

TASK-AGNOSTIC PRE-TRAINING AND TASK-GUIDED FINE-TUNING FOR VERSATILE DIFFUSION PLANNER

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models have demonstrated their capabilities in modeling trajectories of multi-tasks. However, existing multi-task planners or policies typically rely on task-specific demonstrations via multi-task imitation, or require task-specific reward labels to facilitate policy optimization via Reinforcement Learning (RL). They heavily rely on task-specific labeled data, which can be difficult to acquire. To address these challenges, we aim to develop a versatile diffusion planner that can leverage large-scale inferior data that contains task-agnostic sub-optimal trajectories, with the ability to fast adapt to specific tasks. In this paper, we propose **SODP**, a two-stage framework that leverages **Sub-Optimal** data to learn a **Diffusion Planner**, which is generalizable for various downstream tasks. Specifically, in the pre-training stage, we train a foundation diffusion planner that extracts general planning capabilities by modeling the versatile distribution of multi-task trajectories, which can be sub-optimal and has wide data coverage. Then for downstream tasks, we adopt RL-based fine-tuning with task-specific rewards to quickly refine the diffusion planner, which aims to generate action sequences with higher task-specific returns. Experimental results from multi-task domains including Meta-World and Adroit demonstrate that SODP outperforms state-of-the-art methods with only a small amount of data for reward-guided fine-tuning.

1 INTRODUCTION

There has been a long-standing pursuit to develop agents capable of performing multiple tasks (Reed et al., 2022; Lee et al., 2022). Although traditional RL methods have made significant strides in training agents to master individual tasks (Silver et al., 2016; OpenAI et al., 2019), expanding this capability to handle diverse tasks remains a significant challenge due to the diversity of task variants and optimization directions with different rewards. Multi-task RL aims to address this by developing agents via task-conditioned optimization (Yu et al., 2020; Lee et al., 2022) or parameter-compositional learning (Sun et al., 2022; Lee et al., 2023). However, the inherent diversity in task trajectory distributions makes it challenging to model and accommodate modeling across different task structures. Diffusion models (Sohl-Dickstein et al., 2015; Ho et al., 2020), originally designed for generative tasks, provide a powerful framework to address these difficulties. Their capacity to capture complex, multi-modal distributions within high-dimensional data spaces (Podell et al., 2023; Ho et al., 2022; Jing et al., 2022) makes them well suited to represent the broad variability encountered in multi-task environments.

Motivated by this, existing methods have employed diffusion models to mimic expert behaviors derived from human demonstrations on various tasks (Pearce et al., 2023; Xu et al., 2023; Chi et al., 2023). However, acquiring task-specific demonstrations is often time-consuming and costly, especially in environments requiring specialized domain expertise. Alternative approaches attempt to optimize diffusion models with return guidance (He et al., 2024; Liang et al., 2023) or conventional RL paradigm (Wang et al., 2022b), which demands a large volume of data with reward labels for each task. To address the above limitations, we wonder whether a generalized diffusion planner can be learned from a large amount of low-quality trajectories without reward labels, with the ability to adapt quickly to various downstream tasks. We only require the inferior data to comprise a mixture of sub-optimal state-action pairs from various tasks, which can be easily obtained in the real world. In training, the diffusion planner seeks to model the distribution of diverse trajectories with broad coverage, enabling it to acquire generalizable capabilities and allowing the planner to further con-

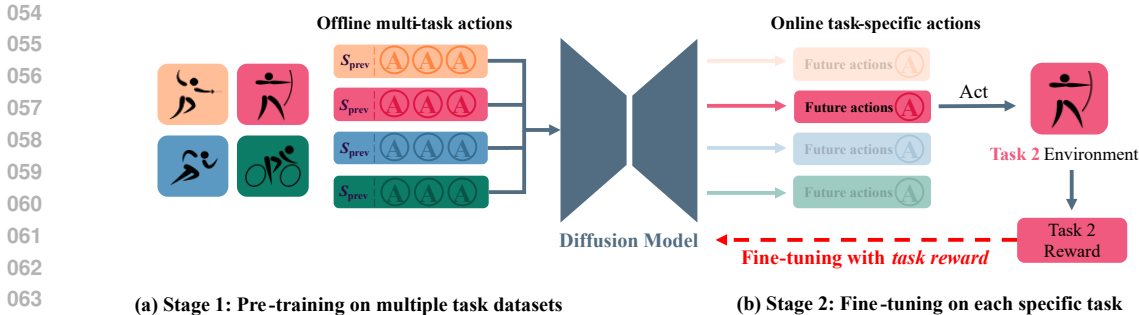


Figure 1: The Overall framework. Different colors represent different tasks. The diffusion model is first pre-trained on a mixed dataset drawn from multiple tasks, and is then fine-tuned for each specific task using task-specific rewards.

centrate on high-reward regions of specific downstream tasks via fast adaptation. An overview of our method is given in Figure 1.

In this paper, we propose a novel framework to utilize Sub-Optimal data to train a **Diffusion Planner (SODP)** that can generalize across a wide range of downstream tasks. SODP consists of two stages: pre-training and fine-tuning. By leveraging a set of trajectories of different tasks for pre-training, we employ action-sequence prediction to capture shared knowledge across tasks. Since the state space may vary between tasks, focusing on the common action space (e.g., end-effector poses of a robot arm) facilitates task generalization. We frame the pre-training stage as a conditional generative problem that generates future actions based on historical states. Then, inspired by the remarkable success of RL-based alignment for LLMs (Ouyang et al., 2022; Glaese et al., 2022), we adopt an RL-based fine-tuning approach to tailor the pre-trained diffusion planner to specific downstream tasks. Specifically, we conduct online interaction based on the pre-trained planner to collect task-specific experiences with reward labels, and perform policy gradients to iteratively refine the predicted action-sequence distribution based on reward feedback of tasks. Through fine-tuning, the diffusion planner can gradually adapt toward generating actions with high task-specific rewards and eventually become optimal for the given task.

Figure 2 illustrates our method. In pre-training, the model captures diverse behavior patterns from the training data, encompassing inferior and mediocre actions. After fine-tuning, the model shrinks the action distribution and concentrates on generating optimal action sequences for a specific task. Our contributions can be summarized as follows. (i) We propose a novel pre-training and fine-tuning paradigm for learning a versatile diffusion planner, which leverages sub-optimal transitions to capture the broad action distributions across tasks, and adopt task-specific fine-tuning to transfer the planner to downstream tasks. (ii) We give an efficient fine-tuning algorithm based on policy gradient for diffusion planners, which progressively shifts the action distribution to concentrate on regions associated with higher task returns. (iii) We conduct extensive experiments using sub-optimal data from state-based Meta-World (Yu et al., 2019) as well as image-based Adroit (Rajeswaran et al., 2017), showcasing its superior performance compared to state-of-the-art approaches.

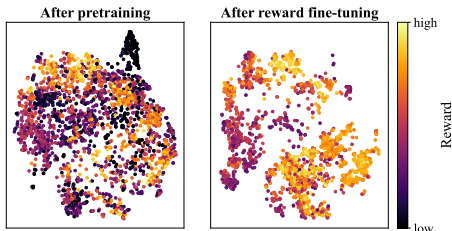


Figure 2: Illustration of SODP in Meta-World *button-press-wall* task. We present trajectories generated by the diffusion model after pre-training and fine-tuning of SODP. The pre-trained model captures a wide range of behaviors, and the fine-tuned model discards the inferior behaviors to coverage to high-reward regions.

2 PRELIMINARIES

Multi-task RL We consider the multi-task RL problem involving N tasks, where each task $\mathcal{T} \sim p(\mathcal{T})$ is represented by a task-specified Markov Decision Process (MDP). Each MDP is defined by a

tuple $(\mathcal{S}^T, \mathcal{A}, P^T, R^T, P_0^T, \gamma)$, where \mathcal{S}^T is the state space of task \mathcal{T} , \mathcal{A} is the global action space, $P^T(s_{t+1}^T | s_t^T, a_t^T) : \mathcal{S}^T \times \mathcal{A} \rightarrow \mathcal{S}^T$ is the transition function, $R^T(s_t^T, a_t^T) : \mathcal{S}^T \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, $\gamma \in (0, 1]$ is the discount factor, and P_0^T is the initial state distribution. We assume that all tasks share a common action space, executed by the same agent, while differing in their respective reward functions, state spaces, and transition dynamics. At each timestep t , the agent perceives a state $s_t^T \in \mathcal{S}^T$, takes an action $a_t^T \in \mathcal{A}$ according to the policy $\pi^T(a_t^T | s_t^T)$, and receives a reward r_t^T . The agent’s objective is to determine an optimal policy that maximizes the expected return across all tasks: $\pi^* = \arg \max_{\pi} \mathbb{E}_{\mathcal{T} \sim p(\mathcal{T})} \mathbb{E}_{a_t \sim \pi^T} \left[\sum_{t=0}^{\infty} \gamma^t r_t^T \right]$.

Diffusion Models Diffusion models (Sohl-Dickstein et al., 2015) are a type of generative model that first add noise to the data \mathbf{x}_0 from a unknown distribution $q(\mathbf{x}_0)$ in K steps through a forward process defined as:

$$q(\mathbf{x}_k | \mathbf{x}_{k-1}) := \mathcal{N}(\mathbf{x}_k; \sqrt{1 - \beta_k} \mathbf{x}_{k-1}, \beta_k \mathbf{I}), \quad (1)$$

where β_k is a predefined variance schedule. Then, a trainable reverse process is constructed as:

$$p_{\theta}(\mathbf{x}_{k-1} | \mathbf{x}_k) := \mathcal{N}(\mathbf{x}_{k-1}; \mu_{\theta}(\mathbf{x}_k, k), \Sigma_k), \quad (2)$$

where $\mu_{\theta}(\mathbf{x}_k, k)$ is the forward process posterior mean as a function of a noise prediction neural network $\epsilon_{\theta}(\mathbf{x}_k, k)$ with a learnable parameter θ (Ho et al., 2020). $\epsilon_{\theta}(\mathbf{x}_k, k)$ can be trained via a surrogate loss as

$$\mathcal{L}_{\text{denoise}}(\theta) := \mathbb{E}_{k \sim [1, K], \mathbf{x}_0 \sim q, \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\|\epsilon - \epsilon_{\theta}(\mathbf{x}_k, k)\|^2 \right]. \quad (3)$$

After training, samples can be generated by first drawing Gaussian noise \mathbf{x}_K and then iteratively denoising \mathbf{x}_K into a noise-free output x_0 over K iterations using the trained model $\epsilon_{\theta}(\mathbf{x}_k, k)$ by

$$\mathbf{x}_{k-1} = \frac{1}{\sqrt{\alpha_k}} \left(\mathbf{x}_k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_{\theta}(\mathbf{x}_k, k) \right) + \sigma_k \mathcal{N}(0, \mathbf{I}), \quad (4)$$

where $\alpha_k := 1 - \beta_k$, $\bar{\alpha}_k := \prod_{s=1}^k \alpha_s$ and $\sigma_k = \sqrt{\beta_k}$.

3 METHOD

We propose SODP, a two-stage framework that leverages large amounts of sub-optimal data to train a diffusion planner that can generalize to downstream tasks. The process is depicted in Figure 3. In the pre-training stage, we train a guidance-free diffusion model to predict future actions based on historical states, using an mixture offline dataset cross tasks without reward labels. In the fine-tuning stage, we refine the pre-trained model using policy gradient to maximize the task-specific rewards, additionally incorporating a regularization term to prevent the model from losing acquired skills.

3.1 PRE-TRAINING WITH LARGE-SCALE SUB-OPTIMAL DATA

Previous works (He et al., 2024) typically model multi-task RL as a conditional generative problem using diffusion models trained on datasets composed of multiple task subsets $\mathcal{D} = \cup_{i=1}^N \mathcal{D}_i$, as:

$$\max_{\theta} \mathbb{E}_{\tau \sim \cup_i \mathcal{D}_i} \left[\log p_{\theta}(\mathbf{x}_0(\tau) | \mathbf{y}(\tau)) \right], \quad (5)$$

which requires additional condition $\mathbf{y}(\tau)$ to guide diffusion model to generate desirable trajectories. For instance, $\mathbf{y}(\tau)$ should contain the return of trajectory $R(\tau)$ and task description Z as prompt. However, the reward label and trajectory description may be scarce or costly to obtain in the real-world. To overcome this challenge, we train a diffusion planner that can learn from offline trajectories transitions (i.e., $\{(s_t, a_t, s_{t+1})\}$) without reward label or task descriptions. Specifically, we model the problem as a guidance-free generation process (Chi et al., 2023):

$$\max_{\theta} \mathbb{E}_{(s_t, \mathbf{a}_t) \sim \cup_i \mathcal{D}_i} \left[\log p_{\theta}(\mathbf{a}_t^0 | s_t) \right]. \quad (6)$$

Here, we represent $\mathbf{x}_0 := \mathbf{a}_t^0 = (a_t, a_{t+1}, \dots, a_{t+H-1})$ as an action sequence, where H is the planning horizon and t is the timestep sampled from dataset \mathcal{D} . As previous work (Chi et al., 2023), we denote s_t as the historical states at timestep t with length T_o , i.e., $s_t := \{s_{t-T_o+1}, \dots, s_{t-1}, s_t\}$. The formulation in Eq. (6) enables the model to learn the broad action-sequence distribution of

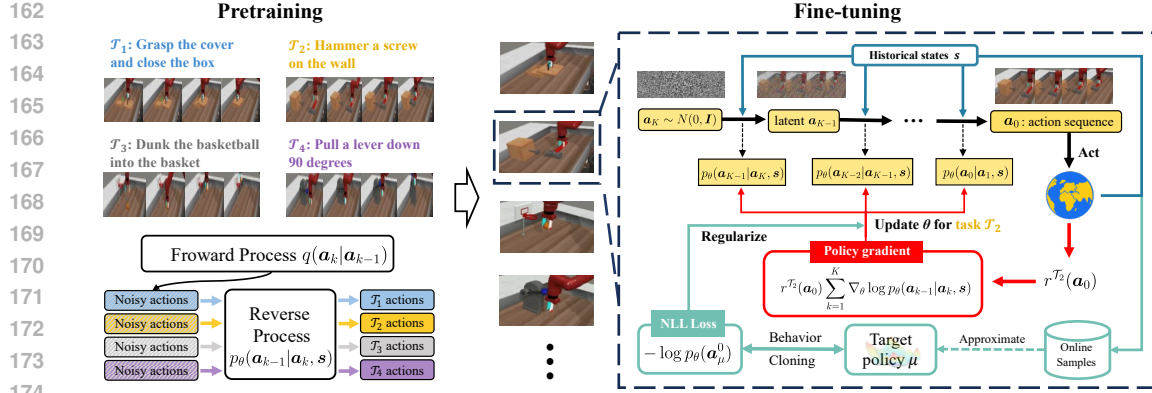


Figure 3: Overview of SODP. We initially pre-train a diffusion model using multi-task transition data to predict action sequences from historical states. Subsequently, we fine-tune the model on downstream tasks using policy gradient methods, incorporating a regularization term to mitigate model degradation.

multi-tasks depending on previous observations, without requiring additional guidance. To train our planning model, we modify Eq. (3) to obtain our pre-training objective as follows:

$$\mathcal{L}_{\text{pre-train}}(\theta) = \mathbb{E}_{k \sim [1, K], (s_t, \mathbf{a}_t^0) \sim D, \epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\|\epsilon - \epsilon_\theta(\mathbf{a}_t^k, \mathbf{s}_t, k)\|^2 \right]. \quad (7)$$

Following Eq. (4), we can generate action sequences through a series of denoising steps:

$$\mathbf{a}_t^{k-1} = \frac{1}{\sqrt{\alpha_k}} \left(\mathbf{a}_t^k - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k}} \epsilon_\theta(\mathbf{a}_t^k, \mathbf{s}_t, k) \right) + \sigma_k \mathcal{N}(0, \mathbf{I}). \quad (8)$$

Unlike other models, the dataset \mathcal{D} we used for the pre-training stage is not restricted to expert trajectories. As shown in Figure 2, we aim to train a foundation model that captures diverse behaviors and learns general capabilities from inferior trajectories, enabling the planner to enhance its representation and action priors through pre-training before learning on downstream tasks.

3.2 REWARD FINE-TUNING FOR DOWNSTREAM TASKS

MDP notation. The fine-tuning stage involves two distinct MDPs: one for RL decision process and the other for the diffusion model denoising process. We use the superscript *diff* (e.g., s_k^{diff} , a_k^{diff}) to denote the MDP associated with diffusion model denoising process, while no superscript is used for the MDP related to the RL process (e.g., s_t , a_t). Additionally, we use $k \in \{K, \dots, 0\}$ to represent the diffusion timestep and $t \in \{1, \dots, T\}$ to represent the trajectory timestep.

We model the denoising process of our pre-trained diffusion planner as a K -step MDP as follows:

$$\begin{aligned} s_k^{\text{diff}} &= (s_t, \mathbf{a}_t^{K-k}), & a_k^{\text{diff}} &= \mathbf{a}_t^{K-k-1}, & P_0^{\text{diff}}(s_0^{\text{diff}}) &= (\delta_{s_t}, \mathcal{N}(0, \mathbf{I})), \\ P^{\text{diff}}(s_{k+1}^{\text{diff}} | s_k^{\text{diff}}, a_k^{\text{diff}}) &= (\delta_{s_t}, \delta_{a_k^{\text{diff}}}), & R^{\text{diff}}(s_k^{\text{diff}}, a_k^{\text{diff}}) &= \begin{cases} r(s_{k+1}^{\text{diff}}) = r(\mathbf{a}_t^0) & \text{if } k = K - 1, \\ 0 & \text{otherwise.} \end{cases}, \\ \pi_\theta^{\text{diff}}(a_k^{\text{diff}} | s_k^{\text{diff}}) &= p_\theta(\mathbf{a}_t^{K-k-1} | \mathbf{a}_t^{K-k}, \mathbf{s}_t), \end{aligned} \quad (9)$$

where s_k^{diff} and a_k^{diff} are the state and action at timestep k , P_0^{diff} and P^{diff} are the initial distribution and transition dynamics, δ is the Dirac delta distribution, R^{diff} is the reward function and $p_\theta(\mathbf{a}_t^{K-k-1} | \mathbf{a}_t^{K-k}, \mathbf{s}_t)$ is the pre-trained diffusion planner. This formulation allows the state transitions in the MDP to be mapped to the denoising process in the diffusion model. The MDP initiates by sampling an initial state $s_0^{\text{diff}} \sim P_0^{\text{diff}}$, which corresponds to sample Gaussian noise \mathbf{a}_t^K at the beginning of the reverse process. At each timestep k , the policy $\pi_\theta^{\text{diff}}(a_k^{\text{diff}} | s_k^{\text{diff}})$ takes an action a_k^{diff} based on current state s_k^{diff} , which corresponds to generate next latent \mathbf{a}_t^{K-k-1} based on current latent \mathbf{a}_t^{K-k} following Eq. (8). The reward remains zero until a noise-free output \mathbf{a}_t^0 is evaluated. Different from

previous text-to-image studies that typically evaluate the final sample using a pre-trained reward model (Black et al., 2023; Fan et al., 2024), we aim to fine-tune the pre-trained diffusion planner to maximize rewards of downstream tasks, which makes constructing reward models for all tasks costly. Therefore, we directly evaluate the generated action sequences in an online RL environment for each specific task \mathcal{T} . Specifically, for any given timestep t , we use the planner to generate future actions $\mathbf{a}_t^0 = (a_t, a_{t+1}, \dots, a_{t+H-1})$ and then execute the first T_a steps. Then we calculate the discounted cumulative reward from the environment to assess the generated sample, expressed as $r(\mathbf{a}_t^0) = \sum_t^{T_a} \gamma^{t-1} r^\mathcal{T}(s_t, a_t)$. We write $r(\mathbf{a}_t^0)$ as shorthand for $r(s_t, \mathbf{a}_t^0)$ for brevity.

Finetuning objective. The objective of fine-tuning our pre-trained diffusion planner is to maximize the expected reward of the generated action sequences for the target downstream task \mathcal{T} , which can be defined as:

$$J^\mathcal{T}(\theta) = \sum_t \mathbb{E}_{p_\theta(\mathbf{a}_t^0 | s_t)} [r^\mathcal{T}(\mathbf{a}_t^0)]. \quad (10)$$

Directly optimizing the objective $J^\mathcal{T}(\theta)$ is intractable since it is infeasible to evaluate the return over all possible actions. Therefore, we utilize policy gradient methods (Sutton et al., 1999), which estimate the policy gradient and apply a stochastic gradient ascent algorithm for updates. The gradient of the objective $J^\mathcal{T}(\theta)$ can be obtained as follows:

$$\nabla_\theta J^\mathcal{T}(\theta) = \sum_t \mathbb{E}_{p_\theta(\mathbf{a}_t^{0:K} | s_t)} \left[r^\mathcal{T}(\mathbf{a}_t^0) \sum_{k=1}^K \nabla_\theta \log p_\theta(\mathbf{a}_t^{k-1} | \mathbf{a}_t^k, s_t) \right]. \quad (11)$$

However, optimizing with Eq. (11) can be computationally intensive, as it requires generating new samples after each optimization step. To enhance sample efficiency and leverage historical sequences, we employ importance sampling, following the approach of proximal policy optimization (PPO) (Schulman et al., 2017), and derive the loss function for reward improvement as follows:

$$\mathcal{L}_{\text{Imp}}^\mathcal{T}(\theta) = \sum_t \mathbb{E}_{p_{\theta_{\text{old}}}(\mathbf{a}_t^{0:K} | s_t)} \left[\sum_{k=1}^K -r^\mathcal{T}(\mathbf{a}_t^0) \max \left(\rho_k(\theta, \theta_{\text{old}}), \text{clip}(\rho_k(\theta, \theta_{\text{old}}), 1 + \epsilon, 1 - \epsilon) \right) \right], \quad (12)$$

where $\rho_k(\theta, \theta_{\text{old}}) = \frac{p_\theta(\mathbf{a}_t^{k-1} | \mathbf{a}_t^k, s_t)}{p_{\theta_{\text{old}}}(\mathbf{a}_t^{k-1} | \mathbf{a}_t^k, s_t)}$ and ϵ is a hyperparameter. Then, we can train our model using $\mathcal{L}_{\text{Imp}}^\mathcal{T}(\theta)$ in an end-to-end manner, which is equivalent to maximizing the objective in Eq. (10).

Regularization term. However, fine-tuning the model solely depending on the reward is insufficient since the model may step too far, which can lead to performance collapse and instability during reward maximization. To address this problem, we introduce a Behavior-Clone (BC) regularization term during the fine-tuning process. Concretely, we aim to constrain our policy θ to closely match a target policy μ , ensuring that θ does not deviate significantly from μ after policy updates. This constraint can be modeled using a negative log-likelihood (NLL) loss as:

$$\min_\theta \mathbb{E}_{\mathbf{a}_\mu^0 \sim p_\mu} [-\log p_\theta(\mathbf{a}_\mu^0)]. \quad (13)$$

Following Ho et al. (2020), we can obtain a surrogate loss to optimize Eq. (13) as follows:

$$\mathcal{L}_{\text{BC}}(\theta) = \mathbb{E}_{k \sim [1, K], \mathbf{a}_\mu^k \sim p_\mu} \left[\|\epsilon(\mathbf{a}_\mu^k, k) - \epsilon_\theta(\mathbf{a}_\mu^k, k)\|^2 \right], \quad (14)$$

where $\epsilon(\mathbf{a}_\mu^k, k)$ represents the ground-truth noise added to \mathbf{a}_μ^k at timestep k , which can be calculated as $\epsilon(\mathbf{a}_\mu^k, k) = \frac{\mathbf{a}_\mu^k - \sqrt{\bar{\alpha}_k} \cdot \mathbf{a}_\mu^0}{\sqrt{1 - \bar{\alpha}_k}}$.

How to select the target policy? Intuitively, an ideal target policy is the optimal policy that generates samples x^* satisfying $\mathcal{C}(x^*) \geq \mathcal{C}(x)$ for all possible x , where $\mathcal{C}(x)$ represents a measure of the performance or quality of the sample, such as the accumulated reward for action sequences. Since μ is unknown during fine-tuning, we approximate it by sampling action sequences \mathbf{a} that satisfy $\mathcal{C}(\mathbf{a}) \approx \mathcal{C}(\mathbf{a}^*)$. In practice, we denote \mathbf{a}^* as the best actions from recent play experience, such as those that yielded the top n highest rewards or successfully completed the given task. We then sample \mathbf{a}^k from these proficient action sequences obtained during online interaction, nearly equivalent to sampling from μ to regularize the fine-tuning process. We also remark that the BC regularizer

is not the only way to incorporate regularization into Eq. (12). For example, a Kullback–Leibler (KL) divergence between the fine-tuned and pre-trained models, or a diffusion pre-train loss can be employed to regularize the fine-tuning process, as shown in text-to-image and text-to-speech generation (Fan et al., 2024; Chen et al., 2024). However, we find these regularization may cause the pre-trained planner trap in sub-optimal regions, hindering performance improvement. We will further discuss them in experiments.

Combining Eq. (12) with Eq. (14), the loss function for reward fine-tuning in downstream tasks $\mathcal{T} \sim p(\mathcal{T})$ is expressed as follows:

$$\mathcal{L}_{\text{fine-tuning}}^{\mathcal{T}}(\theta) = \mathcal{L}_{\text{imp}}^{\mathcal{T}}(\theta) + \lambda \mathcal{L}_{\text{BC}}(\theta), \quad (15)$$

where λ is a weight coefficient. The overall process of pre-training and fine-tuning using SODP is summarized in Alg. 1 in the appendix. Since our goal is to generate complete trajectories rather than individual segments, we utilize a trajectory-level buffer (Zheng et al., 2022) for estimating the target policy μ . Further, to ensure the accuracy of the approximation, we generate several proficient trajectories using the pre-trained model at the beginning of each iteration.

4 RELATED WORK

Diffusion Models in RL. Diffusion models are a leading class of generative models, achieving state-of-the-art performance across a variety of tasks, such as image generation (Ramesh et al., 2021), audio synthesis (Kong et al., 2020; Huang et al., 2023), and drug design (Schneuing et al., 2022; Guan et al., 2024). Recent studies have applied them in imitation learning to model human demonstrations and predict future actions (Li et al., 2024; Reuss et al., 2023). Other approaches have trained conditional diffusion models either as planners (Ajay et al., 2022; Brehmer et al., 2024) or policies (Hansen-Estruch et al., 2023; Kang et al., 2024). However, most of these efforts focus on single-task settings. While some recent works aim to extend diffusion models to multi-task scenarios, they often rely on additional conditions, such as prompts (He et al., 2024) or preference labels (Yu et al., 2024). These methods are limited by their dependence on expert data or explicit task knowledge. In contrast, our method learns broad action-sequence distributions from inferior data to enhance action priors, enabling effective generalization across a range of downstream tasks.

Fine-tuning Diffusion Models. Despite the impressive success of diffusion models, they often face challenges in aligning with specific downstream objectives, such as image aesthetics (Schuhmann et al., 2022), fairness (Shen et al., 2023), or human preference (Xu et al., 2024), primarily due to their training on unsupervised data. Some methods have been proposed to address this issue by directly fine-tuning models using downstream objectives (Prabhudesai et al., 2023; Clark et al., 2023), but they rely on differentiable reward models, which are impractical in RL since accurately modeling rewards with neural networks is quite costly (Kim et al., 2023). Other methods reformulate the denoising process as an MDP and apply policy gradients for fine-tuning (Black et al., 2023; Fan et al., 2024). However, they heavily depend on strong pre-trained models and have proven ineffective in our case. Our goal is to fine-tune a less powerful model that has been trained on inferior data.

Concurrent with our work, DPPO (Ren et al., 2024) also explores reward fine-tuning for refining RL diffusion planners. However, their approach focuses exclusively on single-task settings and allows access to expert demonstrations. In contrast, we train our model on multi-task data without the need for superior demonstrations. Additionally, we analyze the limitations of current regularization methods for versatile RL diffusion models and propose a new regularizer that improves the performance of sub-optimal pre-trained models.

5 EXPERIMENTS

In this section, we conduct experiments to evaluate our proposed method and address the following questions: (1) How does SODP’s performance compare to current methods? (2) Can SODP scale to high-dimensional observation inputs? (3) How does SODP achieve higher rewards during online fine-tuning?

5.1 EXPERIMENTAL SETUP

We evaluate SODP in both state-based and image-based environments. We conduct experiments on the Meta-World benchmark (Yu et al., 2019) for both state-based and image-based tasks. We also perform image-based experiments on the Adroit benchmark (Rajeswaran et al., 2017).

Meta-World. The Meta-World benchmark comprises 50 distinct manipulation tasks, each requiring a Sawyer robot to interact with various objects. These tasks are designed to assess the robot’s ability to handle different scenarios, such as grasping, pushing, pulling, and manipulating objects of varying shapes, sizes, and complexities. While the state space and reward functions differ across tasks, the action space remains consistent. Following recent studies (He et al., 2024; Hu et al., 2024), we extend all tasks to a random-goal setting, referred to as MT50-rand.

Adroit. The Adroit benchmark includes three dexterous manipulation tasks, requiring a 24-degree-of-freedom dexterous hand to solve complex challenges such as in-hand manipulation and tool use. The goals in this environment are also randomized. For Adroit, we use images as the observation to assess whether our method can scale to high-dimensional input.

Datasets. Following previous work (He et al., 2024), for Meta-World, we use a sub-optimal dataset comprising the first 50% of experiences (50M transitions) obtained from the replay buffer of a SAC (Haarnoja et al., 2018) agent during training. To verify the applicability of our method to tasks of varying difficulty levels, we divide the entire dataset into four subsets based on the task categories presented in Seo et al. (2023). For Adroit, we train a VRL3 (Wang et al., 2022a) agent for each task and use the initial 30% experiences (90K transitions) from the converged replay buffer. For Meta-World, all baselines and our pre-training stage are trained on the same dataset. For Adroit, the baselines are trained on expert demonstrations and ours is trained on sub-optimal transitions.

Baselines. For Meta-World, we compare our proposed SODP with the following baselines: (1) **MT-SAC.** Extended SAC with one-hot task ID as additional input. (2) **MTBC.** Extended BC to multi-task learning through network scaling and a task-ID-conditioned actor. (3) **MTIQL.** Extended IQL (Kostrikov et al., 2021) with multi-head critic networks and a task-ID-conditioned actor for multi-task policy learning. (4) **MTDQL.** Extended Diffusion-QL (Wang et al., 2022b) which is similar to MTIQL. (5) **MTDT.** Extended Decision Transformer (DT) (Chen et al., 2021a) to multitask settings by incorporating task ID encoding and state inputs for task-specific learning. (6) **Prompt-DT** (Xu et al., 2022). An extension of DT, which generates actions by utilizing trajectory prompts and reward-to-go signals. (7) **MTDIFF** (He et al., 2024). A diffusion-based approach that integrates Transformer architectures with prompt learning to facilitate generative planning in multitask offline environments. We extend it with a visual extractor in image-based Meta-World experiments. (8) **HarmoDT** (Hu et al., 2024). A DT-based approach that leverages parameter sharing to exploit task similarities while mitigating the adverse effects of conflicting gradients simultaneously. The results for these baselines are directly replicated from those reported in HarmoDT (Hu et al., 2024).

The action space for different tasks in Adroit is different and is incompatible with MTDIFF and HarmoDT. Therefore, we compare SODP with following baselines designed for complex environments: (1) **BCRNN** (Mandlekar et al., 2021). A variant of BC that employs a Recurrent Neural Network (RNN) as the policy network, predicting the sequence of actions based on the sequence of states as input. (2) **IBC** (Florence et al., 2022). Extended BC with energy-based models (EBM) to train implicit behavioral cloning policies. (3) **Diffusion Policy** (Chi et al., 2023). A diffusion-based approach that predicts future action sequences based on historical states. (4) **DP3** (Ze et al., 2024). A visual imitation learning algorithm that incorporates 3D visual representations into diffusion policies, using a point clouds encoder to process visual observations into visual features. The results for these baselines are directly replicated from those reported in DP3 (Ze et al., 2024).

5.2 RESULTS

We use the average success rate across all tasks as the evaluation metric and report the mean and standard deviation of success rates across three seeds. All baselines are trained on sub-optimal data. As shown in Table 1, our method achieves over a 60% success rate when learning from inferior data, outperforming all baseline methods. Compared to the existing state-of-the-art approach, our method demonstrates a 5.9% improvement. Notably, when compared to MTDIFF, the current leading method based on diffusion models, our approach shows a 24.4% improvement.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

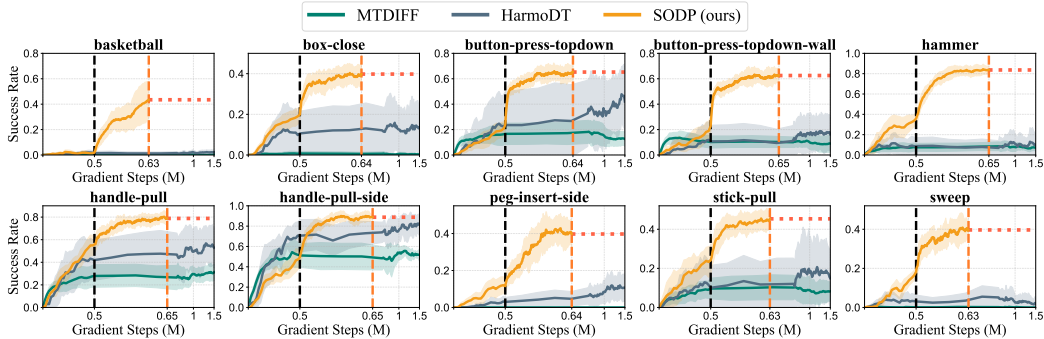


Figure 4: Learning efficiency. We sample 10 tasks and present the learning curves of SODP, MTDIFF, and HarmoDT across five seeds. X-axis represents gradient steps. We pre-train the planner for $5e^5$ steps, followed by fine-tuning with a smaller number of steps. SODP rapidly converges to high success rates, whereas MTDIFF and HarmoDT struggle with some challenging tasks.

MTBC performs the worst, as imitation learning heavily depends on data quality, and directly cloning behaviors from sub-optimal data typically results in inferior performance. In contrast, our method models versatile action distributions from low-quality data and leverages them as priors to guide policy optimization in downstream tasks, leading to improved performance. We conduct additional experiments by augmenting original dataset with online trajectories and results can be found in Appendix C.1.

To further analyze the learning dynamics, we sample 10 tasks and present their learning curves of SODP alongside two leading baselines, MTDIFF and HarmoDT, across five seeds. As shown in Figure 4, SODP rapidly converges to high success rates, surpassing the other two baselines. The pre-training stage equips the planner with comprehensive action distribution priors and allows it to rapidly transfer and enhance these capabilities across a variety of downstream tasks. As a result, the pre-training stage significantly accelerates convergence, leading to more efficient learning in the fine-tuning stage. The two baseline approaches struggle to address complex and challenging tasks such as *basketball* and *hammer*. In contrast, our method effectively guides the model to generate proficient actions, demonstrating the benefits of fine-tuning with policy gradient concerning return maximization. Moreover, while HarmoDT exhibits instability across different random seeds, our method demonstrates robustness against randomness.

Does SODP generalize to high-dimensional observations? We scale our method to image-based observations using the Adroit benchmark by employing a point-cloud encoder from DP3 (Ze et al., 2024) to process the 3D scene represented by point clouds. Specifically, we capture depth images directly from the environment and convert them into point clouds using Open3D (Zhou et al., 2018). These point clouds are then processed by the DP3 Encoder, which maps them into visual features. We then train our diffusion planner following the same procedure in Algorithm 1 except the input states are visual features. Following DP3 (Ze et al., 2024), We compute the average of the highest 5 evaluation success rates during training and report the mean and std across 3 seeds. As shown in Table 2, our method achieves an 8.2% improvement across all tasks. Since *hammer* is more challenging than

Table 1: Average success rate across 3 seeds on Meta-World 50 tasks with random goals (MT50-rand), using sub-optimal data. Each task is evaluated for 50 episodes.

Method	Meta-World 50 Tasks
MTSAC	42.67 \pm 0.12
MTBC	34.53 \pm 1.25
MTIQL	43.28 \pm 0.90
MTDQL	17.33 \pm 0.03
MTDT	42.33 \pm 1.89
Prompt-DT	48.40 \pm 0.16
MTDIFF-P	48.67 \pm 1.32
MTDIFF-P-ONEHOT	48.94 \pm 0.95
HarmonDT-R	53.80 \pm 1.07
HarmonDT-M	57.20 \pm 0.73
HarmonDT-F	57.20 \pm 0.68
SODP (ours)	60.56\pm0.14

Table 2: Average success rate across 3 seeds on Adroit 3 tasks. IBC and BCRNN are extended by incorporating the DP3 point cloud encoder, resulting in IBC+3D and BCRNN+3D.

Algorithm \ Task	Adroit			Average
	Hammer	Door	Pen	
BCRNN	0 \pm 0	0 \pm 0	9 \pm 3	3.0
BCRNN+3D	8 \pm 14	0 \pm 0	8 \pm 1	5.3
IBC	0 \pm 0	0 \pm 0	9 \pm 2	3.0
IBC+3D	0 \pm 0	0 \pm 0	10 \pm 1	3.3
Diffusion Policy	48 \pm 17	50 \pm 5	25 \pm 4	31.7
Simple DP3	100 \pm 0	58 \pm 4	46 \pm 5	68.0
DP3	100 \pm 0	62 \pm 4	43 \pm 6	68.3
SODP (ours)	67 \pm 6	96 \pm 1	59 \pm 4	73.9

door, our method may need more insightful priors from pre-training to achieve better performance. Experiments on image-based Meta-World can be found in Appendix C.5.

5.3 EFFECTIVENESS OF BC REGULARIZATION

To demonstrate the effectiveness of our BC regularization, we conduct an ablation study on fine-tuning same pre-trained model with our BC regularization and other variants. We consider following variants:

- **SODP w/o regularization.** This variant is similar to DDPO (Black et al., 2023) and DPPO (Ren et al., 2024), which fine-tunes the model directly using Eq. (12) without any regularization.
- **SODP_kl.** This variant is similar to DPOK (Fan et al., 2024), with the addition of a KL regularization term to constrain the divergence between the fine-tuned model and the pre-trained model.
- **SODP_pl.** This variant is similar to DLPO (Chen et al., 2024), incorporating the original diffusion pre-training loss (PL) into the fine-tuning objective to prevent the model from deviation.

The details of these variants are presented in Appendix E and more ablation studies on different fine-tuning methods can be found in Appendix C.2. Figure 5 demonstrates the effectiveness of our regularization in achieving a higher success rate. We observe that directly fine-tuning the model without any regularization results in the worst performance, with a decline in success rate, as the model may degrade the capabilities acquired from pre-training due to the lack of constraints. However, adding KL and PL is insufficient, as they cause oscillations near the pre-trained model. This aligns with the original intent of these regularizers, which is to prevent excessive deviation. This is reasonable for methods like DPOK and DLPO, which utilize pre-trained models such as Stable Diffusion (Rombach et al., 2022) and WaveGrad2 (Chen et al., 2021b). These models already exhibit strong generative capabilities without fine-tuning, and the goal is to make slight adjustments to align them with more fine-grained attributes, such as aesthetic scores and human preferences.

In contrast, our model is pre-trained on sub-optimal data and lacks the ability to solve complex tasks. We expect it to develop new skills for completing these tasks through fine-tuning. However, directly applying KL regularization to the pre-trained model leads to conservative policies that heavily rely on the existing capability, thereby confining the model to a sub-optimal region. While PL regularization allows some slight exploration, it is uncontrolled and random. Consequently, we observe that the KL regularization almost remains unchanged and the PL regularization slightly increases the performance in *basketball* but decreases in other tasks.

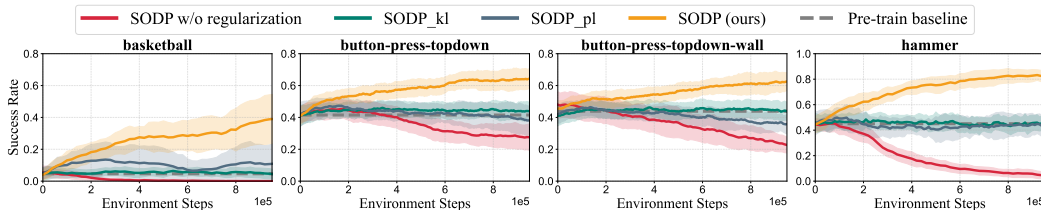
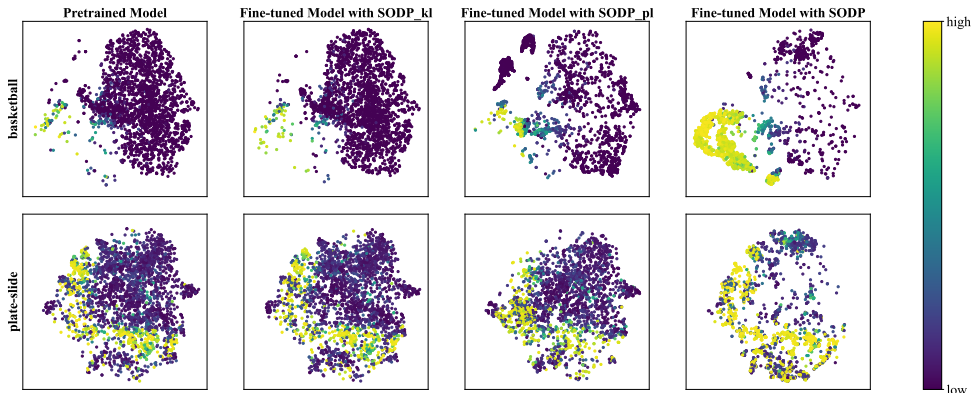


Figure 5: Learning efficiency for different regularization. X-axis represents environment steps. Performance declines significantly without any regularization. Both KL and PL regularization confine the model to sub-optimal regions. In contrast, our BC regularization effectively guides the model away from these sub-optimal areas, facilitating the attainment of optimal actions.

Visualization. We hypothesize that the effectiveness of our BC regularization lies in two aspects: (i) it ensures that our model can reuse the skills it has acquired, thereby preventing a decline in performance; (ii) It guides our model to effectively explore optimal regions due to the utilization of optimal μ as the target policy. To demonstrate this, we visualize trajectories of using the actions generated by our planner using t-SNE (Van der Maaten & Hinton, 2008). As shown in Figure 6, the trajectory distribution after fine-tuning with KL regularization closely resembles the original pre-training distribution, indicating that the model is reusing learned actions and lacks exploration into new regions. The exploration in PL is unstructured as it may lead to worse regions (e.g., the upper-left region in *basketball*). In contrast, our method demonstrates superior exploration capabilities to discover new, high-reward regions based on acquired knowledge (e.g., the lower-left region in

486 *basketball* and the bottom region in *plate-slide*). Meanwhile, the model can derive valuable insights
 487 from pre-trained knowledge by exploiting discovered high-reward actions (e.g. the central region in
 488 *plate-slide*) while discarding low-reward actions.

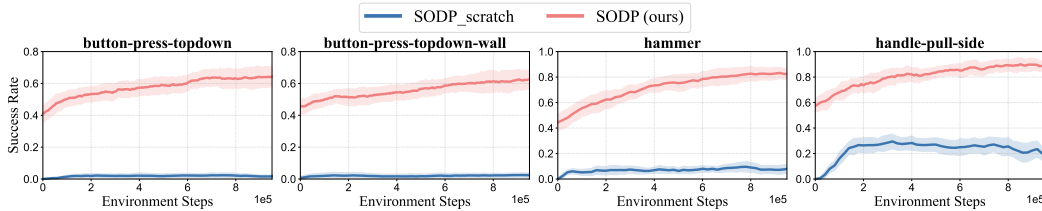


503 Figure 6: Visualization of trajectories using generated actions for different regularization. KL and
 504 PL regularization results in conservative policies with distributions closely resembling the original.
 505 Our BC regularization retains pre-trained knowledge while effectively discovering new actions that
 506 can lead to high rewards.

508 5.4 EFFECTIVENESS OF PRE-TRAINING

509 We investigate the impact of pre-training. We compare the performance of SODP with a version
 510 trained from scratch (**SODP_scratch**). For **SODP_scratch**, we use the same rollouts generated by
 511 the pre-trained model to approximate the target policy and initialize the replay buffer.

512 Figure 7 shows that fine-tuning the planner from scratch results in worse performance. Without
 513 pre-training, the planner lacks an action prior to guide its behavior, leading to stagnation as it strug-
 514 gles to move towards high reward regions. Additionally, it becomes unstable, as the limited useful
 515 knowledge is easily disrupted by a large number of ineffective trials.



516 Figure 7: Effectiveness of pre-training. X-axis represents environment steps. Fine-tuning from
 517 scratch struggles to identify high-reward actions due to the lack of representation prior. In contrast,
 518 pre-training allows the planner to extract useful knowledge, guiding fine-tuning by refining the prior
 519 distribution towards more effective behaviors.

525 6 CONCLUSION

526 We propose SODP, a novel framework for training a versatile diffusion planner using sub-optimal
 527 data. By effectively combining pre-training and fine-tuning, we capture broad behavioral patterns
 528 drawn from large-scale multi-task transitions and then rapidly adapt them to achieve higher perfor-
 529 mance in specific downstream tasks. During fine-tuning, we introduce a BC regularization method,
 530 which preserves the pre-trained model’s capabilities while guiding effective exploration. Experi-
 531 ments demonstrate that SODP achieves superior performance across a wide range of challenging
 532 manipulation tasks. In future work, we aim to develop embodied versatile agents that can effectively
 533 learn to solve real-world tasks using inferior data.

540 ETHICS STATEMENT

541
542 All procedures in this paper were conducted in accordance with the ICLR Code of Ethics (<https://iclr.cc/public/CodeOfEthics>).
543

544 REPRODUCIBILITY STATEMENT

545
546 We have provided all the implementation details necessary to reproduce our experiments in Ap-
547 pendix B, and the dataset used is the same as the one proposed by He et al. (2024).
548

549 REFERENCES
550

- 551 Anurag Ajay, Yilun Du, Abhi Gupta, Joshua Tenenbaum, Tommi Jaakkola, and Pulkit Agrawal.
552 Is conditional generative modeling all you need for decision-making? *arXiv preprint*
553 *arXiv:2211.15657*, 2022.
- 554 Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion
555 models with reinforcement learning. *arXiv preprint arXiv:2305.13301*, 2023.
- 556 Johann Brehmer, Joey Bose, Pim De Haan, and Taco S Cohen. Edgi: Equivariant diffusion for
557 planning with embodied agents. *Advances in Neural Information Processing Systems*, 36, 2024.
- 558 Jingyi Chen, Ju-Seung Byun, Micha Elsner, and Andrew Perrault. Reinforcement learning for fine-
559 tuning text-to-speech diffusion models. *arXiv preprint arXiv:2405.14632*, 2024.
- 560 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel,
561 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
562 modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021a.
- 563 Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, Najim Dehak, and
564 William Chan. Wavegrad 2: Iterative refinement for text-to-speech synthesis. *arXiv preprint*
565 *arXiv:2106.09660*, 2021b.
- 566 Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shu-
567 ran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint*
568 *arXiv:2303.04137*, 2023.
- 569 Kevin Clark, Paul Vicol, Kevin Swersky, and David J Fleet. Directly fine-tuning diffusion models
570 on differentiable rewards. *arXiv preprint arXiv:2309.17400*, 2023.
- 571 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
572 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Reinforcement learning for fine-
573 tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36,
574 2024.
- 575 Pete Florence, Corey Lynch, Andy Zeng, Oscar A Ramirez, Ayzaan Wahid, Laura Downs, Adrian
576 Wong, Johnny Lee, Igor Mordatch, and Jonathan Tompson. Implicit behavioral cloning. In
577 *Conference on Robot Learning*, pp. 158–168. PMLR, 2022.
- 578 Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Mari-
579 beth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, et al. Improving alignment of
580 dialogue agents via targeted human judgements. *arXiv preprint arXiv:2209.14375*, 2022.
- 581 Jiaqi Guan, Xiangxin Zhou, Yuwei Yang, Yu Bao, Jian Peng, Jianzhu Ma, Qiang Liu, Liang Wang,
582 and Quanquan Gu. Decompdiff: diffusion models with decomposed priors for structure-based
583 drug design. *arXiv preprint arXiv:2403.07902*, 2024.
- 584 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
585 maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*
586 *ence on machine learning*, pp. 1861–1870. PMLR, 2018.
- 587 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine.
588 Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint*
589 *arXiv:2304.10573*, 2023.

- 594 Haoran He, Chenjia Bai, Kang Xu, Zhuoran Yang, Weinan Zhang, Dong Wang, Bin Zhao, and Xue-
595 long Li. Diffusion model is an effective planner and data synthesizer for multi-task reinforcement
596 learning. *Advances in neural information processing systems*, 36, 2024.
597
- 598 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
599 *neural information processing systems*, 33:6840–6851, 2020.
600
- 601 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
602 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
603 video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.
- 604 Shengchao Hu, Ziqing Fan, Li Shen, Ya Zhang, Yanfeng Wang, and Dacheng Tao. Harmodt: Har-
605 mony multi-task decision transformer for offline reinforcement learning. In *International Confer-*
606 *ence on Machine Learning*, 2024.
607
- 608 Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin
609 Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced
610 diffusion models. In *International Conference on Machine Learning*, pp. 13916–13932. PMLR,
611 2023.
- 612 Bowen Jing, Gabriele Corso, Jeffrey Chang, Regina Barzilay, and Tommi Jaakkola. Torsional dif-
613 fusion for molecular conformer generation. *Advances in Neural Information Processing Systems*,
614 35:24240–24253, 2022.
615
- 616 Bingyi Kang, Xiao Ma, Chao Du, Tianyu Pang, and Shuicheng Yan. Efficient diffusion policies for
617 offline reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- 618 Changyeon Kim, Jongjin Park, Jinwoo Shin, Honglak Lee, Pieter Abbeel, and Kimin Lee. Pref-
619 erence transformer: Modeling human preferences using transformers for rl. *arXiv preprint*
620 *arXiv:2303.00957*, 2023.
621
- 622 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
623 2014.
624
- 625 Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile
626 diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020.
- 627 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-
628 learning. *arXiv preprint arXiv:2110.06169*, 2021.
629
- 630 Kuang-Huei Lee, Ofir Nachum, Mengjiao Sherry Yang, Lisa Lee, Daniel Freeman, Sergio Guar-
631 rama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, et al. Multi-game decision trans-
632 formers. *Advances in Neural Information Processing Systems*, 35:27921–27936, 2022.
633
- 634 Suyoung Lee, Myungsik Cho, and Youngchul Sung. Parameterizing non-parametric meta-
635 reinforcement learning tasks via subtask decomposition. *Advances in Neural Information Pro-*
636 *cessing Systems*, 36:43356–43383, 2023.
- 637 Xiang Li, Varun Belagali, Jinghuan Shang, and Michael S Ryoo. Crossway diffusion: Improv-
638 ing diffusion-based visuomotor policy via self-supervised learning. In *2024 IEEE International*
639 *Conference on Robotics and Automation (ICRA)*, pp. 16841–16849. IEEE, 2024.
640
- 641 Zhixuan Liang, Yao Mu, Mingyu Ding, Fei Ni, Masayoshi Tomizuka, and Ping Luo. Adaptdiffuser:
642 Diffusion models as adaptive self-evolving planners. *arXiv preprint arXiv:2302.01877*, 2023.
- 643 Ajay Mandlekar, Danfei Xu, Josiah Wong, Soroush Nasiriany, Chen Wang, Rohun Kulkarni, Li Fei-
644 Fei, Silvio Savarese, Yuke Zhu, and Roberto Martín-Martín. What matters in learning from offline
645 human demonstrations for robot manipulation. *arXiv preprint arXiv:2108.03298*, 2021.
646
- 647 Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models.
In *International conference on machine learning*, pp. 8162–8171. PMLR, 2021.

- 648 OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew,
649 Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider,
650 Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba,
651 and Lei Zhang. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019.
652
- 653 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
654 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-
655 low instructions with human feedback. *Advances in neural information processing systems*, 35:
656 27730–27744, 2022.
- 657 Tim Pearce, Tabish Rashid, Anssi Kanervisto, Dave Bignell, Mingfei Sun, Raluca Georgescu, Ser-
658 gio Valcarcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, et al. Imitating human
659 behaviour with diffusion models. *arXiv preprint arXiv:2301.10677*, 2023.
- 660 Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual
661 reasoning with a general conditioning layer. In *Proceedings of the AAAI conference on artificial*
662 *intelligence*, volume 32, 2018.
- 663 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
664 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
665 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 666
- 667 Mihir Prabhudesai, Anirudh Goyal, Deepak Pathak, and Katerina Fragkiadaki. Aligning text-to-
668 image diffusion models with reward backpropagation. *arXiv preprint arXiv:2310.03739*, 2023.
669
- 670 Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel
671 Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement
672 learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- 673 Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen,
674 and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine*
675 *learning*, pp. 8821–8831. Pmlr, 2021.
- 676
- 677 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov,
678 Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al.
679 A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- 680 Allen Z Ren, Justin Lidard, Lars L Ankile, Anthony Simeonov, Pulkit Agrawal, Anirudha Majum-
681 dar, Benjamin Burchfiel, Hongkai Dai, and Max Simchowitz. Diffusion policy policy optimiza-
682 tion. *arXiv preprint arXiv:2409.00588*, 2024.
- 683
- 684 Moritz Reuss, Maximilian Li, Xiaogang Jia, and Rudolf Lioutikov. Goal-conditioned imitation
685 learning using score-based diffusion policies. *arXiv preprint arXiv:2304.02532*, 2023.
- 686 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
687 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
688 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 689
- 690 Arne Schneuing, Yuanqi Du, Charles Harris, Arian Jamasb, Ilya Igashov, Weitao Du, Tom Blun-
691 dell, Pietro Lió, Carla Gomes, Max Welling, et al. Structure-based drug design with equivariant
692 diffusion models. *arXiv preprint arXiv:2210.13695*, 2022.
- 693 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
694 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
695 open large-scale dataset for training next generation image-text models. *Advances in Neural*
696 *Information Processing Systems*, 35:25278–25294, 2022.
- 697 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
698 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
699
- 700 Younggyo Seo, Danijar Hafner, Hao Liu, Fangchen Liu, Stephen James, Kimin Lee, and Pieter
701 Abbeel. Masked world models for visual control. In *Conference on Robot Learning*, pp. 1332–
1344. PMLR, 2023.

- 702 Xudong Shen, Chao Du, Tianyu Pang, Min Lin, Yongkang Wong, and Mohan Kankanhalli. Fine-
703 tuning text-to-image diffusion models for fairness. *arXiv preprint arXiv:2311.07604*, 2023.
704
- 705 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
706 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
707 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 708 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
709 learning using nonequilibrium thermodynamics. In *International conference on machine learn-*
710 *ing*, pp. 2256–2265. PMLR, 2015.
711
- 712 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
713 *preprint arXiv:2010.02502*, 2020.
- 714 Lingfeng Sun, Haichao Zhang, Wei Xu, and Masayoshi Tomizuka. Paco: Parameter-compositional
715 multi-task reinforcement learning. *Advances in Neural Information Processing Systems*, 35:
716 21495–21507, 2022.
- 717 Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient meth-
718 ods for reinforcement learning with function approximation. *Advances in neural information*
719 *processing systems*, 12, 1999.
720
- 721 Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*
722 *learning research*, 9(11), 2008.
723
- 724 Che Wang, Xufang Luo, Keith Ross, and Dongsheng Li. Vrl3: A data-driven framework for visual
725 deep reinforcement learning. *Advances in Neural Information Processing Systems*, 35:32974–
726 32988, 2022a.
- 727 Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy
728 class for offline reinforcement learning. *arXiv preprint arXiv:2208.06193*, 2022b.
729
- 730 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao
731 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.
732 *Advances in Neural Information Processing Systems*, 36, 2024.
- 733 Mengda Xu, Zhenjia Xu, Cheng Chi, Manuela Veloso, and Shuran Song. Xskill: Cross embodiment
734 skill discovery. In *Conference on Robot Learning*, pp. 3536–3555. PMLR, 2023.
735
- 736 Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua Tenenbaum, and Chuang
737 Gan. Prompting decision transformer for few-shot policy generalization. In *international confer-*
738 *ence on machine learning*, pp. 24631–24645. PMLR, 2022.
- 739 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
740 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
741 In *Conference on Robot Learning (CoRL)*, 2019.
- 742 Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn.
743 Gradient surgery for multi-task learning. *Advances in Neural Information Processing Systems*,
744 33:5824–5836, 2020.
745
- 746 Xudong Yu, Chenjia Bai, Haoran He, Changhong Wang, and Xuelong Li. Regularized conditional
747 diffusion model for multi-task preference alignment. *arXiv preprint arXiv:2404.04920*, 2024.
- 748 Yanjie Ze, Gu Zhang, Kangning Zhang, Chenyuan Hu, Muhan Wang, and Huazhe Xu. 3d diffusion
749 policy: Generalizable visuomotor policy learning via simple 3d representations. In *Proceedings*
750 *of Robotics: Science and Systems (RSS)*, 2024.
751
- 752 Qinqing Zheng, Amy Zhang, and Aditya Grover. Online decision transformer. In *international*
753 *conference on machine learning*, pp. 27042–27059. PMLR, 2022.
- 754 Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Open3D: A modern library for 3D data processing.
755 *arXiv:1801.09847*, 2018.

A DERIVATIONS

A.1 DERIVATION OF POLICY GRADIENT IN EQUATION (11)

Assume $p_\theta(\mathbf{a}_t^{0:K}|\mathbf{s}_t)r^\mathcal{T}(\mathbf{a}_t^0)$ and $\nabla_\theta p_\theta(\mathbf{a}_t^{0:K}|\mathbf{s}_t)r^\mathcal{T}(\mathbf{a}_t^0)$ are continuous (Fan et al., 2024), we have:

$$\begin{aligned}
\nabla_\theta J^\mathcal{T}(\theta) &= \nabla_\theta \sum_t \mathbb{E}_{p_\theta(\mathbf{a}_t^0|\mathbf{s}_t)} [r^\mathcal{T}(\mathbf{a}_t^0)] \\
&= \sum_t \left[\nabla_\theta \int r^\mathcal{T}(\mathbf{a}_t^0) \cdot p_\theta(\mathbf{a}_t^0|\mathbf{s}_t) d\mathbf{a}_t^0 \right] \\
&= \sum_t \left[\nabla_\theta \int r^\mathcal{T}(\mathbf{a}_t^0) \cdot \left(\int p_\theta(\mathbf{a}_t^{0:K}|\mathbf{s}_t) d\mathbf{a}_t^{1:K} \right) d\mathbf{a}_t^0 \right] \\
&= \sum_t \left[\int r^\mathcal{T}(\mathbf{a}_t^0) \cdot \nabla_\theta \log p_\theta(\mathbf{a}_t^{0:K}|\mathbf{s}_t) \cdot p_\theta(\mathbf{a}_t^{0:K}|\mathbf{s}_t) d\mathbf{a}_t^{0:K} \right] \\
&= \sum_t \left[\int r^\mathcal{T}(\mathbf{a}_t^0) \cdot \nabla_\theta \log \left(p_K(\mathbf{a}_t^K|\mathbf{s}_t) \prod_{k=1}^K p_\theta(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t) \right) \cdot p_\theta(\mathbf{a}_t^{0:K}|\mathbf{s}_t) d\mathbf{a}_t^{0:K} \right] \\
&= \sum_t \mathbb{E}_{p_\theta(\mathbf{a}_t^{0:K}|\mathbf{s}_t)} \left[r^\mathcal{T}(\mathbf{a}_t^0) \sum_{k=1}^K \nabla_\theta \log p_\theta(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t) \right].
\end{aligned} \tag{16}$$

A.2 DERIVATION OF LOSS FUNCTION IN EQUATION (12)

By using importance sampling approach, we can rewrite Eq. (16) as follows:

$$\sum_t \mathbb{E}_{p_{\theta_{\text{old}}}(\mathbf{a}_t^{0:K}|\mathbf{s}_t)} \left[r^\mathcal{T}(\mathbf{a}_t^0) \sum_{k=1}^K \frac{p_\theta(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t)}{p_{\theta_{\text{old}}}(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t)} \nabla_\theta \log p_\theta(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t) \right] \tag{17}$$

Then, we can get a new objective function corresponding to Eq. (17) as:

$$J_{\theta_{\text{old}}}^\mathcal{T}(\theta) = \max_\theta \sum_t \mathbb{E}_{p_{\theta_{\text{old}}}(\mathbf{a}_t^{0:K}|\mathbf{s}_t)} \left[r^\mathcal{T}(\mathbf{a}_t^0) \sum_{k=1}^K \frac{p_\theta(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t)}{p_{\theta_{\text{old}}}(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t)} \right] \tag{18}$$

Let $\rho_k(\theta, \theta_{\text{old}}) = \frac{p_\theta(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t)}{p_{\theta_{\text{old}}}(\mathbf{a}_t^{k-1}|\mathbf{a}_t^k, \mathbf{s}_t)}$ denote the probability ratio. Based on PPO (Schulman et al., 2017), we clip ρ_k and use the minimum between the clipped and unclipped ratios to derive a lower bound of the original objective (18), which serves as our final objective function:

$$J_{\text{clip}}^\mathcal{T}(\theta) = \max_\theta \sum_t \mathbb{E}_{p_{\theta_{\text{old}}}(\mathbf{a}_t^{0:K}|\mathbf{s}_t)} \left[r^\mathcal{T}(\mathbf{a}_t^0) \sum_{k=1}^K \min \left(\rho_k(\theta, \theta_{\text{old}}), \text{clip}(\rho_k(\theta, \theta_{\text{old}}), 1 + \epsilon, 1 - \epsilon) \right) \right] \tag{19}$$

To refine our pre-trained planner, we employ the negative of objective (19) as the loss function to facilitate reward maximization during fine-tuning.

A.3 DERIVATION OF LOSS FUNCTION IN EQUATION (14)

Directly computing and minimizing the NLL is difficult. However, we can derive an upper bound of Eq. (14) as follows:

$$\begin{aligned}
\mathbb{E}_{\mathbf{a}_\mu^0 \sim p_\mu} [-\log p_\theta(\mathbf{a}_\mu^0)] &\leq \mathbb{E}_{\mathbf{a}_\mu^0 \sim p_\mu} \left[\mathbb{E}_{q(\mathbf{a}_\mu^{1:K} | \mathbf{a}_\mu^0)} \left[-\log \frac{p_\theta(\mathbf{a}_\mu^{0:K})}{q(\mathbf{a}_\mu^{1:K} | \mathbf{a}_\mu^0)} \right] \right] \\
&= \mathbb{E}_{\mathbf{a}_\mu^0 \sim p_\mu} \left[\mathbb{E}_{q(\mathbf{a}_\mu^{1:K} | \mathbf{a}_\mu^0)} \left[-\log p(\mathbf{a}_\mu^K) - \sum_{k=1}^K \log \frac{p_\theta(\mathbf{a}_\mu^{k-1} | \mathbf{a}_\mu^k)}{q(\mathbf{a}_\mu^k | \mathbf{a}_\mu^{k-1})} \right] \right] \quad (20) \\
&= \mathbb{E}_{\mathbf{a}_\mu^0 \sim p_\mu} \left[\sum_{k=2}^K \mathbb{E}_{q(\mathbf{a}_\mu^k | \mathbf{a}_\mu^0)} D_{\text{KL}}[q(\mathbf{a}_\mu^{k-1} | \mathbf{a}_\mu^k, \mathbf{a}_\mu^0) \| p(\mathbf{a}_\mu^{k-1} | \mathbf{a}_\mu^k)] + \right. \\
&\quad \left. D_{\text{KL}}(q(\mathbf{a}_\mu^K | \mathbf{a}_\mu^0) \| p(\mathbf{a}_\mu^K)) - \mathbb{E}_{q(\mathbf{a}_\mu^1 | \mathbf{a}_\mu^0)} [\log p_\theta(\mathbf{a}_\mu^0 | \mathbf{a}_\mu^1)] \right]
\end{aligned}$$

Following previous work (Ho et al., 2020), the optimization of the bound can be simplified as:

$$\arg \min_{\theta} \frac{1}{2\sigma_q^2(k)} \frac{(1 - \alpha_k)^2}{(1 - \bar{\alpha}_k)\alpha_k} \|\epsilon(\mathbf{a}_\mu^k, k) - \epsilon_\theta(\mathbf{a}_\mu^k, k)\|^2 \quad (21)$$

where:

$$\sigma_q^2(k) = \frac{(1 - \alpha_k)(1 - \bar{\alpha}_{k-1})}{1 - \bar{\alpha}_k} \quad (22)$$

Here, $\epsilon_\theta(\mathbf{a}_\mu^k, k)$ is a noise model that learns to predict the source noise $\epsilon(\mathbf{a}_\mu^k, k)$ which determines \mathbf{a}_μ^k from \mathbf{a}_μ^0 .

B THE DETAILS OF SODP

B.1 DIFFUSION POLICY

We use diffusion policy (Chi et al., 2023) to generate future actions. For any given time step t , the model uses the most recent T_o steps of states as input to generate the next T_p action steps. Then, the first T_a steps of these generated actions are executed in the environment without re-planning. In our experiments, we use $T_p = 12, T_o = 2, T_a = 8$ for Meta-World and $T_p = 4, T_o = 2, T_a = 3$ for Adroit.

We employ a CNN-based diffusion policy as our noise model, utilizing a U-net architecture that incorporates Feature-wise Linear Modulation (FiLM) (Perez et al., 2018) to condition on historical states. The implementation is based on the code from <https://github.com/CleanDiffuserTeam/CleanDiffuser>, and we use their default hyper-parameters. For Adroit, we use a simplified backbone provided by Simple DP3 (<https://github.com/YanjieZe/3D-Diffusion-Policy>), which removes some components in the U-net.

B.2 IMPLEMENTATION DETAILS

The pseudo-code of SODP is given in Alg. 1. We describe details of pre-training and fine-tuning as follows:

- For pretraining, we use cosine schedule for β_k (Nichol & Dhariwal, 2021) and set diffusion steps $K = 100$. We pre-train the model for $5e^5$ steps in Meta-World and $3e^3$ steps in Adroit.
- For fine-tuning, we use DDIM (Song et al., 2020) with 10 sampling steps and $\eta = 1$. We fine-tune each task for $1e^6$ steps in Meta-World and $3e^3$ steps in Adroit. Following DPOK (Fan et al., 2024), we perform $p_{\text{step}} \in \{10, 30\}$ gradient steps per episode. We set discount factor $\gamma = 1$ for all tasks.
- We set $N_{\text{init}} \in \{10, 20\}$ for approximating target distribution and $\lambda = 1.0$ as the BC weight coefficient.
- Batch size is set to 256 for both pre-training and fine-tuning.

- We use Adam optimizer (Kingma, 2014) with default parameters for both pre-training and fine-tuning. Learning rate is set to $1e^{-4}$ for pretraining and $1e^{-5}$ for fine-tuning with exponential decay.

Algorithm 1: SODP: Two-stage framework for learning from sub-optimal data

Input: diffusion planner θ , N downstream tasks \mathcal{T}_i , multi-task sub-optimal data $D = \cup_{i=1}^N \mathcal{D}_{\mathcal{T}_i}$, target buffer $\mathcal{B}_{\text{target}}$, replay buffer \mathcal{B} , episode length L , pre-train N_{PT} and fine-tune N_{FT} steps

```

// pre-training model with sub-optimal data
for  $t = 1, \dots, N_{\text{PT}}$  do
  Sample  $(s, \mathbf{a}) \sim D$ , diffusion time step  $k \sim \text{Uniform}(\{1, \dots, K\})$ , noise  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ ;
  Update  $\theta$  using the loss function (7);
// fine-tuning model for downstream tasks
for  $\mathcal{T}_i \in [\mathcal{T}_1, \dots, \mathcal{T}_N]$  do
  Initialization:  $\theta \leftarrow \theta_{\text{PT}}$ ;  $\mathcal{B}, \mathcal{B}_{\text{target}} \leftarrow$  Rollout  $N_{\text{init}}$  proficient trajectories using  $\theta$ ;
  for  $t = 1, \dots, N_{\text{FT}}$  do
    while not end of the episode do
      Obtain samples  $\mathbf{a}_t^{0:K} \sim p_{\theta}(\mathbf{a}_t^{0:K} | s_t)$ ;
      Execute the first  $T_a$  steps and get reward  $r(\mathbf{a}_t^0)$ ;
       $\mathcal{B} \leftarrow \mathcal{B} \cup (s_t, \mathbf{a}_t^{0:K}, r(\mathbf{a}_t^0))$ ;
       $s_t \leftarrow s_{t+T_a}, t \leftarrow t + T_a$ ;
    // approximate target policy  $\mu$ 
    if proficient then
       $\mathcal{B}_{\text{target}} \leftarrow \mathcal{B}_{\text{target}} \cup \{\mathbf{a}_t^{0:K} | t \in \{0, T_a, \dots, L\}\}$ 
      Compute  $\mathcal{L}_{\text{imp}}^{\mathcal{T}_i}$  using batches from  $\mathcal{B}$  according to Eq. (12);
      Compute  $\mathcal{L}_{\text{BC}}^{\mathcal{T}_i}$  using batches from  $\mathcal{B}_{\text{target}}$  according to Eq. (14);
      Update  $\theta$  using the loss function (15);

```

C EXTENDED RESULTS

In this section, we provide our full experimental results:

1. Baselines incorporating our online interaction trajectories as supplementary training data.
2. Ablation studies evaluating various fine-tuning strategies.
3. Analysis of the impact of pre-training dataset quality.
4. Generalizability to previously unseen tasks.
5. Evaluation on image-based Meta-World tasks across 10 environments.

C.1 AUGMENTED TRAINING DATA FOR MTDIFF AND HARMODT

To isolate the influence of data quantity, we conducted fine-tuning for 100k steps per task using SODP, collecting online interaction samples during the fine-tuning process. These samples were then incorporated as a supplementary dataset alongside the original data, expanding the dataset size from 50M to $50\text{M} + 100\text{k} \times 50$. Subsequently, we trained both MTDIFF and HarmoDT on this augmented dataset to ensure consistent data usage across our method and the baseline methods. The experimental results are presented in Table 3, demonstrating that our method continues to outperform the baseline methods under this configuration. For MTDIFF, we employed the default parameters provided by the authors. However, a performance decline was observed on these new datasets, likely due to the increased presence of inferior data introduced during the online interaction phase.

Table 3: Average success rate using augmented sub-optimal data.

Method	Meta-World 50 Tasks
MTDIFF-P	27.06 \pm 0.42
HarmoDT-F	57.37 \pm 0.34
SODP (ours)	59.26\pm0.18

C.2 ABLATION STUDIES EXAMINING OTHER FINE-TUNING APPROACHES

To demonstrate the effectiveness of our online fine-tuning approach, we compare it with two alternative fine-tuning methods: (i) **SODP_off**, which involves fine-tuning using high-quality offline data, and (ii) **SODP_off_scratch**, which performs direct training with high-quality data without pre-training. Specifically, we fine-tuned the pre-trained models for 100k steps across five tasks, collecting 200 successful episodes (equivalent to 100k steps) for each task. These datasets were then used to independently train five models in an offline setting, utilizing the same loss function as in Eq. (15) (**SODP_off**). Additionally, to investigate the impact of pre-training on offline fine-tuning, we trained the model directly without pre-training (**SODP_off_scratch**).

The experimental results, presented in Table 4, report the success rates averaged over three seeds. Without pre-training, the model lacks the necessary action priors to efficiently identify high-reward action distributions. Furthermore, directly fine-tuning with high-quality offline data proves insufficient, as static reward labels may fail to provide adequate guidance in dynamic environments, hindering the model’s ability to facilitate efficient exploration.

Table 4: Average success rate for different fine-tuning approaches.

Tasks	SODP_off	SODP_off scratch	SODP
button-press-topdown	58.67 \pm 0.03	40.67 \pm 0.08	60.67 \pm 0.03
hammer	71.33 \pm 0.05	13.33 \pm 0.06	73.33 \pm 0.03
handle-pull-side	60.67 \pm 0.03	42.67 \pm 0.08	81.67 \pm 0.07
peg-insert-side	25.33 \pm 0.03	0.0 \pm 0.0	32.67 \pm 0.06
handle-pull	66.67 \pm 0.03	31.33 \pm 0.06	75.33 \pm 0.04
Average success rate	56.53 \pm 0.18	25.6 \pm 0.18	64.73 \pm 0.19

To highlight the importance of modeling the diffusion process as a MDP for reward fine-tuning, we consider an alternative approach that directly applies BC during fine-tuning, using only Eq. (14) as the loss function. As shown in Table 5, directly using BC results in poorer performance, as BC lacks reward labels to effectively guide exploration. While BC during the fine-tuning phase enables access to dynamic actions, it is limited to ‘imitation’ rather than ‘evolution,’ as the model is unable to differentiate between good and bad actions.

Table 5: Average success rate for directly BC during fine-tuning.

Task	Directly BC	SODP
button-press-topdown	51.3 \pm 0.05	60.7 \pm 0.03
basketball	21.3 \pm 0.03	41.2 \pm 0.16
stick-pull	26.7 \pm 0.08	50.5 \pm 0.04

C.3 PRE-TRAINING USING NEAR-OPTIMAL DATA

To evaluate the impact of pre-training data quality on fine-tuning performance, we modified the near-optimal dataset provided by He et al. (2024) by retaining only the last 50% of the data. This modification ensured that the total number of transitions remained the same as the sub-optimal data used in the main paper, while significantly increasing the proportion of expert trajectories. We refer to this modified dataset as near-optimal data and pre-trained a model on the Meta-World 10 tasks. Subsequently, we followed the same fine-tuning procedure outlined in the main paper to fine-tune the model on each task. The experimental results are presented in Table 6. Incorporating more optimal data during the pre-training stage leads to better performance, as the model gains more priors about the optimal action distributions.

Table 6: Average success rate achieved after fine-tuning models pre-trained on different datasets.

Tasks	Sub-optimal dataset	Near-optimal dataset
basketball	52.67 \pm 0.03	80.67 \pm 0.03
button-press	88.00 \pm 0.02	89.33 \pm 0.03
dial-turn	80.67 \pm 0.02	74.00 \pm 0.04
drawer-close	100.00 \pm 0.00	100.00 \pm 0.00
peg-insert-side	62.67 \pm 0.02	84.67 \pm 0.02
pick-place	36.67 \pm 0.03	59.33 \pm 0.03
push	33.33 \pm 0.03	50.67 \pm 0.03
reach	68.67 \pm 0.05	95.33 \pm 0.01
sweep-into	60.67 \pm 0.03	75.33 \pm 0.01
window-open	69.33 \pm 0.04	100.0 \pm 0.00
Average success rate	65.27 \pm 0.21	80.93 \pm 0.16

C.4 FINE-TUING ON UNSEEN TASKS

To evaluate the generalizability of SODP, we conduct experiments to fine-tune the model on tasks that were not included in the pre-training dataset. We pre-train a model on the MT-10 dataset (SODP_mt10) and fine-tune it on three tasks that are not present in the pre-training data. Additionally, to investigate the advantages of pre-training on a multi-task dataset versus a single-task dataset, we compare SODP_mt10 with a variant that is pre-trained solely on the *basketball* dataset (SODP_bas). As shown in Table 7, pre-training on multi-task data enhances generalizability to unseen tasks, as multi-task data provide a broader range of action distribution priors compared to single-task data.

Table 7: Average success rate achieved after fine-tuning on unseen tasks.

Unseen tasks	SODP_mt10	SODP_bas
drawer-open	34.7 \pm 0.06	0.0 \pm 0.0
plate-slide-side	55.3 \pm 0.33	0.0 \pm 0.0
handle-pull-side	71.3 \pm 0.13	0.0 \pm 0.0

C.5 EXPERIMENTS IN IMAGE-BASED META-WORLD

To further validate the scalability of SODP in handling high-dimensional observations, we conduct experiments on image-based Meta-World 10 tasks. Since no existing image-based sub-optimal dataset for Meta-World is available, we collect data for the 10 tasks by training separate SAC agents for each task, as done in He et al. (2024), and rendering the environments to obtain image data. We then follow the same procedure as in Adroit to convert the images into point clouds and use the DP3 encoder to extract visual features. For comparison, we consider the following baselines: DP3 and MTDIFF_3D, an extended variant of MTDIFF that employs the same 3D visual encoder used in SODP. The experimental results are presented in Table 8, demonstrating the generalizability of our method to complex inputs.

Table 8: Average success rate of image-based MT-10 tasks.

Methods	Success rate
DP3	32.6 \pm 0.23
MTDIFF_3D	38.0 \pm 0.82
SODP	47.5 \pm 0.18

D THE DETAILS OF BASELINES

We describe the details of baselines used for comparison in our experiments. For Meta-World, we consider following baselines:

- **MTSAC.** The one-hot encoded task ID is incorporated into the original SAC as an additional input.

- 1026 • **MTBC.** The actor network is modeled using a 3-layer MLP with Mish activation. In training and
1027 inference, the scalar task ID is processed through a separate 3-layer MLP with Mish activation to
1028 produce a latent variable z . The input to the actor network is then formed by concatenating the
1029 original state with this latent variable z
- 1030 • **MTIQL.** Similar to MTBC, the actor network incorporates the task ID through a task-aware em-
1031 bedding. A multi-head critic network is employed to estimate the Q -values for each task, with
1032 each head being parameterized by a 3-layer MLP using Mish activation.
- 1033 • **MTDQL.** Similar to MTIQL, a multi-head critic network is utilized to predict the Q -value for
1034 each task, and the original diffusion actor is extended with an additional task ID input.
- 1035 • **MTDT.** The task ID is embedded into a latent variable z of size 12. This latent variable is then
1036 concatenated with the raw state to form the input tokens.
- 1037 • **Prompt-DT.** Actions are generated based on trajectory prompts and the reward-to-go. A GPT-2
1038 transformer model is utilized as the noise network.
- 1039 • **MTDIFF.** Actions are generated by a GPT-based diffusion model that incorporates prompt learn-
1040 ing to capture task knowledge. MTDIFF considers a variant: MTDIFF-ONEHOT, which replaces
1041 the prompt with a one-hot task ID. We borrow the official codes from [https://github.com/
1042 tinnerhrhe/MTDiff](https://github.com/tinnerhrhe/MTDiff) and use their default hyper-parameters.
- 1043 • **HarmoDT.** Incorporate trainable task-specific masks to address gradient conflict by identifying
1044 an optimal harmony subspace of parameters for each task. There are three variants of HarmoDT:
1045 HarmoDT-R, which keeps task masks unchanged; HarmoDT-F and HarmoDT-M utilize differ-
1046 ent methods to weight masks. We borrow the official codes from [https://github.com/
1047 charleshsc/HarmoDT](https://github.com/charleshsc/HarmoDT) and use their default hyper-parameters.

1048 For Adroit, we consider following baselines:

- 1050 • **BCRNN.** A variant of BC that models the policy network as an RNN. The network is trained on
1051 temporal sequences of length H , denoted as $(s_t, a_t, \dots, s_{t+H}, a_{t+H})$, to predict action sequences
1052 based on historical states.
- 1053 • **IBC.** BC is represented as a conditional energy-based modeling problem, where implicit policies
1054 are trained to imitate expert demonstrations.
- 1055 • **Diffusion Policy.** The generation of robot behaviors is formulated as a conditional denoising dif-
1056 fusion process, where the diffusion model predicts action sequences based on given observations
1057 as conditions.
- 1058 • **DP3.** Diffusion Policy is extended by incorporating 3D visual representations. The 3D scenes
1059 from the environment are represented as point clouds, which are then cropped and downsampled
1060 to reduce redundant information. These processed point clouds are passed through an MLP to
1061 generate visual representations, which serve as conditions for the diffusion models.

1062 For image-based Meta-World, we extended MTDIFF by integrating the same 3D visual encoder
1063 used in SODP to extract visual features from input point clouds.

1064 E VARIANTS OF SODP

1065 In Eq. (15), we introduce a BC regularization term to preserve the pre-trained knowledge and demon-
1066 strate its effectiveness compared to two existing regularization approaches presented in DPOK (Fan
1067 et al., 2024) and DLPO (Chen et al., 2024). Specifically, the regularization term \mathcal{L}_{KL} used in DPOK
1068 is expressed as:

$$1069 \mathcal{L}_{\text{KL}}(\theta) = \sum_{k=1}^K \text{KL}(p_{\theta}(x_{k-1}|x_k) || p_{\text{pre}}(x_{k-1}|x_k)). \quad (23)$$

1070 And the regularization term \mathcal{L}_{PL} used in DLPO is expressed as:

$$1071 \mathcal{L}_{\text{PL}}(\theta) = \mathbb{E}_{k \sim [1, K], p_{\theta}(x_{1:K})} \left[\|\epsilon(x_k, k) - \epsilon_{\theta}(x_k, k)\|^2 \right]. \quad (24)$$

1080 These methods can be considered as different approaches to selecting the target policy in line with
1081 our analysis and can be seen as variants of Eq. (14), where DPPO selects $\mu = \theta_{\text{pre-train}}$ and DLPO
1082 selects $\mu = \theta$. The rationale behind their selection is based on the assumption that $\theta \approx \theta_{\text{pre-train}} \approx \theta^*$.
1083 This assumption is reasonable in text-to-image or text-to-speech tasks, as the pre-trained models
1084 they used are already strong and perform exceptionally well even without fine-tuning. However, this
1085 assumption does not apply to our pre-trained planner, as the model is trained on sub-optimal data.
1086 As a result, as shown in Section 5.3, these regularization methods may lead the pre-trained planner
1087 to be stuck in inferior regions, limiting its ability to improve performance.

1088 F COMPARISON TO DPPO

1089 We summarize some similarities and differences between our work and the concurrent work
1090 DPPO (Ren et al., 2024) as follows:

- 1091 • Both DPPO and our approach formulate the diffusion policy denoising process as an MDP and
1092 use policy gradients to fine-tune the model for higher environment rewards.
- 1093 • DPPO demonstrated that reward-based RL fine-tuning promotes effective exploration, which is
1094 consistent with our observations.
- 1095 • While DPPO requires task-specific expert demonstrations for pre-training, our method pre-trains
1096 a foundation model capable of capturing useful behavior patterns from multi-task inferior data.
- 1097 • We show that directly fine-tuning the pre-trained planner without any regularization, as done in
1098 DPPO, fails in the multi-task setting. We further analyze the limitations of current regularization
1099 methods and propose a novel BC regularization term. By employing our regularizer, the pre-
1100 trained model achieves higher success rates after fine-tuning.
- 1101 • Unlike DPPO, we don't employ advantage estimator.

1102 G SINGLE-TASK PERFORMANCE

1103 We evaluate the performance for each task for 50 episodes. We report the average evaluated return
1104 of pre-trained and fine-tuned models in Table 9.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

Table 9: Evaluated return of SODP pre-trained model and fine-tuned model for each task in MT50-rand. We report the mean and standard deviation for 50 episodes for each task.

Tasks	Return of pre-trained model	Return of fine-tuned model
basketball-v2	133.5 ± 100.7	2347.1 ± 580.8
bin-picking-v2	96.8 ± 23.9	602.7 ± 72.8
button-press-topdown-v2	1405 ± 20.3	1679 ± 25.9
button-press-v2	1397 ± 15.6	2452.7 ± 89.3
button-press-wall-v2	1375 ± 10.58	2524.7 ± 18.1
coffee-button-v2	293.2 ± 12.1	451.5 ± 14.4
coffee-pull-v2	39.5 ± 6.2	117.9 ± 23.3
coffee-push-v2	33.8 ± 6.1	273.3 ± 36.6
dial-turn-v2	1217.7 ± 239.3	1557.3 ± 226.7
disassemble-v2	237 ± 117.2	502 ± 164.6
door-close-v2	3347.7 ± 124.9	4116.3 ± 118.6
door-lock-v2	1042.3 ± 94.9	2491 ± 79.5
door-open-v2	2036.3 ± 79.3	2460.3 ± 57.7
door-unlock-v2	1335 ± 46.9	2257.7 ± 323.0
hand-insert-v2	85.9 ± 56.7	449.5 ± 54.9
drawer-close-v2	2468.3 ± 167.2	3953.7 ± 214.3
drawer-open-v2	1656 ± 45.6	2489.7 ± 188.4
faucet-open-v2	2728.7 ± 424.1	4094.7 ± 290.3
faucet-close-v2	2156.7 ± 113.6	3772 ± 70.1
handle-press-side-v2	1919.7 ± 449.5	3478.3 ± 98.0
handle-press-v2	2216.3 ± 182.0	3415.7 ± 221.6
handle-pull-side-v2	1351.7 ± 119.0	2665.7 ± 243.9
handle-pull-v2	1510.7 ± 111.6	2734 ± 64.3
lever-pull-v2	650.7 ± 32.5	1068.8 ± 110.4
peg-insert-side-v2	300.3 ± 122.2	1969.7 ± 237.6
pick-place-wall-v2	596.7 ± 10.6	1175.7 ± 150.1
pick-out-of-hole-v2	38.5 ± 6.3	106.7 ± 7.9
reach-v2	2664.3 ± 77.5	3083.7 ± 149.7
push-back-v2	55.8 ± 26.7	350.9 ± 30.3
push-v2	46.8 ± 35.9	148.9 ± 43.9
pick-place-v2	3.9 ± 0.2	5.9 ± 1.8
plate-slide-v2	1268.7 ± 75.8	2862 ± 234.5
plate-slide-side-v2	826.4 ± 54.3	1929.7 ± 104.5
plate-slide-back-v2	795.8 ± 62.9	1587.7 ± 125.0
plate-slide-back-side-v2	626.7 ± 36.9	1541.3 ± 78.5
soccer-v2	863.8 ± 159.5	1234.2 ± 225.8
push-wall-v2	175.2 ± 35.0	471.7 ± 73.9
shelf-place-v2	260.4 ± 111.5	785.1 ± 81.5
sweep-into-v2	621.0 ± 132.9	1282.3 ± 135.9
sweep-v2	442.3 ± 83.0	1081.7 ± 58.0
window-open-v2	1342.3 ± 60.1	2474.3 ± 266.4
window-close-v2	1087.7 ± 78.9	1816.7 ± 166.3
assembly-v2	282.5 ± 4.3	446.1 ± 26.9
button-press-topdown-wall-v2	1374 ± 16.1	1702.7 ± 71.5
hammer-v2	1678.3 ± 52.5	1907.7 ± 25.4
peg-unplug-side-v2	34.2 ± 2.9	52.9 ± 4.6
reach-wall-v2	3373.7 ± 41.7	3839.7 ± 69.3
stick-push-v2	412.9 ± 95.8	833.5 ± 99.1
stick-pull-v2	1977 ± 155.9	3116.3 ± 58.1
box-close-v2	692.3 ± 22.8	1300.1 ± 53.2