

# SkillDiffuser: Interpretable Hierarchical Planning via Skill Abstractions in Diffusion-Based Task Execution

Zhixuan Liang<sup>1,3</sup> Yao Mu<sup>1,3</sup> Hengbo Ma<sup>2</sup>  
Masayoshi Tomizuka<sup>2</sup> Mingyu Ding<sup>2†</sup> Ping Luo<sup>1,3†</sup>

<sup>1</sup>The University of Hong Kong <sup>2</sup>University of California, Berkeley <sup>3</sup>Shanghai AI Laboratory

{zxliang, ymu, pluo}@cs.hku.hk {hengbo\_ma, tomizuka, myding}@berkeley.edu

<https://skilldiffuser.github.io/>

## Abstract

Diffusion models have demonstrated strong potential for robotic trajectory planning. However, generating coherent trajectories from high-level instructions remains challenging, especially for long-range composition tasks requiring multiple sequential skills. We propose SkillDiffuser, an end-to-end hierarchical planning framework integrating interpretable skill learning with conditional diffusion planning to address this problem. At the higher level, the skill abstraction module learns discrete, human-understandable skill representations from visual observations and language instructions. These learned skill embeddings are then used to condition the diffusion model to generate customized latent trajectories aligned with the skills. This allows generating diverse state trajectories that adhere to the learnable skills. By integrating skill learning with conditional trajectory generation, SkillDiffuser produces coherent behavior following abstract instructions across diverse tasks. Experiments on multi-task robotic manipulation benchmarks like Meta-World and LOReL demonstrate state-of-the-art performance and human-interpretable skill representations from SkillDiffuser. More visualization results and information could be found on our [website](#).

## 1. Introduction

Recent research [6, 7, 18, 19] has demonstrated diffusion models’ superior generative capabilities compared to previous models that help enhance reinforcement learning across various dimensions, including the generation of action trajectories [2, 19], policy representation [5, 50], and data synthesis [14, 23]. However, their ability to generate coherent trajectories for intricate tasks still poses challenges in terms of performance and generalizability, as these tasks often re-

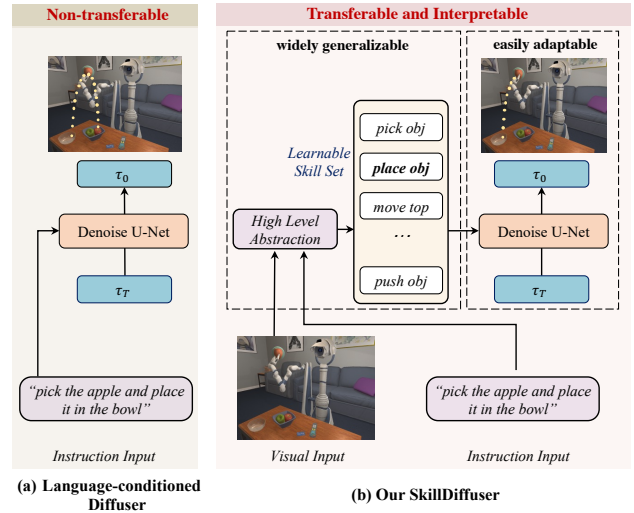


Figure 1. **Comparison of SkillDiffuser and previous language conditioned diffusers.** SkillDiffuser utilizes high-level abstraction to translate visual observations and language instructions into human understandable skills with language grounding. It then enables the low-level diffusion model condition on these skills, not only improving the execution performance of multi-step composition tasks but greatly enhancing the generalization and adaptability of the framework.

quire the fulfillment of abstract instructions that consist of numerous coordination-intensive sequential steps.

Previous approaches [2, 14], such as Decision Diffuser, aim to tackle this challenge by decomposing complex tasks into simpler sub-skills, organized within a predefined skill library. These methods rely on conditioning the diffusion model with one-hot skill representations to generate trajectories for each of these sub-skills. Nevertheless, these techniques encounter difficulties in attempting to autonomously learn end-to-end from a wide range of datasets, which hinders their scalability and ability to achieve end-to-end learning [1]. In addition, without explicitly learning reusable skills, models cannot capture intricate inter-step dependen-

<sup>†</sup>Corresponding authors.

cies and constraints, yielding fragmented and illogical trajectories. Decomposing ambiguous instructions into learnable sub-goals and skills can better enable models to follow logical step sequences, respect task structures, and transfer common procedures (*e.g.* reusable skill abstraction and adaptable skill-based diffusion) between different tasks. This skill-centric paradigm paves the way for diffusion models that can interpret and execute elaborate, abstract instructions necessitating numerous sequential steps.

Inspired by the above observations, we propose SkillDiffuser, a hierarchical planning framework unifying high-level skill learning with low-level conditional diffusion-based execution. As shown in Fig. 1, compared with previous language-conditioned diffusion policies, SkillDiffuser is able to interpret and execute complex instructions end-to-end with higher transferability. SkillDiffuser induces interpretable sub-latent goals by learning reusable skills tailored to the task instructions. The framework conditions a diffusion model on these learned skills to generate customized, coherent trajectories aligned with the overall objectives. By integrating hierarchical skill decomposition with conditional trajectory generation, SkillDiffuser achieves consistent, skill-oriented behavior without relying on a predefined skill library. Moreover, SkillDiffuser is designed to operate solely on visual observations, eliminating the need for robot proprioception (*i.e.* fully observed states). This end-to-end methodology, featuring learnable skills, enables SkillDiffuser to execute abstract instructions across a variety of tasks efficiently.

SkillDiffuser works as follows. **1)** A skill predictor with vector quantization [45] is used for high-level skill learning to distill tasks into discrete and interpretable skills. Rather than forecasting skill durations, we adopt a fixed prediction horizon – predicting skills at regular intervals. This horizon-based discretization process seamlessly integrates visual and linguistic inputs into a cohesive skill set guiding the low-level diffusion model. **2)** For skill-based trajectory generation, we utilize a classifier-free diffusion model as policy, with skills directly embedded as guidance. This setting allows for generating multi-modal state trajectories aligned with skill specifications while avoiding overfitting to a closed dataset. **3)** To enable action inference from predicted states, we train an inverse dynamics network to decode motions between two consecutive frames generated. By separating state prediction from motion decoding, SkillDiffuser yields a fully adaptive framework for directing diverse embodiments via transferable state-space plans.

We evaluate the model’s performance in both skill learning and multi-task planning on the LOReL [29] Sawyer and Meta-World [51] datasets, with crucial experimental settings by considering robotic agents in real-world scenarios must operate with limited state information, chiefly visual observations. Furthermore, we present success rates on un-

seen compositional tasks, the reusability and visualizations of learned skills to illustrate the model does have the ability to abstract high-level skills that are not only effective but also understandable to humans, bringing us closer to intelligent agents that acquire skills in a direct manner.

Our contributions are three-fold: **1)** We propose an end-to-end hierarchical planning framework via skill learning for sub-goal abstraction; **2)** We adopt a classifier-free diffusion model conditioned on learned skills to generate skill-oriented transferable state trajectories; **3)** We demonstrate state-of-the-art performance on complex benchmarks and provide interpretable visualizations of human-understandable skill representations.

## 2. Related Works

### 2.1. Imitation Learning and Multi-task Learning

Imitation learning (IL) has evolved from foundational behavioral cloning to sophisticated multi-task learning frameworks. With traditional approaches relying on supervised learning from expert demonstrations [33, 37, 38], recent advancements have shifted towards learning the reward [16] or Q-function [12] from expert data, enhancing the ability to mimic complex behaviors. A new challenge lies in multi-task IL [41], where imitators are trained across varied tasks, aiming for generalization to new scenarios with task specifications ranging from vector states [28] to vision and language descriptions [8, 11, 13, 48].

Multi-task learning approaches often leverage shared representations to learn a spectrum of tasks simultaneously, enhancing the flexibility and efficiency of the learning process. The Meta-World benchmark [51] assesses multi-task and meta reinforcement learning, highlighting the need for algorithms capable of rapid adaptation. Building on this, the Prompting Decision Transformer [49] showcases few-shot policy generalization using task-specific prompts. And diffusion policy has also been explored in multi-task settings [14], which shows proficiency in generating diverse behaviors across tasks. However, unlike methods that leverage state inputs [14, 49] or access robot proprioception [30], SkillDiffuser uses raw visual inputs only.

### 2.2. Skill Discovery and Hierarchical Learning

Skill learning, the process by which robots acquire new abilities or refine existing ones, is gaining increasing attention due to its pivotal role in enabling autonomous systems to adapt to new tasks, improve the performance over time, and interact naturally with humans and complex environments. Traditionally, this domain was influenced by hand-crafted features and expert demonstrations [32].

With the development of deep learning, Eysenbach *et al.* [9] and A. Sharma *et al.* [40] investigated skill discovery in learning methods, achieving policies conditioned

on learned latent variables and maintaining consistent skill codes throughout trajectories. In the domain of skill learning through language instructions, LISA [13] stands out by sampling multiple skills per trajectory, integrating language conditioning in a unique manner.

Our SkillDiffuser follows this way to extract sub-skills of each instruction at the higher-level. But differently, SkillDiffuser employs an adaptable diffusion policy at the lower-level to condition on different sub-skills to generate different actions, which formulates a creative hierarchical planning framework also advancing research on hierarchical techniques of reinforcement learning [22, 27, 52].

### 2.3. Planning with Diffusion Model

Diffusion models [18] make great breakthroughs in image synthesis [18] recent years and has shown promising results in various generative applications. A seminal work that performs planning with diffusion directly is Diffuser [19], which laid the groundwork for using diffusion models in behavior synthesis. Then, a branch of this kind planning methods achieved state-of-the-art performance in a variety of decision tasks [3, 5, 7, 23, 31]. Among them, the work done by Chi *et al.* in [5] introduces the concept of learning the gradient of the action-distribution score function and iteratively optimizing it, demonstrating its significant potential in visuomotor policy learning. These works have further extended this direction and strengthened the versatility and generalization of diffusion-based planners.

Our approach is inspired by classifier-free diffusion guidance [17], which offers a significant advantage over classifier-guided models [6]. By adopting the classifier-free approach, we can circumvent the challenges associated with training a reward model and  $Q$ -function, which are particularly cumbersome in many real-world planning scenarios of which the complexity is very high. Recent studies have also extended this direction, which use conditional diffusion models to generate customized trajectories. Decision Diffuser [2] is an example which is designed to generate trajectories from a predefined skill library. However, unlike our method, it can't autonomously learn skill abstractions in an end-to-end fashion, which makes it difficult to scale to more tasks. This highlights the necessity for diffusion models with dynamic, learnable skill abstractions, facilitating complex instruction execution.

## 3. Preliminary

### 3.1. Planning with Diffusion over States

As introduced in previous works [2, 19], diffusion model is a promising tool to address the problem of planning in reinforcement learning which is cast as a Markov Decision Process (MDP) [34]. Within the MDP framework  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma)$ , the planning policy aims to iden-

tify an optimal action sequence  $\mathbf{a}_{0:T}^*$  that maximizes the expected cumulative rewards over a finite time horizon  $T$ , governed by the state transition dynamics  $\mathcal{T}$  and reward function  $\mathcal{R}$ .  $\mathcal{S}$  is the state space and  $\mathcal{A}$  is the action space.

By treating the state trajectory as sequence data  $\tau$ , with sequence modeling, diffusion probabilistic models conceptualize planning as an iterative denoising process. The model progressively refines trajectories by reversing a forward diffusion process that is modeled as a Gaussian process, whereby noise is incrementally added to the data, denoted as  $p_\theta(\tau^{i-1} | \tau^i)$ . Training involves minimizing the ELBO of the data's negative log-likelihood, similar to variational Bayesian inference, with the optimization objective:

$$\theta^* = \arg \min_{\theta} -\mathbb{E}_{\tau^0} [\log p_\theta(\tau^0)], \quad (1)$$

where  $p(\tau^N)$  is a standard normal distribution and  $\tau^0$  denotes noiseless sequence data.

For practical implementation, a simplified surrogate loss function is proposed in [18], focusing on predicting the Gaussian mean of the reverse diffusion step:

$$\mathcal{L}_{\text{denoise}}(\theta) = \mathbb{E}_{i, \tau^0 \sim q, \epsilon \sim \mathcal{N}} [\|\epsilon - \epsilon_\theta(\tau^i, i)\|^2]. \quad (2)$$

### 3.2. Classifier-free Diffusion Guidance

On the basis of unconditioned diffusion-based method, in the realm of offline reinforcement learning, a critical endeavor is to generate trajectories with the highest reward-to-go. With the flourishing development of conditional diffusion models [6], classifier-guided approaches embark on this by infusing specific trajectory information (encoded in  $\mathbf{y}(\tau)$ ), such as the return  $\mathcal{J}(\tau^0)$  or designated constraints, into the diffusion process:

$$q(\tau^i | \tau^{i-1}), \quad p_\theta(\tau^{i-1} | \tau^i, \mathbf{y}(\tau)). \quad (3)$$

With assumptions specified in [10], we have

$$\tau^{i-1} \sim \mathcal{N}(\mu_\theta + \alpha \Sigma \nabla_\tau \log p(\mathbf{y}(\tau^i)), \Sigma), \quad (4)$$

where  $\alpha$  is a hyperparameter that adjusts the guidance strength,  $\Sigma$  is the specified covariance of the noise and  $\mu_\theta$  is the learned mean value of noise in unconditional diffusion.

However, the classifier-guided diffusion model requires an accurate estimation of guidance gradient based on the trajectory classifier  $\mathbf{y}(\tau)$ , which may not be feasible and need to introduce a separate dynamic programming procedure to estimate a  $Q$ -function in the training process.

Thus, classifier-free guidance offer an alternative, which augments the trajectory generation process with a guidance signal that amplifies the features of high-reward or optimal characteristics, *i.e.*  $\mathbf{y}(\tau)$ , that are implicitly present in the data. Mathematically, the noise to add during the reverse denoising process is:

$$\hat{\epsilon} = \epsilon_\theta(\tau^i, \emptyset, i) + \omega(\epsilon_\theta(\tau^i, \mathbf{y}, i) - \epsilon_\theta(\tau^i, \emptyset, i)), \quad (5)$$

where  $\omega$  is the guidance scale, and  $\emptyset$  represents the absence of guidance. Setting  $\omega = 0$  removes the classifier-free

guidance towards an unconditional diffusion model, while a large value of  $\omega$  strengthens the influence of the conditional information during trajectory generation.

Also, the loss function to minimize can be rewritten as, 
$$\mathcal{L}_{diff}(\theta) = \mathbb{E}_{i, \tau, \epsilon} \left[ \left\| \epsilon - \epsilon_{\theta}(\tau^i, (1 - \beta)y(\tau^i) + \beta\emptyset, i) \right\|^2 \right], \quad (6)$$

where  $\beta$  is a hyperparameter that controls the probability of dropping specific condition  $y(\tau)$ , enhancing sample diversity while maintaining context relevance.

After training the noise prediction model  $\epsilon_{\theta}$  with the above  $\mathcal{L}(\theta)$ , the trajectory is sampled starting from Gaussian noise and progressively denoised with modified  $\hat{\epsilon}$  through Eq. 5, employing re-parameterization technique.

In summary, through classifier-free guidance, we can modulate the trajectory sampling process to require the generated trajectories more aligned with the desired features represented by  $y(\tau)$ . This process iteratively applies the conditioned noise model to refine the target trajectories that contain future states satisfying the constraints.

## 4. Methodology

### 4.1. Overview

Building upon the motivations discussed above, we present SkillDiffuser, an advanced methodological framework for robust multi-task learning across various robots. This dynamic approach leverages the cooperation of skill learning at the higher level and a conditional diffusion model at the lower level. An overall framework is illustrated in Fig. 2. Notably, we leverage language-grounded representations to abstract skills, thereby rendering the execution of tasks through our diffusion policy both interpretable and comprehensible to humans.

### 4.2. High Level Interpretable Skill Abstractions

In our SkillDiffuser framework, the high level interpretable skill abstraction module plays a crucial role in understanding and executing complex tasks. However, given the multi-task environments we consider, each task with only a single language instruction may be broken down into a sequence of sub-tasks or skills, which are not explicitly delineated within the instruction itself. Furthermore, suppose the offline training dataset is denoted as  $\mathcal{D}$ , it consists of trajectories derived from a sub-optimal policy for various tasks. A trajectory  $\tau = (l, \{i_t, a_t\}_{t=1}^T)$  includes language description  $l$  and a sequence of image observations and actions  $(i_t, a_t)$  over time steps  $T$ , with no reward labels attached. But the trajectories do not indicate the boundaries between sub-tasks as well, which thus requires the proposed methods capable of segmenting and interpreting the tasks into sub-goals in an unsupervised manner.

To address this challenge, we build a skill abstraction component upon a horizon-based skill predictor adapted

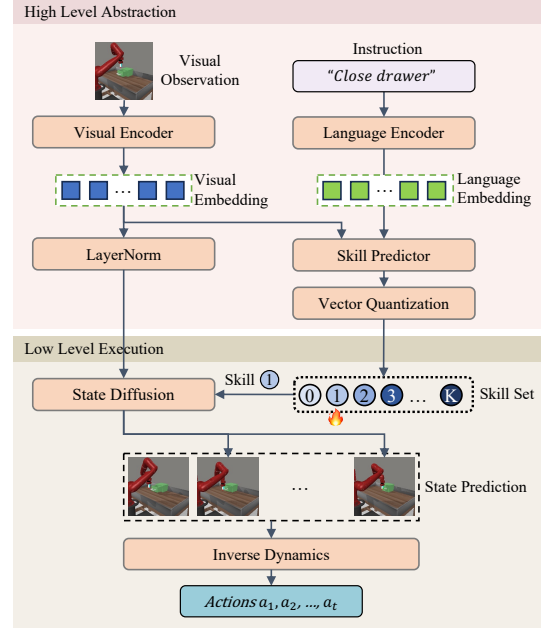


Figure 2. **Overall framework of SkillDiffuser.** It’s a hierarchical planning model that leverages the cooperation of interpretable skill abstractions at the higher level and a skill conditioned diffusion model at the lower level for task execution in a multi-task learning environment. The high-level skill abstraction is achieved through a skill predictor and a vector quantization operation, generating sub-goals (skill set) that the diffusion model employs to determine the appropriate future states. Future states are converted to actions using an inverse dynamics model. This unique fusion enables a consistent underlying planner across different tasks, with the variation only in the inverse dynamics model.

from GPT-2 [35]. It is designed to parse and decompose tasks by fusing visual input and natural language instructions, alongside a Vector Quantization (VQ) sub-module that discretizes the learned skills into a skill set. The specifics of this component are illustrated below.

Firstly, as we utilize only images as robot state information, we transform the images into latent space features with a fixed image encoder (e.g. R3M [30]). For convenience, we denote the image encoder as  $\Phi_{im} : \mathcal{I} \rightarrow \mathcal{R}^I$ , where  $\mathcal{I}$  represents the space of input images and  $\mathcal{R}^I$ , the resultant feature space, serves to condense the visual information into a form that is conducive for high-level semantics. Simultaneously, we use a language encoder to pre-process the natural language instructions, which is formalized as  $\Phi_{lang} : L \rightarrow \mathcal{R}^L$ , with  $L$  being the space of language instructions and  $\mathcal{R}^L$  the language feature space. The outputs of both encoders are then fed into the skill predictor, which operates to integrate these two modalities.

The skill predicting process is as follows: An image  $i_t \in \mathcal{I}$  at time step  $t$  is encoded into a visual embedding  $s_t = \Phi_{im}(i_t)$ . This embedding  $s_t$  is then input to the skill predictor, along with the language instruction  $l \in L$ , through the language encoder’s output  $l_t = \Phi_{lang}(l)$ . The



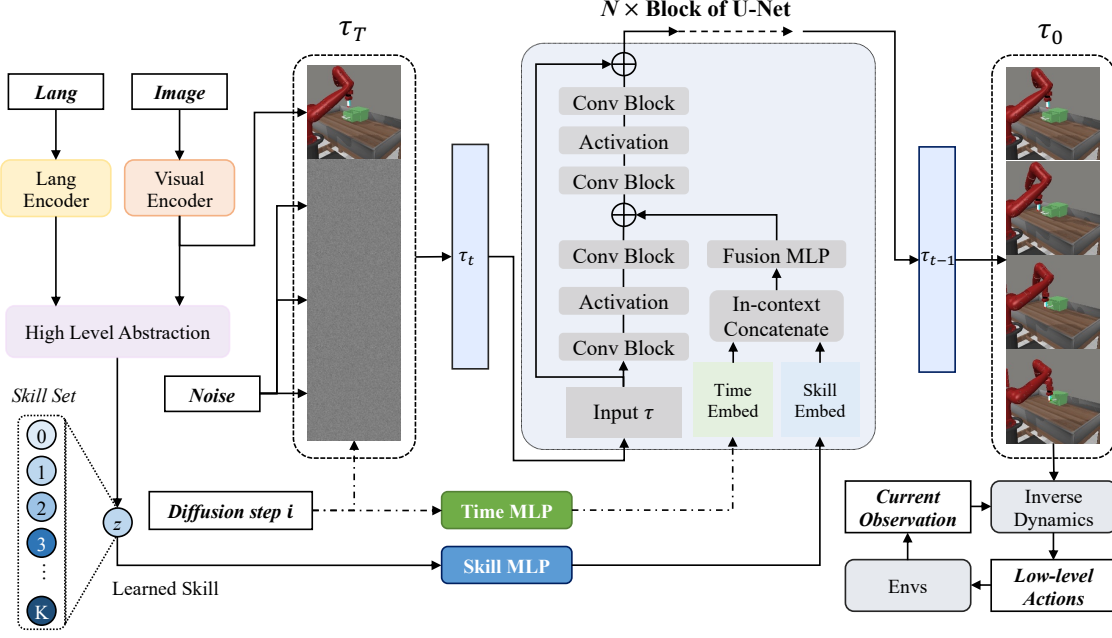


Figure 3. **SkillDiffuser’s low level skill-conditioned diffusion planning model.** Notably, while the schematic here employs images to represent visual features for illustrative purposes, in actual implementation, both the input to and the sampling output of the diffusion model are state embeddings. The current observation is also the feature embedding of current visual observation.

skill predictor, represented as  $f : \mathcal{R}^I \times \mathcal{R}^L \rightarrow \mathcal{C}$ , generates a skill code  $\tilde{z}$  by  $\tilde{z} = f(s_t, l_t)$  that encapsulates task’s requirements interpreted from the visual and language inputs.

After that, Vector Quantization [45] operation  $\mathbf{q}(\cdot)$  is taken to transform these predicted skills into a low-dimensional discrete set  $\mathcal{C}$ . The discrete skill set contains  $K$  skill embeddings  $\{z^1, z^2, \dots, z^K\}$  which represent different potential skills. VQ operation is achieved by mapping the latent  $\tilde{z}$  to its closest entry of the skill set with skill vectors updated to be the moving average of the embeddings  $z$  closest to them, same as [45]. This process is as follows,

$$\tilde{z} = f(\Phi_{im}(\tilde{i}_t), \Phi_{lang}(l_t)), \quad (7)$$

$$z = \mathbf{q}(\tilde{z}) = \arg \min_{z^k \in \mathcal{C}} \|\tilde{z} - z^k\|_2. \quad (8)$$

VQ enforces each learnt skill  $z$  to lie in  $\mathcal{C}$ , which is equal to learning  $K$  prototypes for the language embeddings utilizing  $k$ -means [24] algorithm. This acts as a bottleneck, restricting the flow of language information and aiding in the learning of discrete skill codes. The back-propagation through the non-differentiable quantization is achieved by a straight-through gradient estimator, which simply copies the gradients and enables model to be trained end-to-end.

In our approach, we apply a consistent skill code  $z$  across a defined number of time steps, termed the horizon. This consistent application across multiple horizons adeptly addresses the challenge of varying sub-task durations without altering the horizon itself. Consequently, this strategy not only preserves the flexibility required for diverse task execution but also maintains the model’s architectural stability by avoiding horizon-induced structural changes.

Importantly, the discrete nature of the learned skill codes enhances the interpretability and controllability of the system’s behavior, as each skill code is associated with some human-understandable language phrases. An example is depicted in Fig 6. Through this method, SkillDiffuser can not only learn to perform tasks based on language instructions but also achieve them in a human-interpretable way, allowing for a deeper understanding and control of the decision-making process of specific embodied agents.

### 4.3. Low Level Skill-conditioned Diffusion Planning

As highlighted in Section 1, while existing approaches like Decision Diffuser [2] have introduced conditional diffusion models constrained by skills, their capability is limited to generating trajectories that meet only predefined skills. Consequently, these models fall short of realizing a diffusion model capable of conditioning on an arbitrary learned skill. To overcome these constraints and enable diffusion models to plan over a learned continuous spectrum of skills, we propose an approach that leverages the classifier-free guided diffusion model with the skills embedded.

SkillDiffuser begins by employing a diffusion model to operate over the continuous skill space  $\mathcal{Z}$  learned during high-level skill abstraction. We employ a U-Net to serve as the noise prediction model  $\epsilon_\theta(\cdot)$  and guide it by in-context conditioning. More specifically, we firstly utilize a skill MLP (similar to point-wise feed-forward network [46]), denoted as  $\Lambda$ , to align skill features with denoising model. After that, we fuse the skill embeddings  $\Lambda(z)$  into the state features throughout the residual blocks of U-Net. Details

are depicted in Fig. 3. In this way, we make the diffusion model contextually modulate the influence of skill embeddings at each step. And following Eq. 5, we can synthesize future states attending to these given skills. This is a significant shift from previous static conditioning framework to a more dynamic and adaptive trajectory generation process.

Moreover, following previous works [2, 7, 23], we adopt a state-only diffusion model, which eschews the direct generation of actions in diffusion model. And instead, we utilize another MLP, denoted as  $\Psi(\cdot)$ , to perform inverse dynamics after state generation to infer feasible actions that can achieve transitions between two continuous states. Additionally, we integrate observations from the current frame to provide more detailed information for motion prediction and achieve closed-loop control. Mathematically, it is:

$$\mathbf{a}_t = \Psi([\tilde{\mathbf{s}}_t, \tilde{\mathbf{s}}_{t+1}], \mathbf{i}_t) \quad \text{for } t = 0, \dots, T-1, \quad (9)$$

where  $\tilde{\mathbf{s}}_t$  and  $\tilde{\mathbf{s}}_{t+1}$  are consecutive observation embeddings within  $\tilde{\tau}$  generated by diffusion model,  $\mathbf{i}_t$  is the current observation and  $\mathbf{a}_t$  is the inferred action.

As the resulting action sequence  $\{\mathbf{a}_0, \mathbf{a}_1, \dots, \mathbf{a}_{T-1}\}$  derived from the generated states, it encapsulates the skills to execute the tasks, which allows for remarkable adaptability across multiple tasks. And when faced with a new task, we are only required to change the inverse dynamics model  $\Psi(\cdot)$  specific to the new task’s kinematics, with the architecture and parameters of diffusion model unchanged. Such modularity ensures the generative capabilities of SkillDiffuser are not task-specific but can be leveraged across a diverse range of tasks with varying dynamics. A schematic diagram of the low level module is shown in Fig. 3.

#### 4.4. Training the SkillDiffuser

We engineered a threefold loss function for SkillDiffuser. Firstly, for the inverse dynamics model which is task-specific, we employ a behavior cloning loss [44] to train our inverse MLP emulating expert actions from observed state transitions. This loss, denoted as  $\mathcal{L}_{inv}$ , is formulated as:

$$\mathcal{L}_{inv} = \mathbb{E}_{\tau \sim \mathcal{D}} [\|\mathbf{a} - \Psi(\tilde{\mathbf{s}}, \tilde{\mathbf{s}}', \mathbf{i})\|_2^2], \quad (10)$$

where the notations are similar to Eq. 9.

Correspondingly, the other parts including both skill abstraction and low level execution are task-agnostic. For the high-level skill abstraction module, we apply a vector quantization (VQ) loss to refine the skill predictor. This VQ loss,  $\mathcal{L}_{VQ}$ , ensures the embeddings produced by the skill predictor closely match the skill set vectors, thereby improving the interpretability and consistency of the skill representations. Inspired by VQ-VAE [45], we formulate it as:

$$\mathcal{L}_{VQ} = \mathbb{E}_{\tau} [\|\mathbf{q}(\tilde{\mathbf{z}}) - \tilde{\mathbf{z}}\|_2^2], \quad (11)$$

where  $\tilde{\mathbf{z}}$  follows Eq. 7.

Lastly, for the low-level state-only skill-conditioned diffusion execution, we incorporate a diffusion loss  $\mathcal{L}_{diff}$  as

per Eq. 6, ensuring our model’s state predictions are in line with both the skill guidance and temporal dynamics observed in expert demonstrations. Here, we take  $\mathbf{y}(\tau) = \Lambda(\mathbf{z})$  with  $\mathbf{z}$  derived from Eq. 7 and Eq. 8.

To be noted, we train our SkillDiffuser with two optimizers, one for inverse dynamics model with  $\mathcal{L}_{inv}$  and the other for overall parameters of high-level skill abstraction and low-level planning with  $\mathcal{L}_{VQ} + \lambda \mathcal{L}_{diff}$ , where  $\lambda$  is a loss weight. This carefully constructed loss architecture enables SkillDiffuser to excel in a multi-task environment, generalizing across tasks by simply substituting the inverse dynamics model  $\Psi(\cdot)$  specific to each new task’s requirements.

Additionally, we utilize a pre-trained distilBERT [39] as our language encoder, adopting the configuration consistent with LOReL [29] and LISA [13], while freezing its parameters to guarantee stability in language understanding. And we employ diverse settings to serve as the visual encoder to ensure fair comparison in different datasets. We elaborate the details in corresponding parts of Sec. 5. More training details are shown in Appendix G and we also provide some pseudo-code of our algorithm in Appendix B.

## 5. Experiments

We first present a comprehensive evaluation on the LOReL Sawyer dataset, and then perform the ablation study compared on Meta-World benchmark to illustrate the efficiency of our method. Finally, we visualize the learned skills of our method both on LOReL and Meta-World MT10.

### 5.1. Datasets

**LOReL Sawyer Dataset** [29] which is abbreviated from Language-conditioned Offline Reward learning, is composed of *pseudo-expert* trajectories or *play data* gathered from an arbitrary reinforcement learning policy, annotated with post-hoc crowd-sourced language directives. The LOReL Sawyer dataset encompasses 50k trajectories, each with 20 steps, in a simulated Sawyer robot environment. We assess our approach using the same six tasks as the original paper [29] with paraphrased instructions under five different situations (*i.e.* seen, unseen verb, unseen none, unseen verb + noun and human provided). This comes to a total of 77 instructions for all 6 tasks combined. More details about the dataset can be found in Appendix D.1.

**Meta-World Dataset** [51]. It is a comprehensive benchmark designed for evaluating multi-task and meta reinforcement learning algorithms. It introduces a comprehensive suite of 50 distinct robotic manipulation tasks, all located within a unified tabletop environment controlled by a simulated Sawyer arm. The Multi-Task 10 (MT10) within Meta-World is a subset comprising ten carefully selected tasks, balanced in terms of diversity and complexity. Details of the ten tasks can be found in Appendix D.3.

Task Instruction	Random	LCBC [42]	LCRL [20]	Lang DT [13]	LISA [13]	SkillDiffuser
close drawer	52%	50%	58%	10%	<b>100%</b>	<b>95 ± 3.2%</b>
open drawer	14%	0%	8%	<b>60%</b>	20%	<b>55 ± 13.3%</b>
turn faucet left	24%	12%	13%	0%	0%	<b>55 ± 9.3%</b>
turn faucet right	15%	<b>31%</b>	0%	0%	<b>30%</b>	<b>25 ± 4.4%</b>
move black mug right	12%	<b>73%</b>	0%	20%	60%	18 ± 3.9%
move white mug down	5%	6%	0%	0%	<b>30%</b>	10 ± 1.7%
<b>Average over tasks</b>	20%	29%	13%	15%	40%	<b>43 ± 1.1%</b>

Table 1. **Task-wise success rates on LOReL Sawyer Dataset.** We show our success rates compared to random policy, language-conditioned imitation learning (LCBC), language-conditioned Q-learning (LCRL), a flat non-hierarchical Decision Transformer (Lang-DT), and LISA. The results on each dataset are calculated over 3 seeds. SkillDiffuser outperforms all other methods in terms of average performance over all tasks. Best methods and those within 6% of the best are highlighted in **bold**.

Rephrasal Type	Lang DT	LISA [13]	SkillDiffuser
seen	15	40	<b>43.65</b> ± 4.7
unseen noun	13.33	33.33	<b>36.01</b> ± 6.3
unseen verb	28.33	30	<b>36.70</b> ± 9.5
unseen noun+verb	6.7	20	<b>42.02</b> ± 3.8
human provided	26.98	27.35	<b>40.16</b> ± 2.1
<b>Average</b>	18.07	30.14	<b>39.71</b>

Table 2. **Rephrasal-wise success rates (in %) on LOReL Sawyer.** Results of Lang DT, LISA and our SkillDiffuser are shown here. The standard error is calculated over 3 random seeds.

## 5.2. Evaluation Results on LOReL Sawyer Dataset

**Baselines.** We employ random policy, language conditioned imitation learning (LCBC) [42], language conditioned Q-learning (LCRL) [20], Lang-DT (also known as Flat Baseline in [13]) and state-of-the-art skill-learning method LISA [13] as our baselines. We follow the same settings as LOReL [29] for the first three algorithms, and follow LISA [13] for the last two. The random policy serves as a baseline. And LCBC mimics offline dataset behaviors based on instructions, aligning with previous works focusing on imitation learning to achieve language conditioned behavior. In contrast, LCRL employs reinforcement learning, labeling each episode’s final state with language-instructed rewards. Lang-DT plays as a non-hierarchical benchmark with a Causal Transformer [25], while LISA incorporates a dual-transformer structure, with one for skill prediction and the other for action planning. We do not compare with LOReL planner [29] as it uses MPC on a learned reward function relying on human annotations, while LISA and ours learn with trajectory data only.

**Results.** To ensure fair comparison, our method, SkillDiffuser, is designed to maintain a parameter count similar to baseline models. It employs the same visual and language encoder architecture as used in LISA [13] with meticulously matching of embedding dimensions and the number of heads across the layers of SkillDiffuser.

The results are present in Tables 1 and 2, showing task-wise and rephrasal-wise success rates for LOReL, averaged

Method	Lang DT	LOReL [29]	LISA [13]	SkillDiffuser
<b>Success Rate</b>	13.33 ± 1.3	18.18 ± 1.8	20.89 ± 0.6	<b>25.21</b> ± 2.7

Table 3. **Performance on LOReL multi-step composition tasks.**

over 10 runs with a 20-step time horizon. SkillDiffuser, our approach, achieves the highest average performance in six different tasks, indicating its superior cross-task adaptability particularly when compared to similar approaches which yet are based on Decision Transformer [4], such as LISA [13]. Moreover, SkillDiffuser consistently excels in all rephrased types for LOReL test tasks, outperforming LISA by 9.6% on average. This demonstrates the model’s robustness in handling varied skill representations, marking a notable advancement in skill-conditioned diffusion model.

## 5.3. Performance on LOReL Compositional Tasks

**Settings.** We conduct experiments following the same settings of unseen composition tasks of LISA [13] with 12 composition instructions in Tab. 3. Detailed instructions are listed in Appendix D.2 with such an example that “open drawer and move black mug right”. We extend the max number of episode steps from customary 20 to 40, as LISA.

**Results.** We observe SkillDiffuser achieves 2x the performance of non-hierarchical baseline (i.e. w/o skill abstraction) and also improves about 25% over LISA, highlighting its effectiveness. MPC-based LOReL planners are unable to perform as well in open scenarios like composition tasks.

## 5.4. Ablation Study on Meta-World Dataset

**Settings.** We conduct experiments on Meta-World MT10 benchmark with finely annotated instructions. We also use visual observations only. Details are in Appendix D.3.

**Baselines.** We evaluate our method against three baselines, all modified from existing models. The first, Flat R3M, is adapted from R3M [30] paper’s planner. As the original one utilizes the first four terms of robot proprioception, we eliminate them and make the planner focus exclusively on visual observations. The second baseline is a variant of our SkillDiffuser, lacking the high-level skill abstraction module but retaining the low-level conditional diffusion-based planner,

Method	Lang	Skill Set	Diffuser	Performance
Flat R3M [30]	✗	✗	✗	13.3%
LISA [13]	✓	✓	✗	13.8%
Lang Diffuser	✓	✗	✓	16.7%
SkillDiffuser	✓	✓	✓	<b>23.3%</b>

Table 4. **Ablation study of language skill conditioning on Meta-World.** Results of Flat R3M, Language-condition diffuser, LISA and our SkillDiffuser are shown here. All are averaged over 3 runs.

to assess the impact of skill abstraction. This version integrates a two-layer MLP to predict options from visual and language inputs, functioning as a language-conditioned diffusion planner. The last baseline is our re-implemented LISA [13] for Meta-World Dataset to validate the efficiency of diffusion model. To ensure fairness, we use R3M as the visual encoder for all of these methods on Meta-World.

**Ablation on Language Skill Conditioning.** Table 4 indicates our Meta-World MT10 task is quite different from and much more challenging than previous tasks that use states as observations [49] or take into robot proprioception [30]. We only utilize single-frame visual input and instructions. The Flat R3M method, lacking language conditioning and skill sets can only succeed on tasks like *drawer-close* and *reach* through behavior cloning. Lang-conditioned Diffuser and modified LISA both outperform Flat R3M, suggesting the value of each corresponding module. Our SkillDiffuser, discretizing skills into a skill set, achieves a 6% higher performance than language-conditioned diffuser and a 9.5% higher than LISA, demonstrating the effectiveness of this combinational architecture.

### 5.5. Ablation Study on Reusability of Learned Skills

To evaluate the reusability of our learned skills, we calculate the average number of different skills used for a single instruction and the total number of instructions using each skill of LOReL Sawyer Dataset in Table 5. (With max episode step being 20, we experiment with skill horizon 10.) We observe each single instruction uses 1.55 sub-skills on average and each skill is called multiple times than the number of instructions (75 with 5 eval episodes), verifying the transferability of learnt skills. As suggested in Tab.10 of [13], except a very small horizon will hurt the performance, learning sub-skills to get refined semantics helps perform different actions at different stages. Besides, we also visualize resulting images from applying discrete skills in Appendix E.1 to further validate skills’ interpretability.

### 5.6. Visualization Results of Learned Skill Set

We show the visualization of skill set on LOReL compositional tasks here in Fig. 4 and results on original LOReL dataset in Fig. 5 and Meta-World dataset in Fig. 6 in Appendix C. The visual analysis of our SkillDiffuser’s learned skills on LOReL compositional tasks reveals that out of the 20-size skill set, our method learned 11 skills (e.g. *pull han-*

# of learnt skills	# of inst	# of success	Use 1 skill	Use 2 skills	Average
17	375	144	64	80	1.55
Freq of 17 skills	30, 8, 14, 10, 5, 19, 4, 7, 20, 2, 9, 20, 25, 6, 10, 3, 18				

Table 5. **Average number of different skills used for a single instruction and total number of instructions used for each skill.**

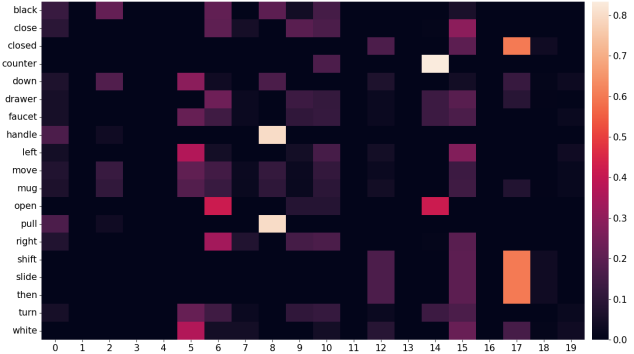


Figure 4. **Visualization of skill heat map on LOReL Sawyer compositional tasks.** We display the word frequency associated with a skill set of size 20 in LOReL, normalized by column. The data’s sparsity and distinct highlights indicate certain language tokens are uniquely linked to specific skills. There are eleven skills learned by our method. (zoom in for best view)

*dle* [skill 0], *open counter* [skill 14], etc.) notably distinguished by their unique word highlights. These bright spots across eleven columns (changing from only one column at initial which corresponds to default BC) in the heatmap underscore the model’s ability to identify and isolate distinct skills from visual inputs, without an explicitly defined skill library. This indicates not only a significant interpretative advancement over previous diffusion-based planning but a successful abstraction of high-level skill representations.

## 6. Conclusion

This paper presents SkillDiffuser, an integrated framework that enables robots to perform tasks from natural language instructions by enabling interpretable skill learning and conditional diffusion planning. It employs vector quantization to learn discrete and comprehensible skill representations directly from visual and linguistic demonstrations. Subsequently, these skills condition a diffusion model to generate state trajectories adhering to the learned skills. Through integrating hierarchical skill decomposition with conditional trajectory generation, SkillDiffuser can comprehend and execute abstract instructions for various manipulation tasks. Extensive experiments on manipulation benchmarks demonstrate state-of-the-art performance, highlighting its effectiveness for multi-step composition tasks and ability to automatically learn interpretable skills.

## Acknowledgements

This paper is partially supported by the National Key R&D Program of China No.2022ZD0161000 and the General Research Fund of Hong Kong No.17200622.



## References

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. 2022. [1](#)
- [2] Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B Tenenbaum, Tommi S Jaakkola, and Pulkit Agrawal. Is conditional generative modeling all you need for decision making? In *The Eleventh International Conference on Learning Representations*, 2022. [1](#), [3](#), [5](#), [6](#)
- [3] Joao Carvalho, An T Le, Mark Baierl, Dorothea Koert, and Jan Peters. Motion planning diffusion: Learning and planning of robot motions with diffusion models. *arXiv preprint arXiv:2308.01557*, 2023. [3](#)
- [4] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021. [7](#)
- [5] Cheng Chi, Siyuan Feng, Yilun Du, Zhenjia Xu, Eric Cousineau, Benjamin Burchfiel, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. In *Proceedings of Robotics: Science and Systems (RSS)*, 2023. [1](#), [3](#)
- [6] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. [1](#), [3](#)
- [7] Yilun Du, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Joshua B Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems*, 2023. [1](#), [3](#), [6](#)
- [8] Yan Duan, Marcin Andrychowicz, Bradly Stadie, OpenAI Jonathan Ho, Jonas Schneider, Ilya Sutskever, Pieter Abbeel, and Wojciech Zaremba. One-shot imitation learning. *Advances in neural information processing systems*, 30, 2017. [2](#)
- [9] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. *arXiv preprint arXiv:1802.06070*, 2018. [2](#)
- [10] William Feller. On the theory of stochastic processes, with particular reference to applications. In *Selected Papers I*, pages 769–798. Springer, 2015. [3](#)
- [11] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation learning via meta-learning. In *Conference on robot learning*, pages 357–368. PMLR, 2017. [2](#)
- [12] Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, and Stefano Ermon. Iq-learn: Inverse soft-q learning for imitation. *Advances in Neural Information Processing Systems*, 34:4028–4039, 2021. [2](#)
- [13] Divyansh Garg, Skanda Vaidyanath, Kuno Kim, Jiaming Song, and Stefano Ermon. Lisa: Learning interpretable skill abstractions from language. *Advances in Neural Information Processing Systems*, 35:21711–21724, 2022. [2](#), [3](#), [6](#), [7](#), [8](#), [5](#)
- [14] Haoran He, Chenjia Bai, Kang Xu, Zhuoran Yang, Weinan Zhang, Dong Wang, Bin Zhao, and Xuelong Li. Diffusion model is an effective planner and data synthesizer for multi-task reinforcement learning. *Advances in neural information processing systems*, 2023. [1](#), [2](#)
- [15] Bernd Heidergott, editor. *Taylor Series Expansions*, pages 179–263. Springer US, Boston, MA, 2007. [1](#)
- [16] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in neural information processing systems*, 29, 2016. [2](#)
- [17] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. [3](#)
- [18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [1](#), [3](#), [2](#)
- [19] Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for flexible behavior synthesis. In *International Conference on Machine Learning*, pages 9902–9915. PMLR, 2022. [1](#), [3](#)
- [20] Yiding Jiang, Shixiang Shane Gu, Kevin P Murphy, and Chelsea Finn. Language as an abstraction for hierarchical deep reinforcement learning. *Advances in Neural Information Processing Systems*, 32, 2019. [7](#)
- [21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. [7](#)
- [22] Alexander C Li, Carlos Florensa, Ignasi Clavera, and Pieter Abbeel. Sub-policy adaptation for hierarchical reinforcement learning. *arXiv preprint arXiv:1906.05862*, 2019. [3](#)
- [23] Zhixuan Liang, Yao Mu, Mingyu Ding, Fei Ni, Masayoshi Tomizuka, and Ping Luo. AdaptDiffuser: Diffusion models as adaptive self-evolving planners. In *International Conference on Machine Learning*, 2023. [1](#), [3](#), [6](#)
- [24] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, pages 281–297. Oakland, CA, USA, 1967. [5](#)
- [25] Valentyn Melnychuk, Dennis Frauen, and Stefan Feuerriegel. Causal transformer for estimating counterfactual outcomes. In *International Conference on Machine Learning*, pages 15293–15329. PMLR, 2022. [7](#)
- [26] Diganta Misra. Mish: A self regularized non-monotonic activation function. In *31st British Machine Vision Conference 2020, BMVC 2020, Virtual Event, UK, September 7-10, 2020*, 2020. [6](#)
- [27] Ofir Nachum, Shixiang Shane Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. *Advances in neural information processing systems*, 31, 2018. [3](#)
- [28] Ashvin V Nair, Vitchyr Pong, Murtaza Dalal, Shikhar Bahl, Steven Lin, and Sergey Levine. Visual reinforcement learning with imagined goals. *Advances in neural information processing systems*, 31, 2018. [2](#)
- [29] Suraj Nair, Eric Mitchell, Kevin Chen, Silvio Savarese, Chelsea Finn, et al. Learning language-conditioned robot behavior from offline data and crowd-sourced annotation. In

- Conference on Robot Learning*, pages 1303–1315. PMLR, 2022. 2, 6, 7, 4
- [30] Suraj Nair, Aravind Rajeswaran, Vikash Kumar, Chelsea Finn, and Abhinav Gupta. R3m: A universal visual representation for robot manipulation. In *Conference on Robot Learning*, pages 892–909. PMLR, 2023. 2, 4, 7, 8, 6
- [31] Fei Ni, Jianye Hao, Yao Mu, Yifu Yuan, Yan Zheng, Bin Wang, and Zhixuan Liang. Metadiffuser: Diffusion model as conditional planner for offline meta-rl. In *International Conference on Machine Learning*, 2023. 3
- [32] Scott Niekum, Sarah Osentoski, George Konidaris, Sachin Chitta, Bhaskara Marthi, and Andrew G Barto. Learning grounded finite-state representations from unstructured demonstrations. *The International Journal of Robotics Research*, 34(2):131–157, 2015. 2
- [33] Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97, 1991. 2
- [34] Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 1994. 3
- [35] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 4
- [36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015. 6
- [37] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 661–668. JMLR Workshop and Conference Proceedings, 2010. 2
- [38] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635. JMLR Workshop and Conference Proceedings, 2011. 2
- [39] V Sanh. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. In *Proceedings of Thirty-third Conference on Neural Information Processing Systems*, 2019. 6
- [40] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. *arXiv preprint arXiv:1907.01657*, 2019. 2
- [41] Avi Singh, Eric Jang, Alexander Irpan, Daniel Kappler, Murtaza Dalal, Sergey Levine, Mohi Khansari, and Chelsea Finn. Scalable multi-task imitation learning with autonomous improvement. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2167–2173. IEEE, 2020. 2
- [42] Simon Stepputtis, Joseph Campbell, Mariano Phielipp, Stefan Lee, Chitta Baral, and Heni Ben Amor. Language-conditioned imitation learning for robot manipulation tasks. *Advances in Neural Information Processing Systems*, 33: 13139–13150, 2020. 7
- [43] Alan Stuart and J Keith Ord. Kendall’s advanced theory of statistics. vol. 1: Distribution theory. *Kendall’s advanced theory of statistics. Vol. 1: Distribution theory*, 1994. 2
- [44] Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 4950–4957, 2018. 6
- [45] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017. 2, 5, 6
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 5
- [47] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on computer vision (ECCV)*, pages 3–19, 2018. 6
- [48] Danfei Xu, Suraj Nair, Yuke Zhu, Julian Gao, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Neural task programming: Learning to generalize across hierarchical tasks. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3795–3802. IEEE, 2018. 2
- [49] Mengdi Xu, Yikang Shen, Shun Zhang, Yuchen Lu, Ding Zhao, Joshua Tenenbaum, and Chuang Gan. Prompting decision transformer for few-shot policy generalization. In *international conference on machine learning*, pages 24631–24645. PMLR, 2022. 2, 8
- [50] Long Yang, Zhixiong Huang, Fenghao Lei, Yucun Zhong, Yiming Yang, Cong Fang, Shiting Wen, Binbin Zhou, and Zhouchen Lin. Policy representation via diffusion probability model for reinforcement learning. *arXiv preprint arXiv:2305.13122*, 2023. 1
- [51] Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning. In *Conference on robot learning*, pages 1094–1100. PMLR, 2020. 2, 6
- [52] Jesse Zhang, Haonan Yu, and Wei Xu. Hierarchical reinforcement learning by discovering intrinsic options. *arXiv preprint arXiv:2101.06521*, 2021. 3

# SkillDiffuser: Interpretable Hierarchical Planning via Skill Abstractions in Diffusion-Based Task Execution

## Supplementary Material

### A. Theoretical Foundation of Classifier-free Diffusion Model for Planning

#### A.1. Review of Classifier-guided Diffusion Model

Firstly, for a given trajectory  $\tau$ , the standard reverse process of an unconditional diffusion probabilistic model is defined by  $p_\theta(\tau^i | \tau^{i+1})$ . This framework is then extended to incorporate conditioning on a specific label  $\mathbf{y}$  (e.g., the reward), which is considered in the context of current-step denoised trajectory  $\tau^i$ , which is represented as  $p_\phi(\mathbf{y} | \tau^i)$ . Consequently, the reverse diffusion process can be reformulated as  $p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y})$ . This approach introduces additional parameters  $\phi$  alongside the original diffusion model parameters  $\theta$ . The parameters  $\phi$  can be viewed as a classifier, that encapsulates the probability of whether a noisy trajectory  $\tau^i$  satisfies the specific label  $\mathbf{y}$ , with a notation of  $p_\phi(\mathbf{y} | \tau^i)$ .

Under the constraints illustrated in [6, 23], we can derive the following theorem with lemma

$$p_{\theta,\phi}(\mathbf{y} | \tau^i, \tau^{i+1}) = p_\phi(\mathbf{y} | \tau^i). \quad (12)$$

**Theorem A.1.** *The conditional sampling probability of reverse diffusion process  $p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y})$  is proportional to unconditional transition probability  $p_\theta(\tau^i | \tau^{i+1})$  multiplied by the classified probability  $p_\phi(\mathbf{y} | \tau^i)$ .*

$$p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y}) = Z p_\theta(\tau^i | \tau^{i+1}) p_\phi(\mathbf{y} | \tau^i) \quad (13)$$

*Proof.*

$$\begin{aligned} p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y}) &= \frac{p_{\theta,\phi}(\tau^i, \tau^{i+1}, \mathbf{y})}{p_{\theta,\phi}(\tau^{i+1}, \mathbf{y})} \\ &= \frac{p_{\theta,\phi}(\mathbf{y} | \tau^i, \tau^{i+1}) p_\theta(\tau^i, \tau^{i+1})}{p_\phi(\mathbf{y} | \tau^{i+1}) p_\theta(\tau^{i+1})} \\ &= \frac{p_{\theta,\phi}(\mathbf{y} | \tau^i, \tau^{i+1}) p_\theta(\tau^i | \tau^{i+1}) p_\theta(\tau^{i+1})}{p_\phi(\mathbf{y} | \tau^{i+1}) p_\theta(\tau^{i+1})} \\ &= \frac{p_{\theta,\phi}(\mathbf{y} | \tau^i, \tau^{i+1}) p_\theta(\tau^i | \tau^{i+1})}{p_\phi(\mathbf{y} | \tau^{i+1})} \\ &= \frac{p_\phi(\mathbf{y} | \tau^i) p_\theta(\tau^i | \tau^{i+1})}{p_\phi(\mathbf{y} | \tau^{i+1})}. \end{aligned} \quad (14)$$

The term  $p_\phi(\mathbf{y} | \tau^{i+1})$  is not directly correlated to  $\tau^i$  at the diffusion timestep  $i$ , thus can be viewed as a constant with notation  $Z$ .  $\square$

On this basis, using Taylor series expansion [15], we can sample trajectories by the modified Gaussian resampling.

**Theorem A.2.** *With a sufficiently large number of reverse diffusion steps, the sampling from reverse diffusion process  $p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y})$  can be approximated by a modified Gaussian resampling. That is*

$$p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y}) \approx \mathcal{N}(\tau^i; \mu_\theta + \Sigma \nabla_\tau \log p_\phi(\mathbf{y} | \tau^i), \Sigma), \quad (15)$$

where  $\mu_\theta = \mu_\theta(\tau^i)$  and  $\Sigma$  are the mean and variance of unconditional reverse diffusion process  $p_\theta(\tau^i | \tau^{i+1})$ .

*Proof.* With the above definition, we can rewrite the transfer probability of the unconditional denoising process as

$$p_\theta(\tau^i | \tau^{i+1}) = \mathcal{N}(\tau^i; \mu_\theta, \Sigma) \quad (16)$$

$$\log p_\theta(\tau^i | \tau^{i+1}) = -\frac{1}{2}(\tau^i - \mu_\theta)^T \Sigma^{-1}(\tau^i - \mu_\theta) + C \quad (17)$$

With a sufficiently large number of reverse diffusion steps, we apply Taylor expansion around  $\tau^i = \mu_\theta$  as

$$\begin{aligned} \log p_\phi(\mathbf{y} | \tau^i) &= \log p_\phi(\mathbf{y} | \tau^i) |_{\tau^i=\mu_\theta} \\ &\quad + (\tau^i - \mu_\theta)^T \nabla_{\tau^i} \log p_\phi(\mathbf{y} | \tau^i) |_{\tau^i=\mu_\theta}. \end{aligned}$$

Therefore, using Eq. 14, we derive

$$\log p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y}) = \log p_\theta(\tau^i | \tau^{i+1}) + \log p_\phi(\mathbf{y} | \tau^i) + C_1$$

$$\begin{aligned} RHS &= -\frac{1}{2}(\tau^i - \mu_\theta)^T \Sigma^{-1}(\tau^i - \mu_\theta) \\ &\quad + (\tau^i - \mu_\theta)^T \nabla \log p_\phi(\mathbf{y} | \tau^i) + C_2 \end{aligned} \quad (18)$$

$$\begin{aligned} RHS &= -\frac{1}{2}(\tau^i - \mu_\theta - \Sigma \nabla \log p_\phi(\mathbf{y} | \tau^i))^T \times \Sigma^{-1} \\ &\quad \times (\tau^i - \mu_\theta - \Sigma \nabla \log p_\phi(\mathbf{y} | \tau^i)) + C_3, \end{aligned}$$

which means,

$$p_{\theta,\phi}(\tau^i | \tau^{i+1}, \mathbf{y}) \approx \mathcal{N}(\tau^i; \mu_\theta + \Sigma \nabla_\tau \log p_\phi(\mathbf{y} | \tau^i), \Sigma) \quad \square$$

#### A.2. Classifier-free Diffusion Model

While classifier guidance successfully achieves conditional guidance during trajectory generation, it is nonetheless reliant on gradients from a separate trained classifier which is hard to obtain in many cases. Classifier-free guidance [18] seeks to eliminate the classifier, which achieves the same effect as classifier guidance, but without such gradients.

First of all, we define the score function of the unconditional diffusion model as

$$\epsilon_\theta(\tau^i) = -\Sigma \nabla_\tau \log p_\phi(\tau^i). \quad (19)$$

Then, through Eq. 13 and 15, the score function of the classifier-guided diffusion model can be expressed as

$$\epsilon_\theta(\tau^i, \mathbf{y}) = \epsilon_\theta(\tau^i) - \alpha \Sigma \nabla_\tau \log p_\phi(\mathbf{y} | \tau^i), \quad (20)$$

where  $\alpha$  is a scale hyper-parameter.

**Theorem A.3.** *Classifier-free guided diffusion model performs sampling with the linear combination of the conditional and unconditional score estimates as,*

$$\hat{\epsilon}_\theta = \hat{\epsilon}_\theta(\tau^i, \mathbf{y}) = (1 - \omega)\epsilon_\theta(\tau^i) + \omega\epsilon_\theta(\tau^i, \mathbf{y}), \quad (21)$$

*implicitly embedding guidance into the score function, with  $\omega$  the scale hyper-parameter.*

*Proof.* Considering there is an implicit classifier denoted as  $\tilde{p}_\phi(\mathbf{y} | \tau^i)$ , with Bayes Rule [43], we can expand it as

$$\tilde{p}_\phi(\mathbf{y} | \tau^i) \propto \tilde{p}_{\theta, \phi}(\tau^i, \mathbf{y}) / p_\theta(\tau^i).$$

Then gradient of this implicit classifier would be

$$\nabla_\tau \log \tilde{p}_\phi(\mathbf{y} | \tau^i) = \nabla_\tau \log \tilde{p}_{\theta, \phi}(\tau^i, \mathbf{y}) - \nabla_\tau \log p_\theta(\tau^i). \quad (22)$$

Substitute Eq. 19 in RHS, we get

$$\begin{aligned} \alpha \Sigma \nabla_\tau \log \tilde{p}_\phi(\mathbf{y} | \tau^i) &= \alpha \Sigma \nabla_\tau \log \tilde{p}_{\theta, \phi}(\tau^i, \mathbf{y}) \\ &\quad - \alpha \Sigma \nabla_\tau \log p_\theta(\tau^i) \\ &= -\alpha \hat{\epsilon}_\theta(\tau^i, \mathbf{y}) + \alpha \epsilon_\theta(\tau^i) \end{aligned}$$

And then substitute Eq. 20 in LHS,

$$\begin{aligned} \epsilon_\theta(\tau^i) - \epsilon_\theta(\tau^i, \mathbf{y}) &= -\alpha \hat{\epsilon}_\theta(\tau^i, \mathbf{y}) + \alpha \epsilon_\theta(\tau^i) \\ \alpha \hat{\epsilon}_\theta(\tau^i, \mathbf{y}) &= \epsilon_\theta(\tau^i, \mathbf{y}) + (\alpha - 1)\epsilon_\theta(\tau^i) \\ \hat{\epsilon}_\theta(\tau^i, \mathbf{y}) &= (1/\alpha)\epsilon_\theta(\tau^i, \mathbf{y}) + (1 - 1/\alpha)\epsilon_\theta(\tau^i) \end{aligned} \quad (23)$$

Let  $\omega = 1/\alpha$ , we obtain,

$$\hat{\epsilon}_\theta(\tau^i, \mathbf{y}) = (1 - \omega)\epsilon_\theta(\tau^i) + \omega\epsilon_\theta(\tau^i, \mathbf{y}),$$

which is equal to Eq. 21.  $\square$

Therefore, in classifier-free diffusion guidance, we only need to train a single neural network to parameterize both conditional score estimator  $\epsilon_\theta(\tau^i, \mathbf{y})$  and unconditional score estimator  $\epsilon_\theta(\tau^i)$ , where for the unconditional model we can set an empty set  $\emptyset$  for the condition identifier  $\mathbf{y}$  when predicting the score, i.e.  $\epsilon_\theta(\tau^i) = \epsilon_\theta(\tau^i, \mathbf{y} = \emptyset)$ . Following the settings of [18], we jointly train the unconditional and conditional models simply by randomly setting  $\mathbf{y}$  to the unconditional class identifier  $\emptyset$  with probability  $\beta$ , which balances off the diversity and the relevance of the conditional label of generated samples.

## B. Pseudo-code of Training SkillDiffuser

As illustrated in Sec. 4.4, we provide the pseudocode for our SkillDiffuser’s training process in Algorithm 1, detailing its sequential stages and core mechanics. Additionally, Algorithm 2 describes the inference process, illustrating its steps of skill abstraction and trajectory generation.

---

### Algorithm 1 Training process of SkillDiffuser

---

**Input:** Dataset  $\mathcal{D}$  of partially observed trajectories with paired language  $\{\tau_\xi = (l, \{\mathbf{i}_t, \mathbf{a}_t\}_{t=0}^{T-1})\}_{\xi=1}^N$ , size of the skill set  $K$  and horizon  $H$ , pre-trained language and visual encoder  $\Phi_{lang}, \Phi_{im}$

- 1: Initialize skill predictor  $f$ , conditional diffusion model  $\mathcal{M}$ , skill embedding model  $\Lambda$  and inverse dynamics model  $\Psi$
- 2: Vector Quantization op  $\mathbf{q}(\cdot)$
- 3: **while** not converged **do**
- 4:   Sample  $\tau = (l, \{\mathbf{i}_t, \mathbf{a}_t\}_{t=0}^{T-1})$
- 5:   Initialize partially observed states  $S = \{\Phi(\mathbf{i}_0)\}$
- 6:   **for**  $k = 0 \dots \lfloor \frac{T}{H} \rfloor$  **do**  $\triangleright$  Sample a skill every  $H$  steps
- 7:      $z \leftarrow \mathbf{q}(f(\Phi_{lang}(l), S))$
- 8:      $\mathcal{L}_{diff} \leftarrow \mathcal{M}_{diff}(S, \Lambda(z))$   $\triangleright$  Diffusing process
- 9:     **for** step  $t = 1 \dots H$  **do**
- 10:        $S \leftarrow S \cup \{\Phi(\mathbf{i}_{kH+t+1})\}$
- 11:        $\tilde{\mathbf{a}}_{kH+t} \leftarrow \Psi([s_{kH+t}, s_{kH+t+1}], \mathbf{i}_{kH+t})$   $\triangleright$  Predict action using inverse dynamics model
- 12:        $\mathcal{L}_{inv} = \mathbb{E}[\|\mathbf{a}_{kH+t} - \tilde{\mathbf{a}}_{kH+t}\|_2^2]$
- 13:       Train  $\Psi$  with objective  $\mathcal{L}_{inv}$
- 14:     **end for**
- 15:     Train  $f, \Lambda$  and  $\mathcal{M}$  with objective  $\mathcal{L}_{VQ} + \lambda \mathcal{L}_{diff}$
- 16:   **end for**
- 17: **end while**

---



---

### Algorithm 2 Inference process of SkillDiffuser

---

**Input:** Initial partial observation  $\mathbf{i}_0$  and the language instruction  $l$ , pre-trained language and visual encoder  $\Phi_{lang}, \Phi_{im}$

**Input:** Trained skill predictor  $f$ , conditional diffusion model  $\mathcal{M}$ , skill embedding model  $\Lambda$  and inverse dynamics model  $\Psi$

- 1: Initialize partially observed states  $S = \{\Phi(\mathbf{i}_0)\}$
- 2: **for**  $k = 0 \dots \lfloor \frac{T}{H} \rfloor$  **do**  $\triangleright$  Sample a skill every  $H$  steps
- 3:    $z \leftarrow \mathbf{q}(f(\Phi_{lang}(l), S))$
- 4:    $S' \leftarrow \mathcal{M}_{denoise}(S, \Lambda(z))$   $\triangleright$  Denoising process
- 5:   **for** step  $t = 1 \dots H$  **do**
- 6:      $\mathbf{a}_{kH+t} \leftarrow \Psi([s_{kH+t}, s'_{kH+t+1}], \mathbf{i}_{kH+t})$
- 7:      $\tilde{\mathbf{s}}_{kH+t+1} \leftarrow \text{Env.step}(\mathbf{a}_{kH+t})$   $\triangleright$  Take action
- 8:      $S \leftarrow S \cup \{\tilde{\mathbf{s}}_{kH+t+1}\}$
- 9:   **end for**
- 10: **end for**

---



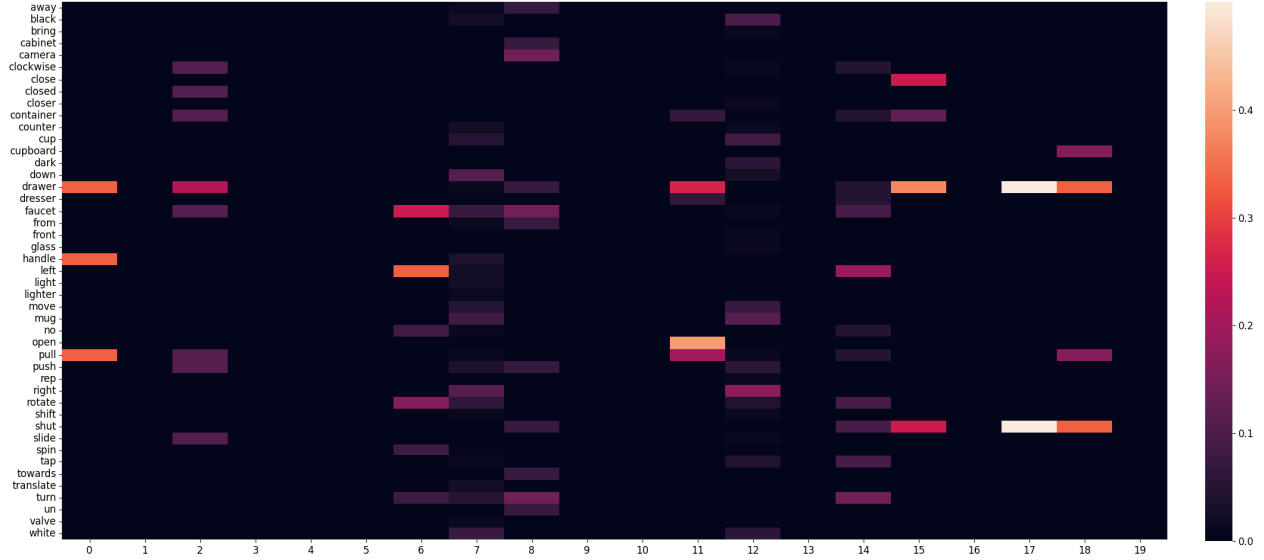


Figure 5. **Visualization of skill heat map on LOReL.** We display the word frequency associated with a skill set of size 20 in LOReL, normalized by column. The data’s sparsity and distinct highlights indicate certain language tokens are uniquely linked to specific skills. There are eleven skills learned by our method.

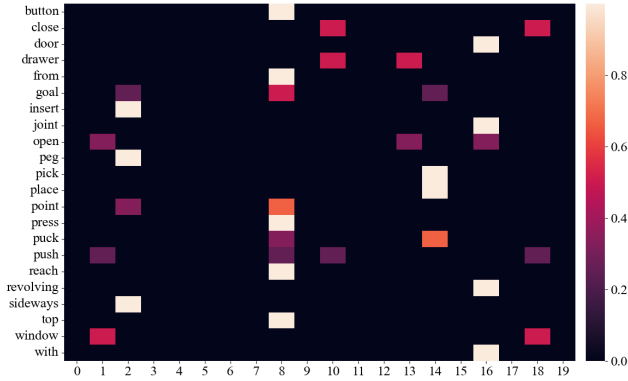


Figure 6. **Visualization of skill heat map on Meta-World Multi-Task 10 (MT10).** There are eight skills learned by our method. (zoom in for best view)

## C. More Visualizations

### C.1. Visualization Results of Learned Skill Set

As mentioned before, we show the visualization results of skill set on LOReL Sawyer Dataset in Fig. 5 and Meta-World Multi-Task 10 (MT10) in Fig. 6. The visualization results show that out of a 20-size skill-set, our SkillDiffuser learned 11 skills for LOReL (e.g. *pull drawer handle* [skill 0], *shut close container drawer* [skill 15], etc.) and 8 skills for Meta-World MT10 (e.g. *open push window* [skill 0], *open door with revolving joint* [skill 16], etc.). The results demonstrate strong skill abstraction abilities. For example, the skill “shut close container drawer” abstracts different expressions like “shut drawer”, “shut container” into one skill semantic. In the heatmap, the presence of distinct bright spots across eleven columns strongly reaf-

firms the model’s capability to discern and pinpoint specific skills from visual inputs, in the absence of a pre-defined skill library. This observation is not just a testament to the model’s enhanced interpretative prowess over conventional diffusion-based planning approaches but also marks a remarkable stride in abstracting high-level skills into representations that are intuitively understandable by humans. Such evidence further validates the model’s proficiency in sophisticated skill identification and representation.

### C.2. Word Cloud of Learned Skills

We further show the word cloud of 8 learned skills of LOReL Sawyer Dataset in Figure 7. From the results, we can find that the model has successfully mastered eight key skills, each closely linked to specific tasks. These skills demonstrate strong robustness to ambiguous language instructions. For instance, skill 4 effectively abstracts the skill of “open a drawer” from ambiguous expressions such as “open a container”, “pull a dresser”, “pull a drawer” and random combinations of these words. Similarly, skill 6 extracts the skill of “turn a faucet to the left”. This analysis indicates our method’s resilience to varied and poorly defined language inputs, confirming our SkillDiffuser can competently interpret and act upon a wide range of linguistic instructions, even those that are ambiguous or incomplete. These findings provide new perspectives and methodological guidance for future research in similar fields, especially in handling complex tasks with ambiguous language instructions. We also provide the word cloud of learned skills from Meta-World MT10 dataset in Fig. 8.

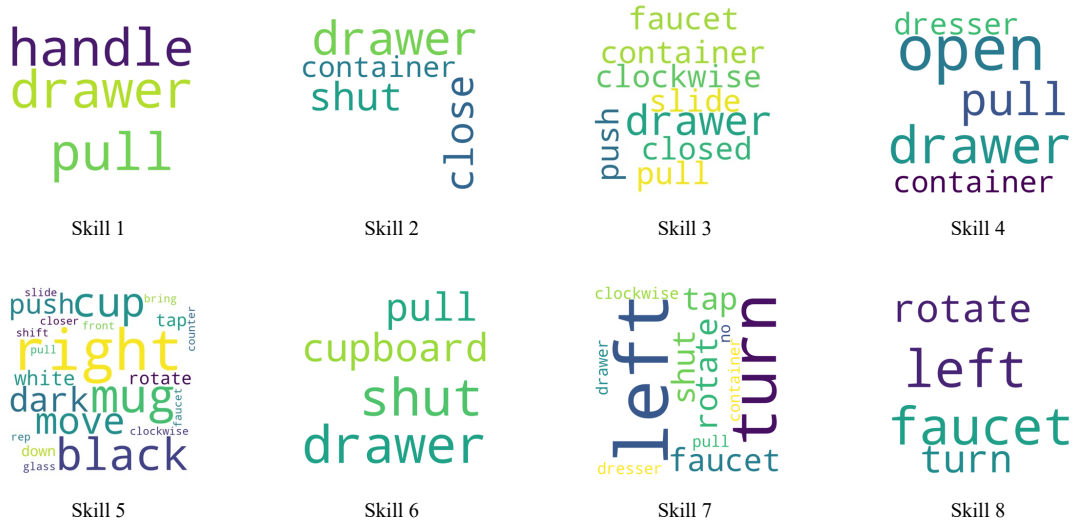


Figure 7. **Word cloud of learned skills in LOReL Sawyer Dataset.** We show eight of them here with the size corresponding to the word frequency in one skill.



Figure 8. **Word cloud of learned skills in Meta-World MT10 Dataset.** We show eight of them here with the size corresponding to the word frequency in one skill.

## D. Dataset Descriptions

### D.1. LOReL Sawyer Dataset

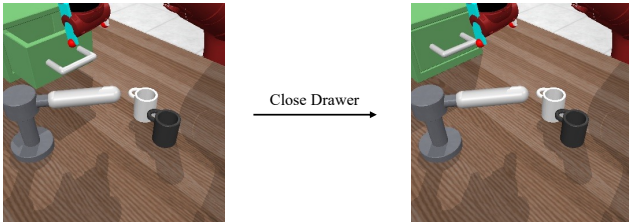


Figure 9. **A sample instance of LOReL Sawyer Dataset.** The start and goal images correspond to the instruction “close drawer”.

Language-conditioned Offline Reward Learning dataset, abbreviated as LOReL [29], contains trajectories originating from a reinforcement learning buffer which is generated by a random policy. The trajectories are sub-optimal and have language annotations through crowd-sourcing. Overall, the dataset encompasses approximately 50,000 language-annotated trajectories, each within a simulated environment featuring a Sawyer robot arm, with every demonstration extending over 20 discrete steps. A typical LOReL Sawyer environment is shown in Fig. 9. We assess our approach using the same set of instructions as those outlined in the original paper [29] which are described with their objectives in Tab. 6. These evaluation tasks are along with var-

Task	Description
Closing the Drawer	Involves the robot’s precise manipulation of a drawer to close it, testing spatial dynamics understanding and fine motor control.
Opening the Drawer	Requires the robot to open a drawer, emphasizing its capability in tasks that necessitate pulling and spatial navigation.
Turning the Faucet Left	Assesses the robot’s precision in rotational movements for turning a faucet to the left, a nuanced everyday action.
Turning the Faucet Right	Tests the robot’s adaptability in mirrored instructions, involving turning the faucet right, similar to the left turning task but in the opposite direction.
Pushing the Black Mug Right	Requires the robot to push a specific object (black mug) to the right, testing its skills in object recognition and directional movement.
Pushing the White Mug Down	Involves pushing a different object (white mug) downward, further evaluating the robot’s ability to differentiate objects and execute varied motion commands.

Table 6. Overview of tasks in LOReL Sawyer Dataset.

Instructions
open drawer and move black mug right
pull the handle and move black mug down
move white mug right
move black mug down
close drawer and turn faucet right
close drawer and turn faucet left
turn faucet left and move white mug down
turn faucet right and close drawer
move white mug down and turn faucet left
close the drawer, turn the faucet left and move black mug right
open drawer and turn faucet counterclockwise
slide the drawer closed and then shift white mug down

Table 7. LOReL composition tasks

ious rephrases of instructions which modify either the noun (“unseen noun”), the verb (“unseen verb”), both (“unseen noun+verb”), or entail a complete rewrite of the task (“human provided”), leading to a total of 77 distinct instructions for all six tasks. This structure of tasks and rephrases enables a comprehensive assessment of the robot’s ability to interpret and execute a wide range of language-based commands within the simulated environment.

## D.2. LOReL Composition Tasks

We follow the same settings as LISA [13] to create 12 new composition tasks through combining original evaluation instructions as shown in Tab. 7.

Additionally, we also incorporate tasks such as “move white mug right” and “move black mug down” to explore the composition of skills related to colors (e.g., black and white) and directions (e.g., right and down). This aims to explore whether such skills can be combined to fulfill complex instructions.

## D.3. Meta-World Dataset

The Meta-World dataset establishes a new benchmark in the field of multi-task and meta-reinforcement learning, offer-

Task Identifier	Language Instruction
window-close	push and close a window
window-open	push and open a window
door-open	open a door with a revolving joint
peg-insert-side	insert a peg sideways to the goal point
drawer-open	open a drawer
pick-place	pick a puck, and place the puck to the goal
reach	reach the goal point
button-press-topdown	press the button from the top
push	push the puck to the goal point
drawer-close	push and close a drawer

Table 8. Annotated instructions for Meta-World MT10 tasks.

ing 50 unique robotic manipulation tasks. These tasks range from simple to complex operations, providing researchers with a diverse testing ground. Each task is meticulously designed to ensure both challenge and common structural features that can be leveraged in multi-task and meta-learning algorithms. This design makes Meta-World an ideal choice for assessing the effectiveness and adaptability of algorithms in complex and variable task environments.

Particularly, the Multi-Task 10 (MT10) subset comprises 10 carefully selected tasks, where algorithms are trained and subsequently tested on the same set of tasks. As shown in Fig. 13, MT10 challenges algorithms’ learning and generalization capabilities in a multi-task environment, with the aim to evaluate the consistency and efficiency of algorithms in mastering multiple tasks, as well as their adaptability and robustness in the face of diverse tasks. As there is currently no widely-recognized instruction labeling of MT10, we provide our annotations here in Tab. 8.

We sample 100 trajectories for each task of MT10 and form the expert dataset of 1000 trajectories. We have released our dataset with image observations on <https://skilldiffuser.github.io>.



Figure 10. **Resulting images from applying skill 11 of Fig. 5.** The black dashed line is a horizontal reference and please pay attention to the red oval region. (zoom in for best view)

## E. More Ablations

### E.1. Ablation Study on Skill Interpretability

Resulting Images from Applying Discrete Skills We visualize resulting images from applying skill 11 of Fig. 5 which has grounding of “open, drawer, pull, dresser, container” (ranked from high frequency to low ones), consistent with its actual actions in Fig. 10. We can clearly observe a behavior of pulling the drawer. And we would like to clarify not all skills have clear semantic or action correspondences, while some do.

### E.2. Ablation Study on Condition Guidance Weight

Classifier-free guidance is widely used in generative model domain for its ability to act as temperature control when setting guidance weight above 1 during inference. In all of our experiments, we set guidance weight to 1.2 by default. But we also conduct ablation study on the condition guidance weight here in Tab. 9. From the results, we find the guidance weight slightly greater than 1 helps the planner’s performance, while excessive weight hurts.

Guidance Weight	1.0	1.2	1.8	3.0	5.0
Success Rate on Seen Tasks	39.33%	46.67%	38.86%	39.03%	33.50%

Table 9. Ablation on guidance weight. (5 episodes over 3 seeds.)

## F. More Results

### F.1. Task-wise Performance on LOReL Dataset

We further demonstrate the performance of our method and other baselines on LOReL Sawyer dataset in Fig 11 and 12. As can be seen from the figures, especially from Fig. 12, our method’s average performance on 5 rephrases is nearly 10 percentage points higher than the previous SOTA, which demonstrates its strong robustness against ambiguous language instructions.

### F.2. Task-wise Performance on Meta-World

We also provide the task-wise success rates on Meta-World MT10 dataset in Fig. 14, achieved by Flat R3M [30], Language-conditioned Diffuser and SkillDiffuser. The average performance is shown separately in the right figure.

From our experimental outcomes, it is clear to observe that our SkillDiffuser demonstrates commendable performance, particularly excelling in tasks involving mirrored instructions. SkillDiffuser exhibits an average performance enhancement of over 5% than previous language-conditioned Diffuser, which highlights the model’s advanced capability in understanding complex and ambiguous instructions compared to traditional methods. It showcases SkillDiffuser’s superior use of hierarchical architecture that employs interpretable skill learning for diffusion-based planners to better generate future trajectories.

## G. Implementation Details

### G.1. Hyper-parameters

Generally, we follow the settings illustrated in [13] with details specified in the following Tab. 10.

Hyper-parameter	LOReL	Meta-World
Skill Predictor Transformer Layers	1	1
Skill Predictor Embedding Dim	128	128
Skill Predictor Transformer Heads	4	4
Skill Set Code Dim	16	16
Skill Set Size	20	20
Dropout	0.1	0.1
Batch Size	256	64
Skill Predictor Learning Rate	1e-6	1e-5
Conditional Diffuser Learning Rate	1e-3	5e-3
Condition Guidance Weight	1.2	1.2
Inverse Dynamics Model Learning Rate	1e-3	5e-4
Diffuser Loss Weight	0.005	0.01
Horizon	8	8
VQ EMA Update	0.99	0.99
Skill Predictor and Diffuser Optimizer	Adam	Adam
Inverse Dynamics Model Optimizer	Adam	Adam

Table 10. **Hyper-parameters of SkillDiffuser.**

### G.2. Architecture Details

1. We use 1 layer Transformer network for the skill predictor and follow the implementation of VQ-VAE [45] to achieve VQ operation.
2. The size of skill set is set to 20 and the planning horizon is set to 8 for all implementations.
3. A temporal U-Net [36] with 6 repeated residual blocks is employed to model the noise  $\epsilon_\theta$  of the diffusion process. Each block is comprised of two temporal convolutions, each followed by group norm [47], and a final Mish non-linearity [26]. Timestep and skill embeddings are generated by two separate single fully-connected layer and added to the activation output after the first temporal convolution of each block.



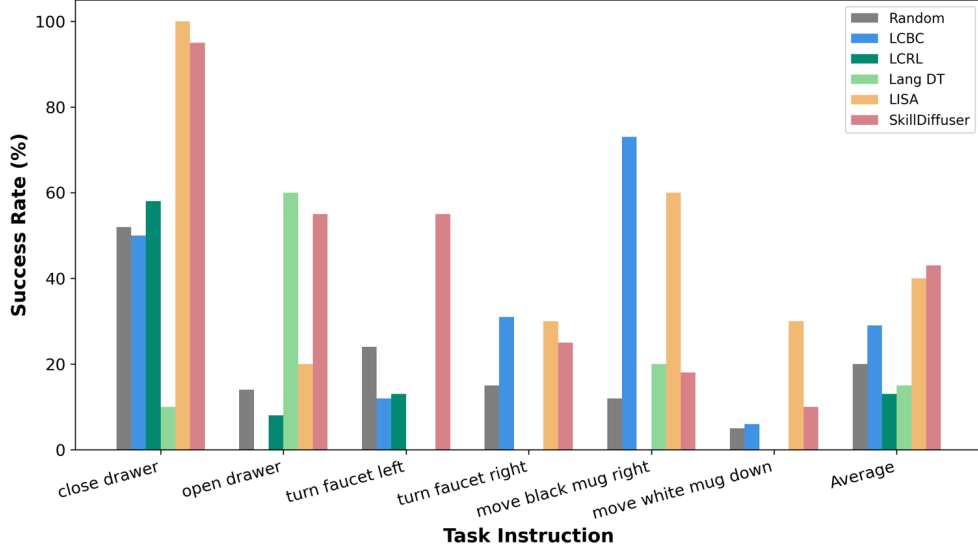


Figure 11. **Task-wise success rates (in %) on LOReL Sawyer Dataset.**

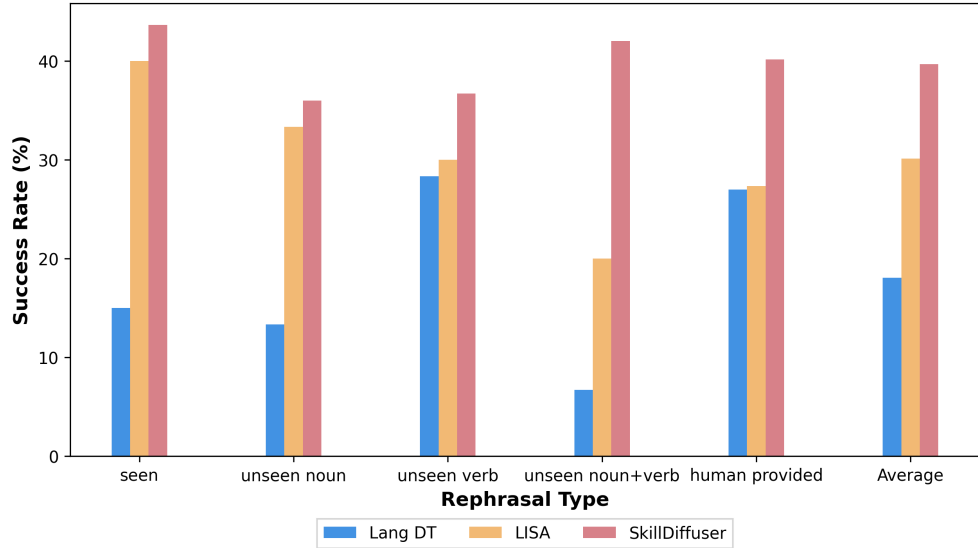


Figure 12. **Rephrasal-wise success rates (in %) on LOReL Sawyer Dataset.**

### G.3. Training Details

1. We train our model with one NVIDIA A100 Core Tensor GPU for about 45 hours in LOReL Sawyer dataset and about 24 hours in Meta-World MT10 dataset (1000 trajectories in total).
2. In both LOReL and Meta-World dataset, the skill predictor and diffusion model are trained with Adam optimizer [21] using a learning rate of  $1 \times 10^{-3}$  for the diffusion model,  $1 \times 10^{-6}$  for the LOReL skill predictor while  $1 \times 10^{-5}$  for Meta-World skill predictor. We only update parameters of Meta-World skill predictor every ten iterations. The inverse dynamics model is updated with Adam optimizer as well.
3. The batch size is set to 256 for LOReL Sawyer dataset and 64 for Meta-World MT10 dataset.
4. The training steps of the diffusion model are 5K for LOReL Sawyer dataset and 8K for Meta-World MT10 dataset. And the training epochs of the skill predictor are 500 for both datasets.
5. The planning horizon  $T$  of diffusion model is set to 100 and the denoising steps are set to 200 for all tasks.

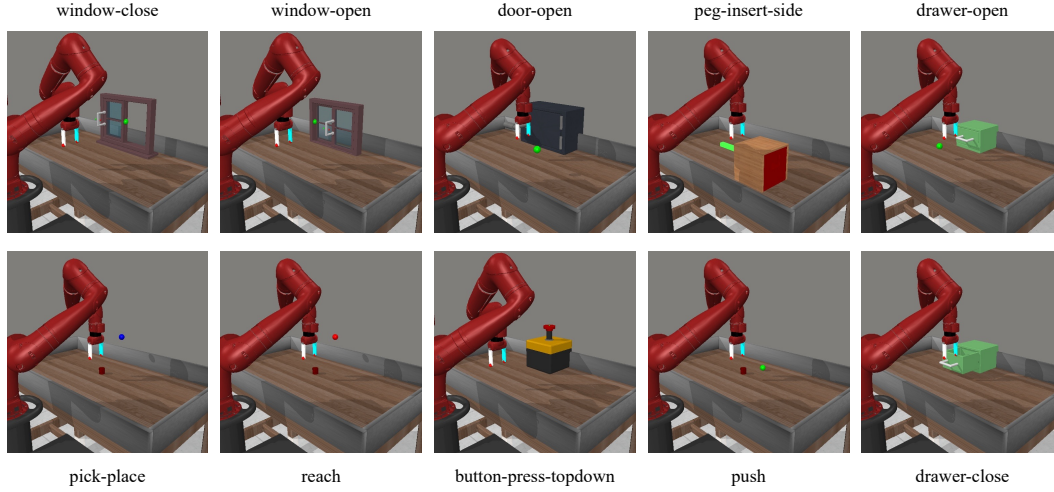


Figure 13. Partially visual observations of all the 10 tasks in Meta-World MT10 Dataset.

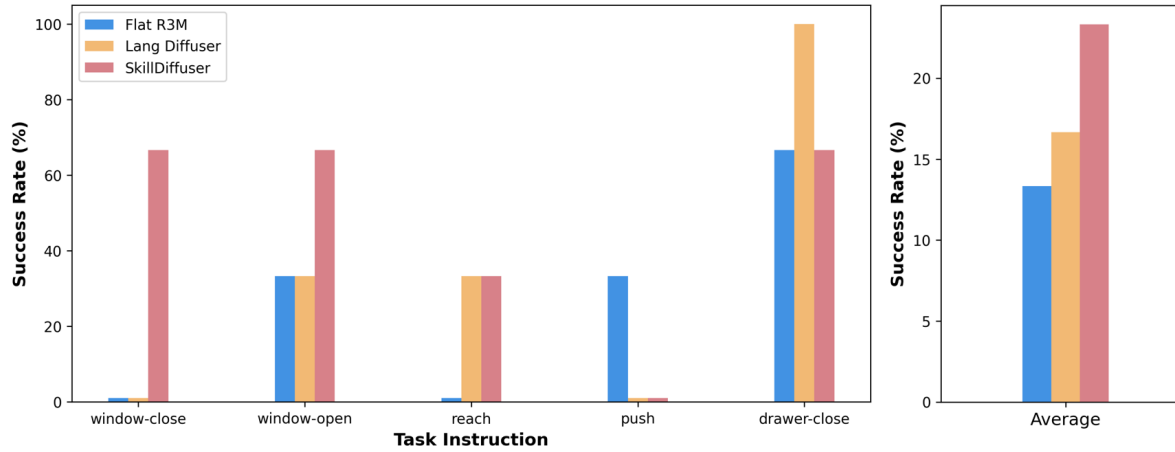


Figure 14. Task-wise success rates (in %) on Meta-World MT10 Dataset.