstrations because the behavior policies with high success rates are not available. For all tasks, the robot initialization pose, object position, and goal position are randomly generated for each episode. The robot is a 7-DoF Franka Panda robot with a parallel-jaw gripper. The input RGBD images $o \in \mathbb{R}^{2 \times 128 \times 128 \times 4}$ are from two cameras, including a hand-eye camera and a top fixed camera. The input robot proprioceptive states $m \in \mathbb{R}^{38}$ include joint positions, joint velocities, robot base position,

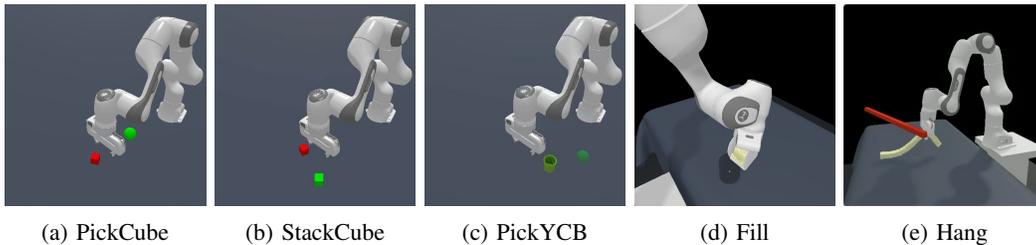


Fig. 3: We conducted experiments on five tasks on the ManiSkill2 benchmark including three rigid-body tasks and two soft-body tasks. The rigid-body tasks include a) PickCube, b) StackCube, and c) PickYCB. The soft-body tasks include d) Fill and e) Hang. The green circles in PickCube and PickYCB represent the goal position.

TABLE I: Comparisons in terms of success rates of the five tasks on ManiSkill2 benchmark. BC and TF-BC are the baselines that can only use expert demonstrations.

Method	PickCube	StackCube	PickYCB	Fill	Hang
BC [12]	82.6	80.0	43.6	24.8	14.4
TF-BC [36]	82.2	81.6	42.4	25.2	16.0
SWBT-base	83.6	84.0	44.4	26.0	16.0
SWBT	85.6	88.0	50.4	29.6	28.4

TABLE II: Comparisons in terms of success rates of the five tasks on ManiSkill2 benchmark. All baselines are proposed for the offline IL setting with imperfect demonstrations.

Method	PickCube	StackCube	PickYCB	Fill	Hang
DT [20]	30.8	78.8	36.8	4.4	6.0
DemoDICE [15]	76.4	78.8	39.6	16.0	8.8
DWBC [17]	32.8	64.0	26.8	10.4	4.0
SWBT	85.6	88.0	50.4	29.6	28.4

goal position, and end-effector position. The output action $a \in \mathbb{R}^7$ is the delta target end-effector pose. We use the success rate to measure the performance of an algorithm. Please refer to the Maniskill2 [47] for more details.

B. Evaluation Results in Simulation

We comprehensively compared SWBT to many state-of-the-art offline imitation learning algorithms including Behavior Cloning (BC) [12], Transformer-based Behavior Cloning (TF-BC) [36], Decision Transformer (DT) [20], DemoDICE [15], and Discriminator-Weighted Behavioral Cloning (DWBC) [17]. When training the DT, we set the goal as 1 for expert demonstrations and 0 for imperfect demonstrations and hope this signal can help the transformer network distinguish the behavior mode. We trained all methods for 100k gradient steps and evaluated 50 episodes every 5k steps. The final results are calculated as the average success rates of the last 5 checkpoints over 3 random seeds.

In Table I, we compared SWBT with methods BC and TF-BC that can only use expert demonstrations. To verify the effectiveness of the pre-training process, we trained **SWBT-base** that only leverages the imperfect demonstrations in the pre-training process. Table I shows that SWBT outperformed all baselines consistently on all 5 tasks by a large margin. These results suggest that the high-quality imperfect demon-

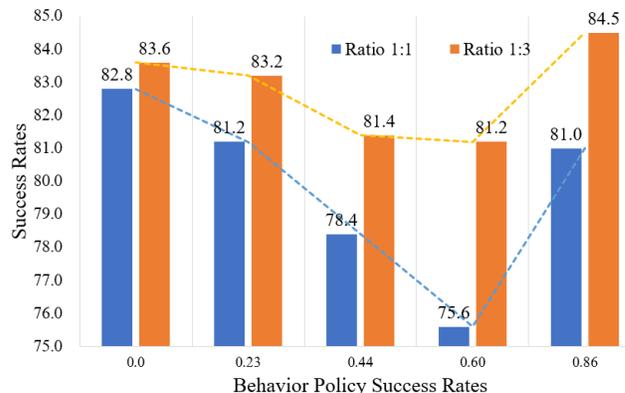


Fig. 4: Comparisons of different quantity and quality of the imperfect demonstrations. The x-axis represents the success rates of the behavior policies that collect the imperfect demonstrations. The y-axis shows the success rates of SWBT. The green and red lines represent the ratios of the expert and imperfect demonstrations are 1:1 and 1:3, respectively.

strations can significantly help improve the final success rate. In addition, the SWBT-base also improved performance on all tasks compared to BC and TF-BC, which shows that the pre-training process can learn expressive features by combining the powerful transformer and imperfect demonstrations.

Table II compared SWBT with methods DT, DemoDICE, and DWBC that simultaneously use both the expert and imperfect demonstration. We can observe that SWBT still surpassed all baseline methods a lot. It is difficult for DemoDICE and DWBC to assign accurate weights for each imperfect high-dimensional transition without a pertaining process. We also found that the training process of DWBC is very unstable without careful hyperparameter tuning. All of the above reasons lead to a decrease in the success rates of these algorithms, even compared to BC. In contrast, SWBT can improve performance by accurately calculating the quality scores and filtering out high-quality demonstrations.

C. Ablation Study in Simulation

Impact of the imperfect demonstration quantity and quality. To verify how the imperfect demonstration quantity and quality influence the final performance of our method, we conducted extensive experiments on the StackCube task. We collected five different level imperfect demonstrations

TABLE III: Comparisons of different quality score threshold β on the StackCube task.

Threshold β	0.0	0.50	0.70	0.80	0.90	0.95	0.99
success rate	17.2	21.6	43.8	66.4	88.0	84.4	83.6
reserved data	180k	57k	35k	16k	8.8k	3.2k	0.3k

TABLE IV: Comparisons of different similarity functions on the StackCube task.

	Seg.+ <i>Cosine</i>	Last+ <i>Cosine</i>	Seg.+ L_2	Last+ L_2
success rate	77.2	83.2	82.0	88.0

using different behavior policies with success rates of 0.00, 0.23, 0.44, 0.60, and 0.86, respectively. For each level, we collected imperfect demonstrations in two ratios, 1:1 and 1:3, where the former 1 represents the 300 expert demonstrations (i.e., 300 and 900 imperfect demonstrations).

Figure 4 shows the success rates with different imperfect demonstration quantities and qualities on the StackCube task. An interesting finding is that those low-quality and high-quality imperfect demonstrations can benefit the performance compared to the medium-quality imperfect demonstrations. We believe that low-quality demonstrations contain rich explorations of the environment and therefore can help extract more informative features in the pre-training process. Then they are discarded in the fine-tuning stage and do not interfere with the results. In contrast, high-quality imperfect demonstrations directly help the process of the fine-tuning stage. The medium-quality demonstrations neither provide rich environmental information nor directly enhance policy learning, so they achieve the most minor improvement.

Impact of quality score threshold β . We evaluated the impact of various quality score thresholds β on the StackCube task. We collected 300 expert trajectories and 900 imperfect trajectories (i.e., 180k transitions) consisting of three groups of 300 trajectories, collected by three level agents with success rates of 0.0, 0.60, and 0.86, respectively.

As shown in Table III, we can observe that as the quality score threshold value β increases, the reserved demonstrations become less and less. The highest success rate on the StackCube task is with β as 0.90. When β is less than 0.90, the performance increases as β increases because more and more low-quality demonstrations are discarded. When β is larger than 0.90, we argue that the performance drops as β increases because more high-quality demonstrations are discarded, degrading the task performance.

Impact of different similarity calculation functions. We evaluated different combinations for the calculation of the similarity distance. We examined L_2 distance and *Cosine* similarity distance with the whole segment features denoted as ‘‘Seg.’’ and the last time step feature denoted as ‘‘Last’’ on the StackCube task. As shown in Table IV, ‘‘Last+ L_2 ’’ achieved the highest performance, while the results of using cosine similarity are relatively lower. We argue that the last time step of the trajectory segment is the most important step

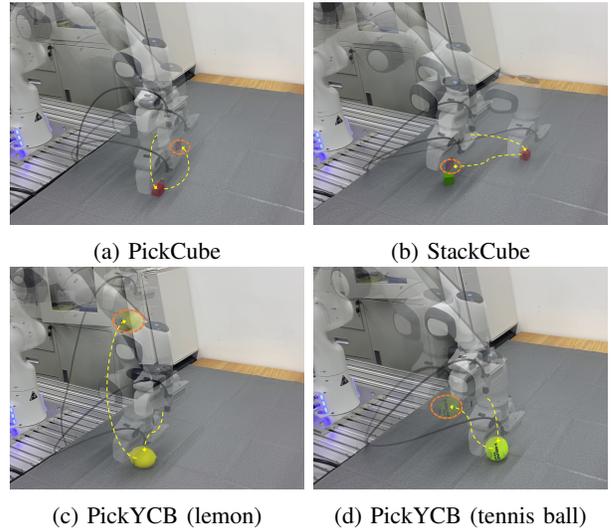


Fig. 5: We conducted a digital twin system [48] following the ManiSkill2 for real-world experiments on the PickCube, StackCube, and PickYCB tasks. Experiments show that SWBT can be applied to real-world applications successfully.

for predicting the next action, and the attention mechanism in the transformer has extracted the historical information.

D. Real-World Experiments

We applied SWBT to real-world applications by following the setting in ManiSkill2 benchmark and conducting a digital twin system [48]. We trained the policy in the simulation and deployed the model on the real Franka Panda robot. The input of the model includes 2 RGBD images from the simulator and robot proprioceptive states from the real robot. Then the generated actions were applied to both the simulated and real-world robot arms. As shown in Figure 5, we conducted experiments on three rigid body tasks including PickCube, StackCube, and PickYCB. Experiments demonstrate that SWBT can be applied to real-world applications successfully. Please refer to the supplementary video for the manipulation process and more details.

V. CONCLUSION AND LIMITATION

To the best of our knowledge, we are the first to use both expert and imperfect demonstrations simultaneously to train the robot manipulation policy in an offline imitation learning setting. Concretely, we propose a novel framework, named **Similarity Weighted Behavior Transformer (SWBT)**, that calculates accurate quality scores using a pre-trained transformer and then does weighted behavior cloning with the high-quality imperfect demonstrations. The experiments on the ManiSkill2 benchmark and real-world applications demonstrated that SWBT can efficiently and correctly leverage imperfect demonstrations to boost the final performance.

Although SWBT can effectively utilize imperfect demonstrations, it still has a potential limitation. SWBT is extremely data-intensive due to the transformer architecture, and thus it is challenging to apply SWBT to tasks with sparse data, e.g., few-shot learning.

REFERENCES

- [1] S. Gu, E. Holly, T. Lillicrap, and S. Levine, "Deep reinforcement learning for robotic manipulation with asynchronous off-policy updates," in *2017 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2017.
- [2] T. Zhang, Z. McCarthy, O. Jow, D. Lee, X. Chen, K. Goldberg, and P. Abbeel, "Deep imitation learning for complex manipulation tasks from virtual reality teleoperation," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018.
- [3] E. Johns, "Coarse-to-fine imitation learning: Robot manipulation from a single demonstration," in *2021 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2021.
- [4] J. Thumm and M. Althoff, "Provably safe deep reinforcement learning for robotic manipulation in human environments," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022.
- [5] 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.
- [6] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song, "Diffusion policy: Visuomotor policy learning via action diffusion," in *Proceedings of Robotics: Science and Systems (RSS)*, 2023.
- [7] A. Y. Ng, D. Harada, and S. Russell, "Policy invariance under reward transformations: Theory and application to reward shaping," in *Icml*. Citeseer, 1999.
- [8] M. Kwon, S. M. Xie, K. Bullard, and D. Sadigh, "Reward design with language models," in *The Eleventh International Conference on Learning Representations*, 2023.
- [9] R. Rahmatizadeh, P. Abolghasemi, L. Bölöni, and S. Levine, "Vision-based multi-task manipulation for inexpensive robots using end-to-end learning from demonstration," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018.
- [10] E. Jang, A. Irpan, M. Khansari, D. Kappler, F. Ebert, C. Lynch, S. Levine, and C. Finn, "Bc-z: Zero-shot task generalization with robotic imitation learning," in *Conference on Robot Learning*. PMLR, 2022.
- [11] B. Zitkovich, T. Yu, S. Xu, P. Xu, T. Xiao, F. Xia, J. Wu, P. Wohlhart, S. Welker, A. Wahid, quan vuong, V. Vanhoucke, H. Tran, R. Soricut, A. Singh, J. Singh, P. Sermanet, P. R. Sanketi, G. Salazar, M. S. Ryoo, K. Reymann, K. Rao, K. Pertsch, I. Mordatch, H. Michalewski, Y. Lu, S. Levine, L. Lee, T.-W. E. Lee, I. Leal, Y. Kuang, D. Kalashnikov, R. Julian, N. J. Joshi, A. Irpan, brian ichter, J. Hsu, A. Herzog, K. Hausman, K. Gopalakrishnan, C. Fu, P. Florence, C. Finn, K. A. Dubey, D. Driess, T. Ding, K. M. Choromanski, X. Chen, Y. Chebotar, J. Carbajal, N. Brown, A. Brohan, M. G. Arenas, and K. Han, "RT-2: Vision-language-action models transfer web knowledge to robotic control," in *7th Annual Conference on Robot Learning*, 2023.
- [12] D. A. Pomerleau, "Alvinn: An autonomous land vehicle in a neural network," *Advances in neural information processing systems*, 1988.
- [13] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang, R. Julian *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," in *Conference on Robot Learning*. PMLR, 2023.
- [14] Y.-H. Wu, N. Charoenphakdee, H. Bao, V. Tangkaratt, and M. Sugiyama, "Imitation learning from imperfect demonstration," in *International Conference on Machine Learning*. PMLR, 2019.
- [15] G.-H. Kim, S. Seo, J. Lee, W. Jeon, H. Hwang, H. Yang, and K.-E. Kim, "Demodice: Offline imitation learning with supplementary imperfect demonstrations," in *International Conference on Learning Representations*, 2021.
- [16] L. Yu, T. Yu, J. Song, W. Neiswanger, and S. Ermon, "Offline imitation learning with suboptimal demonstrations via relaxed distribution matching," in *Proceedings of the AAAI conference on artificial intelligence*, 2023.
- [17] H. Xu, X. Zhan, H. Yin, and H. Qin, "Discriminator-weighted offline imitation learning from suboptimal demonstrations," in *International Conference on Machine Learning*. PMLR, 2022.
- [18] P. Wu, A. Majumdar, K. Stone, Y. Lin, I. Mordatch, P. Abbeel, and A. Rajeswaran, "Masked trajectory models for prediction, representation, and control," *arXiv preprint arXiv:2305.02968*, 2023.
- [19] Y. Sun, S. Ma, R. Madaan, R. Bonatti, F. Huang, and A. Kapoor, "SMART: Self-supervised multi-task pretraining with control transformers," in *International Conference on Learning Representations*, 2023.
- [20] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, "Decision transformer: Reinforcement learning via sequence modeling," *Advances in neural information processing systems*, 2021.
- [21] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [22] J. Ho and S. Ermon, "Generative adversarial imitation learning," *Advances in neural information processing systems*, 2016.
- [23] 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*. PMLR, 2022.
- [24] J. Ho, A. Jain, and P. Abbeel, "Denosing diffusion probabilistic models," *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [25] I. Kostrikov, O. Nachum, and J. Tompson, "Imitation learning via off-policy distribution matching," in *International Conference on Learning Representations*, 2020.
- [26] M. Sun, A. Mahajan, K. Hofmann, and S. Whiteson, "Softdice for imitation learning: Rethinking off-policy distribution matching," *arXiv preprint arXiv:2106.03155*, 2021.
- [27] G. Swamy, S. Choudhury, Z. S. Wu, and J. A. Bagnell, "Of moments and matching: Trade-offs and treatments in imitation learning," *arXiv preprint arXiv:2103.03236*, 2021.
- [28] F. Jarboui and V. Perchet, "Offline inverse reinforcement learning," *arXiv preprint arXiv:2106.05068*, 2021.
- [29] K. Zolna, A. Novikov, K. Konyushkova, C. Gulcehre, Z. Wang, Y. Aytar, M. Denil, N. de Freitas, and S. Reed, "Offline learning from demonstrations and unlabeled experience," *arXiv preprint arXiv:2011.13885*, 2020.
- [30] L. Wang, W. Zhang, X. He, and H. Zha, "Supervised reinforcement learning with recurrent neural network for dynamic treatment recommendation," in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018.
- [31] K. Brantley, W. Sun, and M. Henaff, "Disagreement-regularized imitation learning," in *International Conference on Learning Representations*, 2019.
- [32] D. Brown, W. Goo, P. Nagarajan, and S. Niekum, "Extrapolating beyond suboptimal demonstrations via inverse reinforcement learning from observations," in *International conference on machine learning*. PMLR, 2019.
- [33] D. S. Brown, W. Goo, and S. Niekum, "Better-than-demonstrator imitation learning via automatically-ranked demonstrations," in *Conference on robot learning*. PMLR, 2020.
- [34] V. Tangkaratt, B. Han, M. E. Khan, and M. Sugiyama, "Variational imitation learning with diverse-quality demonstrations," in *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [35] M. Du, S. Nair, D. Sadigh, and C. Finn, "Behavior retrieval: Few-shot imitation learning by querying unlabeled datasets," *arXiv preprint arXiv:2304.08742*, 2023.
- [36] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, 2017.
- [37] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, 2020.
- [38] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021.
- [39] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," *arXiv preprint arXiv:2304.02643*, 2023.
- [40] Z. Wang, Y. Li, X. Chen, S.-N. Lim, A. Torralba, H. Zhao, and S. Wang, "Detecting everything in the open world: Towards universal object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.
- [41] S. Stepputtis, J. Campbell, M. Phielipp, S. Lee, C. Baral, and H. Ben Amor, "Language-conditioned imitation learning for robot manipulation tasks," *Advances in Neural Information Processing Systems*, 2020.
- [42] C. Lynch and P. Sermanet, "Language conditioned imitation learning over unstructured data," *arXiv preprint arXiv:2005.07648*, 2020.

- [43] H. Kim, Y. Ohmura, and Y. Kuniyoshi, "Transformer-based deep imitation learning for dual-arm robot manipulation," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021.
- [44] W. Liu, C. Paxton, T. Hermans, and D. Fox, "Structformer: Learning spatial structure for language-guided semantic rearrangement of novel objects," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022.
- [45] H. Bharadhwaj, J. Vakil, M. Sharma, A. Gupta, S. Tulsiani, and V. Kumar, "Roboagent: Generalization and efficiency in robot manipulation via semantic augmentations and action chunking," *arXiv preprint arXiv:2309.01918*, 2023.
- [46] R. S. Sutton, A. G. Barto *et al.*, *Introduction to reinforcement learning*. MIT press Cambridge, 1998.
- [47] J. Gu, F. Xiang, X. Li, Z. Ling, X. Liu, T. Mu, Y. Tang, S. Tao, X. Wei, Y. Yao, X. Yuan, P. Xie, Z. Huang, R. Chen, and H. Su, "Maniskill2: A unified benchmark for generalizable manipulation skills," in *International Conference on Learning Representations*, 2023.
- [48] K. Xia, C. Sacco, M. Kirkpatrick, C. Saidu, L. Nguyen, A. Kircaliali, and R. Harik, "A digital twin to train deep reinforcement learning agent for smart manufacturing plants: Environment, interfaces and intelligence," *Journal of Manufacturing Systems*, 2021.
- [49] F. Xiang, Y. Qin, K. Mo, Y. Xia, H. Zhu, F. Liu, M. Liu, H. Jiang, Y. Yuan, H. Wang, L. Yi, A. X. Chang, L. J. Guibas, and H. Su, "SAPIEN: A simulated part-based interactive environment," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [50] A. Rajeswaran, V. Kumar, A. Gupta, G. Vezzani, J. Schulman, E. Todorov, and S. Levine, "Learning Complex Dexterous Manipulation with Deep Reinforcement Learning and Demonstrations," in *Proceedings of Robotics: Science and Systems (RSS)*, 2018.
- [51] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [52] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, "The ycb object and model set: Towards common benchmarks for manipulation research," in *2015 International Conference on Advanced Robotics (ICAR)*, 2015.