
Language Control Diffusion: Efficiently Scaling through Space, Time, and Tasks

Anonymous Authors¹

Abstract

Training generalist agents is difficult across several axes, requiring us to deal with high-dimensional inputs (space), long horizons (time), and multiple and new tasks. Recent advances with architectures have allowed for improved scaling along one or two of these dimensions, but are still prohibitive computationally. In this paper, we propose to address all three axes by leveraging Language to Control Diffusion models as a hierarchical planner conditioned on language (LCD). We effectively and efficiently scale diffusion models for planning in extended temporal, state, and task dimensions to tackle long horizon control problems conditioned on natural language instructions. We compare LCD with other state-of-the-art models on the CALVIN language robotics benchmark and find that LCD outperforms other SOTA methods in multi task success rates while dramatically improving computational efficiency with a single task success rate (SR) of 87% against the previous best of 82.7%. We show that LCD can successfully leverage the unique strength of diffusion models to produce coherent long range plans while addressing their weakness at generating low-level details and control.

1. Introduction

Generalist agents are characterized by the ability to plan over long horizons, understand and respond to human feedback, and generalize to new tasks based on that feedback. Language conditioning is an intuitive way to achieve human-robot alignment for task specification, with built-in structure allowing for generalization to new tasks, as described by a new language command. There have been many recent developments to leverage language for robotics and downstream decision-making and control (Ahn et al., 2022; Li et al., 2022a; Mees et al., 2022a; Lynch & Sermanet, 2020a; Jiang et al., 2019; Colas et al., 2020). However, while language is useful for task specification and generalization, it will not necessarily help with planning over long horizons.

One promising candidate that has emerged for long horizon planning is denoising diffusion models. Text to image diffusion models (Nichol et al., 2022; Ramesh et al., 2022; Ho et al., 2020; 2022c) have recently been able to successfully take advantage of large language models (LLMs) to generate incredibly detailed scenes. In addition, diffusion models have recently been proposed and successfully applied for low-dimensional, long-horizon planning and offline reinforcement learning (RL) (Janner et al., 2022; Wang et al., 2022b). In particular, Diffuser (Janner et al., 2022) has emerged as an especially promising planner with several properties suited for language conditioned control: flexibility in task specification, temporal compositionality, and ability to scale to long-horizon settings.

Current methods have a few pitfalls. Many existing language-conditioned control methods assume access to a high level discrete action space (e.g. switch on the stove, walk to the kitchen) provided by a lower level skill oracle (Ahn et al., 2022; Huang et al., 2022a; Jiang et al., 2022; Li et al., 2022a). The generative LM will typically decompose some high level language instruction into a set of predefined skills, which are then executed by a control policy or oracle. A fixed set of predefined skills may preclude the ability to generalize to novel environments and tasks. Further, one of the most important design considerations for these models is deciding what level of abstraction the communication protocol between the LM and the policy should take. A tradeoff arises from this question, where one must decide between the amount of reasoning to distribute over the LM versus the underlying policy. We can increase the burden on the LM by having it output text at a low level of abstraction that is extremely simple to follow, or we can increase the burden on the policy by having it execute text at a higher level of abstraction. Much prior work has been done on increasing the reasoning load of the language encoder by leveraging the capabilities of LLMs to generate planning code (Liang et al., 2022), calculate affordances of separate skills (Ahn et al., 2022), or perform chain of thought reasoning during rollout (Huang et al., 2022b).

Diffusion models are also computationally prohibitive to run in high-dimensional input spaces and tend to be inaccurate at low-level control. Specifically, Diffuser (Janner

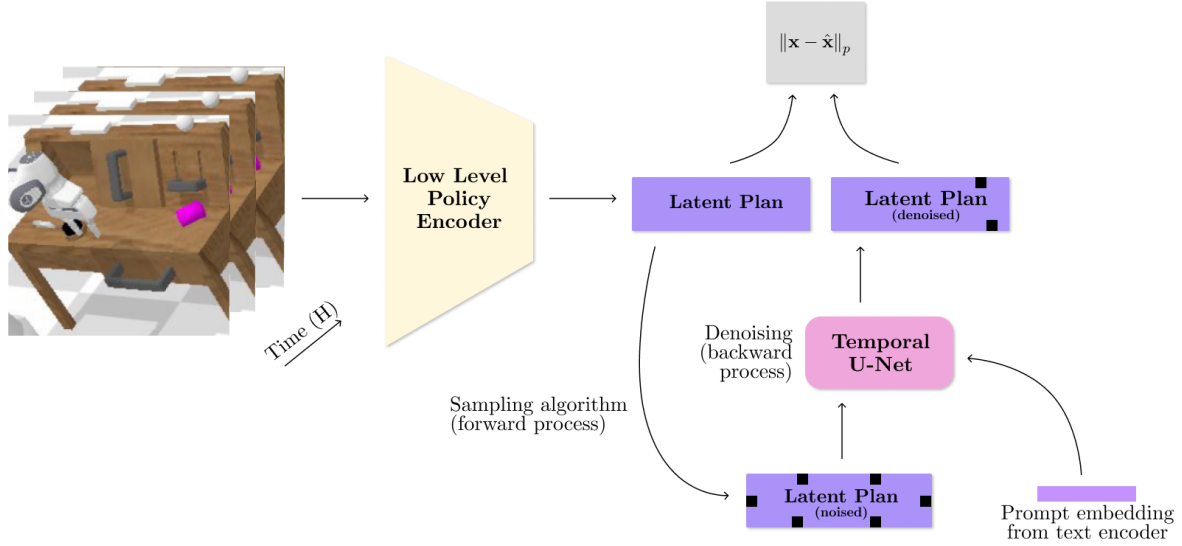


Figure 1: An overview of our training pipeline. The low level policy encoder is used to encode a subsampled horizon of RGB observations into a lower dimensional latent space, which will be used later on as goals for a goal-conditioned policy. We then noise the plan according to a uniformly sampled timestep from the diffusion process’ variance schedule, and train a Temporal U-Net to reverse this process conditioned on an natural language instruction from a frozen upstream large language model, effectively learning how to conditionally denoise the latent plan. To train the U-Net, one can simply use the p -norm between the predicted latent plan and the ground truth latent plan as the loss. We use $p = 1$ in practice following (Janner et al., 2022).

et al., 2022) does not work directly out of the box when scaling to pixels. This is perhaps unsurprising, as a high-dimensional Diffuser is effectively a text-to-video diffusion model, and even internet-scale video diffusion models have demonstrated a poor understanding of physics and temporal coherence over fine details (Ho et al., 2022a; Singer et al., 2022). This is due to the fact that the training objective for generative models has no notion of the underlying task, meaning they must model the entire scene with equal weighting. This includes potentially task-irrelevant details that may hinder or even prevent solving the actual control problem (Zhang et al., 2020; Wang et al., 2022a), which leads to catastrophic failure in control where fine details and precision are essential. Additionally, training a video diffusion model is generally computationally expensive and may take several days or weeks to train, which leaves such an approach out of reach to most researchers. This problem is exacerbated when considering that at inference time that multiple forward passes (often >20) are required for generating even a single state and action, meaning full text to video diffusion models are inefficient and likely impractical for the real-time sampling demands of robotics and control.

We address both of these issues by proposing the usage of a hierarchical diffusion model. By utilizing the representation of a goal-conditioned policy as a low level policy (LLP), we can effectively solve the representation issue by outputting states in a low-dimensional goal space directly into the

LLP. This also allows us to arbitrarily scale the difficulty of the low-level policy learning problem by controlling the horizon length from which the goal state is set from the current state, which has remained relatively underexplored for this application. This direction is promising, as it avoids defining the communication layer altogether and enables generating actions directly from high level text without a human-defined communication protocol in between. This hierarchical approach allows us to scale the diffusion model along three orthogonal axes:

- **Spatial dimension** through a low-dimensional representation that has been purposely optimized for control.
- **Time dimension** through a temporal abstraction enabled by utilizing a goal conditioned low level policy (LLP), as the LLP can use goal states several timesteps away from the current state.
- **Task dimension** through language, as the diffusion model acts as a powerful interpreter for plugging any large language model into control.

In addition, the entire pipeline is extremely fast and simple to train, as we utilize DDIM (Song et al., 2020a), a temporal abstraction on the horizon, as well as a low-dimensional representation for generation. We are able to achieve an average of 87.0% success rate across all tasks on the challenging

CALVIN benchmark. Additionally, we elucidate where diffusion models for text-to-control work well and highlight their limitations. Finally, we explore learning the alignment and grounding between language and state-actions via hierarchical diffusion for improving on the long-horizon and task compositionality properties of language conditioned planning. In summary, we list our core contributions below:

- We propose an effective method for scaling diffusion models for multitask long horizon control from pixels through the utilization of hierarchy
- A successful instantiation of our language control diffusion model, which achieves performance that outperforms the state of the art on the challenging CALVIN benchmark

2. Background

2.1. Reinforcement Learning

We formulate the environment as a Markov decision process (MDP) $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{R}, \gamma, p)$, with state space \mathcal{S} , action space \mathcal{A} , reward function \mathcal{R} , discount factor γ , and transition dynamics p . At each time step t , agents observe a state $s \in \mathcal{S} \subseteq \mathbb{R}^n$, take an action $a \in \mathcal{A} \subseteq \mathbb{R}^m$, and transition to a new state s' with reward r following $s', r \sim p(\cdot, \cdot | s, a)$. The goal of RL is then to learn either a deterministic policy $\pi : \mathcal{S} \mapsto \mathcal{A}$ or a stochastic policy where $a \sim \pi(\cdot | s)$ with $a \in \mathcal{A}$ that maximises the policy objective, $J(\pi) = \mathbb{E}_{a \sim \pi, s' \sim p} \sum_{t=0}^{\infty} \gamma^t r_t$.

2.2. Language Conditioned RL

We consider a language-conditioned RL setting where we assume that the true reward function \mathcal{R} is unknown, and must be inferred from a natural language instruction $L \in \mathcal{L}$. Formally, let \mathcal{F} be the function space of \mathcal{R} . Then the goal becomes learning a mapping from the language instruction to a reward function $\psi : \mathcal{L} \mapsto \mathcal{F}$, and maximizing the policy objective conditioned on the reward function $\psi(L)$, $J(\pi(\cdot | s, \mathcal{R})) = \mathbb{E}_{a \sim \pi, s \sim p} \sum_{t=0}^{\infty} \gamma^t r_t$. This formulation can be seen as a contextual MDP (Hallak et al., 2015). Note that the space of tasks that can be specified by language is much larger than that of reward, due to the Markov restriction of the latter (Abel et al., 2021). For example, “pour the milk” and “pour the milk after five o’clock” are both valid instructions, but are indistinguishable from a reward function if the state does not contain temporal information. We assume access to a prior collected dataset \mathcal{D} of N annotated trajectories $\tau_i = \langle (s_0, a_0, \dots, s_T), L_i \rangle$. The language conditioned policy π_β , or the behavior policy, is defined to be the policy that generates the aforementioned dataset. This general setting can be further restricted to prohibit environment interaction, which recovers offline RL or imitation

learning (IL). In this paper, we assume access to a dataset of expert trajectories, such that $\pi_\beta = \text{optimal policy } \pi^*$. Although this may seem like a strong assumption, the setting still poses a challenging learning problem as many unseen states will be encountered during evaluation and must be generalized to. Several prior methods have failed in this setting (Ross et al., 2011; De Haan et al., 2019).

2.3. Goal Conditioned Imitation Learning and Hierarchical RL

Goal-conditioned reinforcement learning (RL) is a subfield of RL that focuses on learning policies that can achieve specific goals or objectives, rather than simply maximizing the cumulative reward. This approach has been extensively studied in various forms in the literature (Arora & Doshi, 2018; Sodhani et al., 2021; Li et al., 2022a; Harrison et al., 2017; Schaul et al., 2015), although we focus on the hierarchical approach (Sutton et al., 1999). Following (Nachum et al., 2018), we formulate the framework in a two-level manner where a high level policy $\pi_{\text{hi}}(\tilde{a} | s)$ samples a goal state $g = \tilde{a}$ every c time steps, which we refer to as the temporal stride. A low level policy (LLP) $\pi_{\text{lo}}(a | s_t, g_t)$ then attempts to reach this state, which can be trained through hindsight relabelling (Andrychowicz et al., 2017) the offline dataset \mathcal{D} . In addition, as we consider high dimensional control from pixels, we factorize our low level policy formulation into $\pi_{\text{lo}}(a | s_t, g_t) := \phi(\mathcal{E}(s_t), g_t)$ where $z := \mathcal{E}(s_t)$ defines an encoder function that maps s into a latent representation z and ϕ translates the representation z into a distribution over actions a given a goal g . Note that ϕ is an arbitrary function, and can be made as complex as an another hierarchical model or as simple as a matrix.

3. Language Control Diffusion

3.1. Diffusion Training Objective

We now describe our problem formulation and framework in detail. Since we assume that our dataset is optimal, the policy objective reduces to imitation learning,

$$\min_{\pi} \mathbb{E}_{s, \mathcal{R} \sim \mathcal{D}} [D_{\text{KL}}(\pi_\beta(\cdot | s, \mathcal{R}), \pi(\cdot | s, \mathcal{R}))]. \quad (1)$$

As we tackle the problem from a planning perspective, we define a trajectory generator as \mathcal{P} and switch the atomic object from actions to trajectories τ . Thus we aim to minimize the following KL:

$$\begin{aligned} & \min_{\mathcal{P}} D_{\text{KL}}(\mathcal{P}_\beta(\tau | \mathcal{R}), \mathcal{P}(\tau | \mathcal{R})) \\ &= \min_{\mathcal{P}} \mathbb{E}_{\tau, \mathcal{R} \sim \mathcal{D}} [\log \mathcal{P}_\beta(\tau | \mathcal{R}) - \log \mathcal{P}(\tau | \mathcal{R})]. \end{aligned} \quad (2)$$

This can be reformulated into the following diffusion training objective:

Algorithm 1 Hierarchical Conditional Diffusion Training

Input: baseline goal-conditioned policy $\pi_{\text{lo}} := \phi(\mathcal{E}(s_t), g_t)$, diffusion variance schedule α_t , temporal stride c , language model ρ
Output: trained hierarchical policy $\pi(a_t|s_t) := \pi_{\text{lo}}(a_t|s_t, \pi_{\text{hi}}(g_t|s_t))$, where g_t is sampled from π_{hi} as first state in τ_c every c time steps

- 1: Collect dataset $\mathcal{D}_{\text{onpolicy}}$ by rolling out trajectories $\tau \sim \pi_{\text{lo}}, \rho$
- 2: Instantiate π_{hi} as diffusion model $\epsilon_{\theta}(\tau_{\text{noisy}}, t, \rho(L))$
- 3: **repeat**
- 4: Sample mini-batch $(\tau, L) = B$ from $\mathcal{D}_{\text{onpolicy}}$.
- 5: Subsample $\tau_c = (\mathcal{E}(s_0), \mathcal{E}(s_c), \mathcal{E}(s_{2c}), \dots, \mathcal{E}(s_T))$.
- 6: Sample diffusion step $t \sim \text{Uniform}(\{1, \dots, T\})$, noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 7: Update high level policy π_{hi} with gradient $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\alpha_t}\tau_c + \sqrt{1 - \alpha_t}\epsilon, t, \rho(L))\|_p$
- 8: **until** converged

$$\mathbb{E}_{\tau_0, \epsilon} [\|\epsilon_t - \epsilon_{\theta}(\sqrt{\alpha_t}\tau_0 + \sqrt{1 - \alpha_t}\epsilon_t, t)\|^2]. \quad (3)$$

We refer to Appendix C for a more detailed derivation of our objective.

3.2. Near Optimality Guarantees

In order to theoretically justify our approach and show that we can safely do diffusion in the LLP encoder latent space without loss of optimality up to some ϵ , we show that a near optimal policy is always recoverable with the LLP policy encoder under some mild conditions: that our low-level policy has closely imitated our expert dataset and that the transition dynamics are smooth.

Let $\pi_{\text{lo}}(s) := \phi \circ \mathcal{E}(s)$ be a deterministic low-level policy. Similar to Nachum et al. (2019), we can then define the *sub-optimality* of the action and state space induced by $\tilde{\mathcal{A}} = \tilde{\mathcal{S}} = \mathcal{E}(s)$ to be the difference between the best possible high level policy $\pi_{\text{hi}}^* = \arg\max_{\pi \in \Pi} J(\pi_{\text{hi}})$ within this new abstract latent space and the best policy in general.

$$\text{SubOpt}(\mathcal{E}, \phi) = \sup_{s \in \tilde{\mathcal{S}}} V^{\pi^*}(s) - V^{\pi_{\text{hi}}^*}(s). \quad (4)$$

Intuitively, this captures how much potential performance we stand to lose by abstracting our state and action space.

Proposition 3.1. *If the transition function $p(s'|s, a)$ is Lipschitz continuous with constant K_f and $\sup_{s \in \mathcal{S}, a \in \mathcal{A}} |\pi_{\text{lo}}(s) - a^*| \leq \epsilon$, then*

$$\text{SubOpt}(\mathcal{E}, \phi) \leq C\epsilon, \quad (5)$$

where $C = \frac{2\gamma}{(1-\gamma)^2} R_{\max} K_f \text{dom}(P(s'|s, a))$.

We give a detailed proof in Appendix D.

3.3. Low-Level Policy

We adopt the HULC architecture (Mees et al., 2022a), which is in turn modeled from multi-context imitation learning (MCIL) (Lynch & Sermanet, 2020b). MCIL is a method of training agents to perform a task by imitating multiple experts in different contexts. The goal is to learn a policy that can generalize to unseen situations by training on a diverse set of expert demonstrations. The agent learns to adapt its behavior based on the context in which it is operating, rather than just memorizing the specific actions of the expert in a single context. One of their key insights is leveraging the use of a shared goal space between language and vision through transfer learning methods (Tan et al., 2018). This approach can improve the robustness and flexibility of the trained agent. We utilize this shared latent goal space as the space for diffusion, primarily because it enjoys low dimensionality (32 dimensions), offering an efficient communication layer between the high level and low level policy.

In HULC, the authors propose an improved version of Multi-Context Imitation Learning (MCIL) by decomposing control into a hierarchical approach. Note that this means that our method is a successful instantiation of a three-level hierarchical model, as we generate the goal states with diffusion *into their high-level policy*. HULC utilizes hierarchy by generating global discrete latent plans and learning local policies that are conditioned on this global plan. Their method also focuses on several key components that have a significant impact on performance, including the architectures used to encode sequences in relabeled imitation learning, the representation of latent distributions, the alignment of language and visual representations, and data augmentation and optimization techniques.

3.4. Model Architecture

It is computationally infeasible to operate directly on the pixel space. Instead, we do planning in a latent state space by utilizing the representation of a goal-conditioned low-level policy. We adopt T5 (Raffel et al., 2020) as our textual encoder, which has similarly found success in the text to image diffusion model Imagen (Ho et al., 2022a). Specifically, we use the Transformer variant (Vaswani et al., 2017) as the text encoder with the architecture modifications described in Radford et al. (2019).

We utilize the temporal U-Net (Janner et al., 2022) which performs 1D convolution only across the time dimension, rather than the 2D convolution typical in text-to-image generation. This is motivated by our wish to preserve equivariance along the time dimension but not the state-action dimension. In addition, we modify the architecture in Janner et al. (2022) by adding conditioning via cross attention in a fashion that resembles the latent diffusion model. When we

instantiate the model, we also found that conditioning on the first state through cross attention to be simpler than passing in through inpainting as is done in the original Diffuser, and faster.

We use DDIM (Song et al., 2020a) during inference for increased computational efficiency and faster planning. DDIM uses strided sampling and is able to capture almost the same level of fidelity as typical DDPM sampling (Ho et al., 2020) with an order of magnitude speedup. For rolling out the latent plans generated by the denoiser, we resample a new sequence of goals g with the frequency of c , until either the task is completed successfully or the maximum timestep is reached. Our low-level policy takes over control between samples, with goals generated by the high level policy as input.

4. Experiments

4.1. Dataset and Metric

In our experiments we aim to answer the following questions: 1) Does diffusion work well for language-conditioned RL? 2) Is a hierarchical instantiation of the diffusion model really necessary? Finally, 3) how much efficiency is gained by planning in a latent space? We evaluate on the CALVIN benchmark (Mees et al., 2022b). After pretraining our low-level policy on all data, we freeze the policy encoder and train the diffusion U-Net using the frozen encoder. This setup is motivated by the fact that modeling a changing data distribution with the diffusion model is not theoretically sound. We roll out all of our evaluated policies for 100 trajectories on all 34 tasks, and all comparisons are evaluated in their official repository¹.

4.2. Baselines

As introduced in Section 3.3, we compare against the prior state of the art, HULC and MCIL (Mees et al., 2022a; Lynch & Sermanet, 2020b). To ensure a fair comparison, we rigorously follow their evaluation procedure and build directly from their codebase. All baseline results are taken from Mees et al. (2022a).

4.3. Performance of LCD

We outperform prior methods on the challenging multitask long-horizon control (MT-LHC) benchmark, and improve on the strongest prior model’s average performance on horizon length one tasks by 6.1% as shown in 1. In order to further elucidate why our method works better than HULC, we do a deeper analysis on several failure cases that we observed and show how LCD corrects them. In addition, we significantly outperform HULC in 15 of the 34 tasks,

outperform HULC in 23 of the 34 tasks, and improve on the average success rate of single tasks by 3.33% as shown in Table 3. Note that the average success rate differs from Table 1, as the distribution of tasks for MT-LHC is not uniform due to the fact that the Calvin benchmark filters out infeasible trajectories. A visualization of the first task MT-LHC distribution is provided for future work in Appendix I.

Table 1: Our main result: a comparison of success rates between our diffusion model and prior benchmarks on multitask long-horizon control (MT-LHC) for 34 disparate tasks. We report the mean and standard deviation across 3 seeds for our method with each seed evaluating for 1000 episodes.

Horizon	GCBC	MCIL	HULC (LLP only)	Diffuser	Ours
1	64.7(4.0)	76.4(1.5)	82.6(2.6)	1	88.7(1.5)
2	28.4(6.2)	48.8(4.1)	64.6(2.7)	2	69.9(2.8)
3	12.2(4.1)	30.1(4.5)	47.9(3.2)	3	54.5(5.0)
4	4.9(2.0)	18.1(3.0)	36.4(2.4)	2	42.7(5.2)
5	1.3(0.9)	9.3(3.5)	26.5(1.9)	2	32.2(5.2)
Avg	1.11(0.3)	1.82(0.2)	2.6(0.1)	2	2.88(0.19)

Table 2: Wall clock times for training. All numbers are taken from our own experiments and server for reproducibility and fairness, including baselines. Baselines are run with 8 2080ti GPUs following the original author guidelines, whilst our experiments are run with a single 2080ti. Latent dims denotes the size of the latent space that we perform the diffusion generation in. We compare against two variants of Diffuser. Diffuser-1D is the same model as presented in Table 1 which utilizes a VAE trained from scratch on the dataset, whilst Diffuser-2D utilizes a large pretrained VAE from Stable Diffusion (Rombach et al., 2022). Inference time (sec) refers to the average amount of time taken in seconds to produce an action.

	Ours (HLP only)	HULC	Ours (full)	Diffuser-1D	Diffuser-2D
Training (hrs)	13.3	82	95.3	20.8	49.2
Inference time (sec)	0.333	0.005	0.336	1.11	5.02
Model size	20.1M	47.1M	67.8M	74.7M	125.5M
Latent dims	32	N/A	32	512	16x16x4
Avg ∇ updates/sec	6.25	.5	6.25	4	2.1

4.4. High-Level Policy: MLP vs. Diffusion

We give the MSE training loss plot between Multilayered Perceptron (MLP) and Diffusion (LCD) in Figure 2. The MSE here measures the error between the next goal state prediction from current goal state.

4.5. Robustness to Hyperparameters

In Figure 3 we show that our diffusion model method is robust to several hyperparameters including the number of diffusion steps, a frame offset o , and the hidden dimensions of the model. The number of diffusion steps is typically a key parameter determining the quality of the representations.

¹<https://github.com/lukashermann/hulc>

Table 3: Comparison of success rates (SR) across single tasks, evaluated with HULC and LCD. Differences of over 5% in SR are bolded. We again report the performance for the mean and standard deviation across 3 seeds for all methods, with runs taken from the MT-LHC benchmark in Table 1.

Task	HULC	Ours
Close Drawer	99.50(0.36)	95.97(3.46)
Lift Blue Block Drawer	86.97(13.01)	90.83(9.47)
Lift Blue Block Slider	58.80(6.99)	78.77(5.20)
Lift Blue Block Table	72.73(7.89)	76.07(2.03)
Lift Pink Block Drawer	67.93(9.39)	76.27(22.80)
Lift Pink Block Slider	71.50(4.54)	79.30(2.01)
Lift Pink Block Table	68.00(5.60)	70.47(3.07)
Lift Red Block Drawer	91.67(11.79)	77.00(4.71)
Lift Red Block Slider	71.97(2.57)	79.70(7.32)
Lift Red Block Table	60.87(5.84)	69.43(2.77)
Move Slider Left	96.27(1.14)	96.27(1.96)
Move Slider Right	99.27(0.24)	97.67(2.19)
Open Drawer	97.83(1.43)	98.07(0.50)
Place In Drawer	95.93(1.37)	95.90(1.79)
Place In Slider	70.90(5.44)	81.10(0.14)
Push Blue Block Left	60.10(9.69)	71.93(6.82)
Push Blue Block Right	27.60(9.10)	33.10(2.86)
Push Into Drawer	72.13(5.74)	75.47(5.29)
Push Pink Block Left	66.73(17.27)	64.80(11.98)
Push Pink Block Right	57.67(4.46)	51.17(6.04)
Push Red Block Left	64.17(4.71)	55.30(16.32)
Push Red Block Right	28.73(16.07)	37.83(3.27)
Rotate Blue Block Left	59.63(7.38)	70.03(2.56)
Rotate Blue Block Right	62.57(4.17)	72.77(6.76)
Rotate Pink Block Left	71.53(6.99)	75.07(9.15)
Rotate Pink Block Right	55.93(4.29)	64.00(4.84)
Rotate Red Block Left	67.43(5.15)	81.23(8.86)
Rotate Red Block Right	69.53(5.44)	74.57(2.11)
Stack Block	36.83(2.84)	32.00(3.56)
Turn Off Led	98.30(1.87)	96.87(1.80)
Turn Off Lightbulb	98.33(1.25)	95.67(3.78)
Turn On Led	98.73(0.93)	97.63(1.73)
Turn On Lightbulb	98.97(0.40)	95.60(1.58)
Unstack Block	94.53(4.53)	90.80(3.72)
Unstack Block	94.53(4.53)	90.80(3.72)
Avg Success Rate	70.68(2.65)	74.01(2.64)

However, we found that the performance of our method is relatively insensitive to the number of function evaluations. Additionally, we consider augmenting our dataset by adding a frame offset o during the sampling of our goal state s_{c+o} for potentially improving generalization. However, we find that this effect is more or less negligible. Finally, our method is also robust with respect to the number of parameters, and instantiating a larger model does not induce overfitting.

4.6. Summary of critical findings

Diffuser fails to plan in original state space

After investigating the performance of diffuser, our findings indicate that diffuser fails to successfully replan in high dimensional state spaces when using a Variational Autoencoder (VAE). Based on our observations, we hypothesize that the failure of diffuser to replan successfully is due to

Table 4: Ablation of our method by comparing against a simple three layer MLP as high level policy. We use the same methodology as in Table 1, and report the mean and standard deviation across 3 seeds for our method with each seed evaluating for 100 episodes. Even given the ground truth validation language embeddings, the MLP underperforms the original model and LCD.

Horizon Length	MLP	Ours
1	82.1(4.0)	88.7(1.5)
2	61.0(6.2)	69.9(2.8)
3	49.2(4.1)	54.5(5.0)
4	32.1(2.0)	42.7(5.2)
5	25.6(0.9)	32.2(5.2)
Avg. length	2.59(0.01)	2.81(0.12)



Figure 2: Training Loss Comparison. We find that the DM strictly lower bounds the MLP in error, which corresponds to the observed drop in performance when doing evaluation with MLP. With this ablation, we show that predicting the next goal state is non-trivial learning problem and that utilizing the diffusion model is necessary to effectively model this distribution. We further validate this by evaluating the MLP in 4.

the diffusion objective enabling interpolation between low density regions in the VAE’s latent space by the nature of the training objective smoothing the original probability distribution of trajectories. In practice, this interpolation between low density regions correspond to physically infeasible trajectories. This suggests that interpolation between low density regions in the VAE’s latent space is a significant factor in the failure of diffuser to plan and replan successfully. Our results highlight the need for better representation, and for further research of scaling diffusion models to high dimensional control.

Representation is essential for effective diffusion

We find that having a good representation is essential for

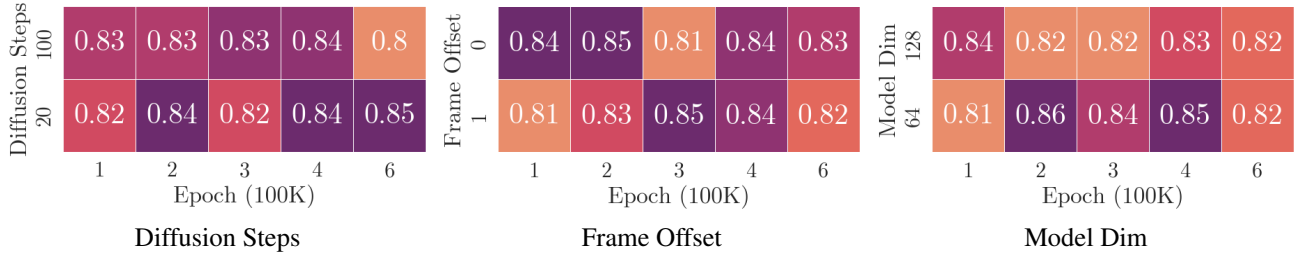


Figure 3: Single task success rates from MT-LHC across different epochs and hyperparameters. Our diffusion model is robust to various hyperparameters such as diffusion steps, frame offset, hidden dimensions, and number of parameters. While the number of diffusion steps usually affects the quality of representations, our method is relatively insensitive to it. Adding a frame offset during sampling to augment the dataset has little effect on generalization. Furthermore, a larger model does not lead to overfitting.

effective diffusion in general control settings. A good representation greatly eases the learning difficulty on the diffusion model by reducing the number of dimensions which in turn increases the information density, making it easier for the model to learn the underlying dynamics of the system. The good representation likely massages the regions of low density latent space between feasible trajectories into areas that still correspond to physically feasible regions when decoded by the low level policy. This enables the model to generalize well to unseen situations, improving its robustness and flexibility. Additionally, by having a good representation, the diffusion model is able to successfully generalize to new scenarios, which is crucial for the model’s effectiveness in a real-world application. Our results highlight the importance of latent representation in the diffusion process.

The encoder representation of deep goal-conditioned RL models can be effectively used for hierarchical RL

Our results imply that encoder representation of deep reinforcement learning models have the potential to be effectively utilized in hierarchical reinforcement learning. In modern deep policy learning from pixels, we propose an effective task-aware representation can be extracted by using the latent space of an intermediate layer as an encoder. Although we have demonstrated the effectiveness of this approach by instantiating it successfully in a hierarchical diffusion model, this approach is general and can be applied to any generative model in hierarchical RL. our results suggest that using a shared policy encoder between the high and low level policies can improve the effectiveness and efficiency of generative modelling in hierarchical RL.

5. Related Work

Text-to-Control Models Text-to-control models, also known as language-conditioned policies (Lynch & Sermanet, 2020a; Jiang et al., 2019; Colas et al., 2020) have been explored in the reinforcement learning community for

Table 5: Task generalization of LCD on a collection of five held out tasks. Note that the success rates given here differ from Table 3, as the initial state distribution also varies. The initial states here come directly from CALVIN’s initial state generator, whereas the initial states in 3 may arise from the last state of a prior task.

Task	Ours
Lift Pink Block Table	55.00(16.33)
Lift Red Block Slider	88.33(8.50)
Push Red Block Left	35.00(7.07)
Push Into Drawer	90.00(10.80)
Rotate Blue Block Right	36.67(14.34)
Avg SR	61.00(7.79)

improving generalization to novel tasks and new environments (Hermann et al., 2017; Bahdanau et al., 2018; Hill et al., 2019; Colas et al., 2020). Although they are not the only way to incorporate external knowledge in the form of text to decision making tasks (Luketina et al., 2019; Zhong et al., 2019), they remain one of the most popular as they are simple to evaluate and enjoy great interest from several jointly related fields including Vision Language Navigation (VLN) (Wang et al., 2020; Zheng et al., 2022) and Embodied AI (Duan et al., 2022). However, language grounding in RL remains notoriously difficult, as language vastly underspecifies all possible configurations of a corresponding state (Quine, 1960). Modern text to control models often still struggle with long-horizon language commands and misinterpret the language instruction.

Hu et al. (2019) attempt to solve a long-horizon Real-Time Strategy (RTS) game with a hierarchical method utilizing language as the communication interface between the high level and low level policy, whilst Harrison et al. (2017) consider training a language encoder-decoder for policy shaping through critique, Hanjie et al. (2021) utilize an attention module conditioned on textual entities for strong generalization and better language grounding. Zhong et al.

(2022) propose a model-based objective to deal with sparse reward settings with language descriptions and Mu et al. (2022) also tackle the sparse reward settings through the usage of intrinsic motivation with bonuses added on novel language descriptions. However, only Hu et al. (2019) consider the long-horizon setting, whom do not consider a high-dimensional state space. Jang et al. (2022) carefully examine the importance of diverse and massive data collection in enabling task generalization through language and propose a FiLM conditioned CNN-MLP model (Perez et al., 2018). Much work has applied using more data and compute for control (Brohan et al., 2022; Jaegle et al., 2021; Reed et al., 2022). However, none of these methods consider the usage of diffusion models as the medium between language and RL.

Diffusion Models Diffusion models such as DALL-E 2 (Ramesh et al., 2022) and GLIDE (Nichol et al., 2022) have recently shown promise as generative models, with state-of-the-art text-to-image generation results demonstrating a surprisingly deep understanding of semantic relationships between objects and the high fidelity generation of novel scenes. Stable diffusion, an instantiation of latent diffusion (Rombach et al., 2022), has also achieved great success and is somewhat related to our method as they also consider performing the denoising generation in a smaller latent space, albeit with a variational autoencoder (Kingma & Welling, 2013) rather than a low level policy encoder.

A key driver of this recent success in generative modeling is the usage of *classifier-free guidance*, which is amenable to the RL framework through the usage of language as a reward function. The language command is a condition for generation, which can be thought of as a reward function, as both are simply different forms of task specification. This perspective has direct connections to control as inference (Levine, 2018).

Given the success of denoising diffusion probabilistic models (Ho et al., 2020) in text-to-image synthesis (Sohl-Dickstein et al., 2015a), the diffusion model has been further explored in both discrete and continuous data domains, including image and video synthesis (Ho et al., 2022b;c), text generation (Li et al., 2022b), and time series (Rasul et al., 2021). Video generation are especially relevant to this work, as they are a direct analogue of diffusion planning model Diffuser (Janner et al., 2022) in pixel space without actions.

Diffuser first proposed to transform decision making into inpainting and utilize diffusion models to solve this problem, which much work has followed up on (Dai et al., 2023; Chi et al., 2023; Ajay et al., 2022). Specifically, they diffuse the state and actions jointly for imitation learning and goal-conditioned reinforcement learning through constraints specified through classifier guidance, and utilize Diffuser

for solving long-horizon and task compositional problems in planning. Instead of predicting the whole trajectory for each state, (Wang et al., 2022b) apply the diffusion model to sample a single action at a time conditioned by state. However, neither of these works consider control from pixels, or utilizing language for generalization.

6. Conclusion

Learning atomic sub-skills is critical to solving the multi-task planning problem and enabling our method to scale to more complex and open environments through learned state and temporal abstractions. We explore reformulating the language-conditioned planning process as text-to-plan generation to better learn the alignment between language and state-action pairs. We further show that the strengths of diffusion models can be leveraged for long horizon plans, and their weakness at low-level detail generation can be managed through learning low-level, goal-conditioned policies with imitation learning. Experiments and qualitative analysis demonstrate the simplicity and effectiveness of our model, showing that LCD is able to achieve state-of-the-art performance on a competitive language-conditioned control benchmark from rich observations, over long-horizons, and generalize to new scenarios.

References

- Abel, D., Dabney, W., Harutyunyan, A., Ho, M. K., Littman, M., Precup, D., and Singh, S. On the Expressivity of Markov Reward. In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P. S., and Vaughan, J. W. (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 7799–7812. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/4079016d940210b4ae9ae7d41c4a2065-Paper.pdf>.
- Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan, K., Hausman, K., Herzog, A., Ho, D., Hsu, J., Ibarz, J., Ichter, B., Irpan, A., Jang, E., Ruano, R. J., Jeffrey, K., Jesmonth, S., Joshi, N. J., Julian, R., Kalashnikov, D., Kuang, Y., Lee, K.-H., Levine, S., Lu, Y., Luu, L., Parada, C., Pastor, P., Quiambao, J., Rao, K., Rettinghouse, J., Reyes, D., Sermanet, P., Sievers, N., Tan, C., Toshev, A., Vanhoucke, V., Xia, F., Xiao, T., Xu, P., Xu, S., Yan, M., and Zeng, A. Do as i can, not as i say: Grounding language in robotic affordances, 2022. URL <https://arxiv.org/abs/2204.01691>.
- Ajay, A., Du, Y., Gupta, A., Tenenbaum, J., Jaakkola, T., and Agrawal, P. Is conditional generative modeling all you need for decision-making? *arXiv preprint arXiv:2211.15657*, 2022.
- Andrychowicz, M., Wolski, F., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B., Tobin, J., Pieter Abbeel, O., and Zaremba, W. Hindsight experience replay. *Advances in neural information processing systems*, 30, 2017.
- Arora, S. and Doshi, P. A survey of inverse reinforcement learning: Challenges, methods and progress, 2018. URL <https://arxiv.org/abs/1806.06877>.
- Bahdanau, D., Hill, F., Leike, J., Hughes, E., Hosseini, A., Kohli, P., and Grefenstette, E. Learning to understand goal specifications by modelling reward. *arXiv preprint arXiv:1806.01946*, 2018.
- Brohan, A., Brown, N., Carbajal, J., Chebotar, Y., Dabis, J., Finn, C., Gopalakrishnan, K., Hausman, K., Herzog, A., Hsu, J., et al. Rt-1: Robotics transformer for real-world control at scale. *arXiv preprint arXiv:2212.06817*, 2022.
- Chi, C., Feng, S., Du, Y., Xu, Z., Cousineau, E., Burchfiel, B., and Song, S. Diffusion policy: Visuomotor policy learning via action diffusion. *arXiv preprint arXiv:2303.04137*, 2023.
- Colas, C., Akakzia, A., Oudeyer, P.-Y., Chetouani, M., and Sigaud, O. Language-conditioned goal generation: a new approach to language grounding for rl. *LAREL Workshop 2020*, abs/2006.07043, 2020.
- Dai, Y., Yang, M., Dai, B., Dai, H., Nachum, O., Tenenbaum, J., Schuurmans, D., and Abbeel, P. Learning universal policies via text-guided video generation. *arXiv preprint arXiv:2302.00111*, 2023.
- De Haan, P., Jayaraman, D., and Levine, S. Causal confusion in imitation learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Duan, J., Yu, S., Tan, H. L., Zhu, H., and Tan, C. A survey of embodied ai: From simulators to research tasks. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2):230–244, 2022.
- Hallak, A., Di Castro, D., and Mannor, S. Contextual markov decision processes. *arXiv preprint arXiv:1502.02259*, 2015.
- Hanjie, A. W., Zhong, V. Y., and Narasimhan, K. Grounding language to entities and dynamics for generalization in reinforcement learning. In *International Conference on Machine Learning*, pp. 4051–4062. PMLR, 2021.
- Harrison, B., Ehsan, U., and Riedl, M. O. Guiding reinforcement learning exploration using natural language. *arXiv preprint arXiv:1707.08616*, 2017.
- Hermann, K. M., Hill, F., Green, S., Wang, F., Faulkner, R., Soyer, H., Szepesvari, D., Czarnecki, W. M., Jaderberg, M., Teplyashin, D., et al. Grounded language learning in a simulated 3D world. *arXiv preprint arXiv:1706.06551*, 2017.
- Hill, F., Lampinen, A., Schneider, R., Clark, S., Botvinick, M., McClelland, J. L., and Santoro, A. Emergent systematic generalization in a situated agent, 2019.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. *arXiv preprint arxiv:2006.11239*, 2020.
- Ho, J., Chan, W., Saharia, C., Whang, J., Gao, R., Gritsenko, A., Kingma, D. P., Poole, B., Norouzi, M., Fleet, D. J., et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022a.
- Ho, J., Saharia, C., Chan, W., Fleet, D. J., Norouzi, M., and Salimans, T. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022b. URL <http://jmlr.org/papers/v23/21-0635.html>.
- Ho, J., Salimans, T., Gritsenko, A., Chan, W., Norouzi, M., and Fleet, D. J. Video diffusion models, 2022c. URL <https://arxiv.org/abs/2204.03458>.

- Hu, H., Yarats, D., Gong, Q., Tian, Y., and Lewis, M. Hierarchical decision making by generating and following natural language instructions. *Advances in neural information processing systems*, 32, 2019.
- Huang, W., Abbeel, P., Pathak, D., and Mordatch, I. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. *ICML*, 2022a.
- Huang, W., Xia, F., Xiao, T., Chan, H., Liang, J., Florence, P., Zeng, A., Tompson, J., Mordatch, I., Chebotar, Y., et al. Inner monologue: Embodied reasoning through planning with language models. *arXiv preprint arXiv:2207.05608*, 2022b.
- Jaegle, A., Borgeaud, S., Alayrac, J.-B., Doersch, C., Ionescu, C., Ding, D., Koppula, S., Zoran, D., Brock, A., Shelhamer, E., et al. Perceiver io: A general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021.
- Jang, E., Irpan, A., Khansari, M., Kappler, D., Ebert, F., Lynch, C., Levine, S., and Finn, C. Bc-z: Zero-shot task generalization with robotic imitation learning. In *Conference on Robot Learning*, pp. 991–1002. PMLR, 2022.
- Janner, M., Du, Y., Tenenbaum, J. B., and Levine, S. Planning with diffusion for flexible behavior synthesis, 2022. URL <https://arxiv.org/abs/2205.09991>.
- Jiang, Y., Gu, S., Murphy, K., and Finn, C. *Language as an Abstraction for Hierarchical Deep Reinforcement Learning*. Curran Associates Inc., Red Hook, NY, USA, 2019.
- Jiang, Y., Gupta, A., Zhang, Z., Wang, G., Dou, Y., Chen, Y., Fei-Fei, L., Anandkumar, A., Zhu, Y., and Fan, L. Vima: General robot manipulation with multimodal prompts. *arXiv preprint arXiv:2210.03094*, 2022.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Levine, S. Reinforcement learning and control as probabilistic inference: Tutorial and review. *arXiv preprint arXiv:1805.00909*, 2018.
- Li, S., Puig, X., Du, Y., Wang, C., Akyurek, E., Torralba, A., Andreas, J., and Mordatch, I. Pre-trained language models for interactive decision-making. *arXiv preprint arXiv:2202.01771*, 2022a.
- Li, X. L., Thickstun, J., Gulrajani, I., Liang, P., and Hashimoto, T. B. Diffusion-lm improves controllable text generation, 2022b. URL <https://arxiv.org/abs/2205.14217>.
- Liang, J., Huang, W., Xia, F., Xu, P., Hausman, K., Ichter, B., Florence, P., and Zeng, A. Code as policies: Language model programs for embodied control. *arXiv preprint arXiv:2209.07753*, 2022.
- Luketina, J., Nardelli, N., Farquhar, G., Foerster, J., Andreas, J., Grefenstette, E., Whiteson, S., and Rocktäschel, T. A survey of reinforcement learning informed by natural language. *arXiv preprint arXiv:1906.03926*, 2019.
- Lynch, C. and Sermanet, P. Grounding language in play. *CoRR*, abs/2005.07648, 2020a. URL <https://arxiv.org/abs/2005.07648>.
- Lynch, C. and Sermanet, P. Language conditioned imitation learning over unstructured data. *arXiv preprint arXiv:2005.07648*, 2020b.
- Mees, O., Hermann, L., and Burgard, W. What matters in language conditioned robotic imitation learning over unstructured data, 2022a. URL <https://arxiv.org/abs/2204.06252>.
- Mees, O., Hermann, L., Rosete-Beas, E., and Burgard, W. Calvin: A benchmark for language-conditioned policy learning for long-horizon robot manipulation tasks. *IEEE Robotics and Automation Letters*, 2022b.
- Mu, J., Zhong, V., Raileanu, R., Jiang, M., Goodman, N., Rocktäschel, T., and Grefenstette, E. Improving intrinsic exploration with language abstractions. *arXiv preprint arXiv:2202.08938*, 2022.
- Nachum, O., Gu, S. S., Lee, H., and Levine, S. Data-efficient hierarchical reinforcement learning. *Advances in neural information processing systems*, 31, 2018.
- Nachum, O., Gu, S., Lee, H., and Levine, S. Near-optimal representation learning for hierarchical reinforcement learning. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=Hlemus0qF7>.
- Nichol, A. Q., Dhariwal, P., Ramesh, A., Shyam, P., Mishkin, P., McGrew, B., Sutskever, I., and Chen, M. GLIDE: Towards photorealistic image generation and editing with text-guided diffusion models. In Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., and Sabato, S. (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 16784–16804. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/nichol22a.html>.

- Perez, E., Strub, F., De Vries, H., Dumoulin, V., and Courville, A. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Quine, W. V. O. Word and object mit press. *Cambridge MA*, 1960.
- Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., and Sutskever, I. Language models are unsupervised multitask learners. 2019.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., and Liu, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., and Chen, M. Hierarchical text-conditional image generation with clip latents, 2022. URL <https://arxiv.org/abs/2204.06125>.
- Rasul, K., Seward, C., Schuster, I., and Vollgraf, R. Autoregressive denoising diffusion models for multivariate probabilistic time series forecasting. 2021. doi: 10.48550/ARXIV.2101.12072. URL <https://arxiv.org/abs/2101.12072>.
- Reed, S., Zolna, K., Parisotto, E., Colmenarejo, S. G., Novikov, A., Barth-Maron, G., Gimenez, M., Sulsky, Y., Kay, J., Springenberg, J. T., et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- Ross, S., Gordon, G., and Bagnell, D. 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*, pp. 627–635. JMLR Workshop and Conference Proceedings, 2011.
- Schaul, T., Horgan, D., Gregor, K., and Silver, D. Universal value function approximators. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1312–1320, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/schaul15.html>.
- Singer, U., Polyak, A., Hayes, T., Yin, X., An, J., Zhang, S., Hu, Q., Yang, H., Ashual, O., Gafni, O., et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022.
- Sodhani, S., Zhang, A., and Pineau, J. Multi-task reinforcement learning with context-based representations, 2021. URL <https://arxiv.org/abs/2102.06177>.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2256–2265, Lille, France, 07–09 Jul 2015a. PMLR. URL <https://proceedings.mlr.press/v37/sohl-dickstein15.html>.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pp. 2256–2265. PMLR, 2015b.
- Song, J., Meng, C., and Ermon, S. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020a.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020b.
- Sutton, R. S., Precup, D., and Singh, S. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2): 181–211, 1999.
- Tan, C., Sun, F., Kong, T., Zhang, W., Yang, C., and Liu, C. A survey on deep transfer learning. In *Artificial Neural Networks and Machine Learning–ICANN 2018: 27th International Conference on Artificial Neural Networks, Rhodes, Greece, October 4-7, 2018, Proceedings, Part III* 27, pp. 270–279. Springer, 2018.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.
- Wang, T., Du, S. S., Torralba, A., Isola, P., Zhang, A., and Tian, Y. Denoised mdps: Learning world models better than the world itself. *arXiv preprint arXiv:2206.15477*, 2022a.

- Wang, X., Huang, Q., Celikyilmaz, A., Gao, J., Shen, D., Wang, Y.-F., Wang, W. Y., and Zhang, L. Vision-language navigation policy learning and adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 43 (12):4205–4216, 2020.
- Wang, Z., Hunt, J. J., and Zhou, M. Diffusion policies as an expressive policy class for offline reinforcement learning, 2022b. URL <https://arxiv.org/abs/2208.06193>.
- Welling, M. and Teh, Y. W. Bayesian learning via stochastic gradient langevin dynamics.
- Weng, L. What are diffusion models? *lilianweng.github.io*, Jul 2021. URL <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>.
- Zhang, A., McAllister, R., Calandra, R., Gal, Y., and Levine, S. Learning invariant representations for reinforcement learning without reconstruction. *arXiv preprint arXiv:2006.10742*, 2020.
- Zheng, K., Chen, X., Jenkins, O. C., and Wang, X. E. Vlmbench: A compositional benchmark for vision-and-language manipulation. *arXiv preprint arXiv:2206.08522*, 2022.
- Zhong, V., Rocktäschel, T., and Grefenstette, E. Rtfm: Generalising to novel environment dynamics via reading. *arXiv preprint arXiv:1910.08210*, 2019.
- Zhong, V., Mu, J., Zettlemoyer, L., Grefenstette, E., and Rocktäschel, T. Improving policy learning via language dynamics distillation. *arXiv preprint arXiv:2210.00066*, 2022.

Appendix

A. Visualizations of Qualitative Results

Please refer to our anonymous website <https://sites.google.com/view/541-lcd/home> for more qualitative results in video format. And we release our code implementation at anonymous github <https://anonymous.4open.science/r/LCD-ICML-9058>.

B. Hyper-parameter settings and training details

For all methods we proposed in Table 1, Table 2, Table 3, and Table 6, we obtain the mean and standard deviation of each method across 3 seeds. Each seed contains the individual training process and evaluates the policy for 100 episodes.

B.1. HP and training details for methods in Table 1 and Table 3.

Model	Module	Hyperparameter	Value
HULC	Trainer	Max Epochs	100
		β for KL Loss	0.01
		λ for Contrastive Loss	3
		Optimizer	Adam
		Learning Rate	2e-4
	Model	Transformer Hidden Size	2048
		Language Embedding Size	384
LCD	Gaussian Diffusion	Action Dimension	32
		Action Weight	10
		Loss Type	L2
		Observation Dimension	32
		Diffusion Steps	20
		Model Dimension	64
	Trainer	EMA Decay	0.995
		Label Frequency	200000
		Sample Frequency	1000
		Batch Size	512
		Learning Rate	2e-4
		Train Steps	250k
		Number of Steps Per Epoch	10000
		Normalizer	Gaussian Normalizer
		Frame Offset	0

Table 6: Hyperparameters for our methods in Table 1 and Table 3.

C. Training Objective Derivation

To model this, we turn to diffusion models (Weng, 2021), whom we borrow much of the following derivation from. Inspired by non-equilibrium thermodynamics, the common forms of diffusion models (Sohl-Dickstein et al., 2015b; Ho et al., 2020; Song et al., 2020a) propose modeling the data distribution $p(\boldsymbol{\tau})$ as a random process that steadily adds increasingly varied amounts of Gaussian noise to samples from $p(\boldsymbol{\tau})$ until the distribution converges to the standard normal. We denote the forward process as $f(\boldsymbol{\tau}_t|\boldsymbol{\tau}_{t-1})$, with a sequence of variances $(\beta_0, \beta_1 \dots \beta_T)$. We define $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$.

$$f(\boldsymbol{\tau}_{1:T}|\boldsymbol{\tau}_0) = \prod_{t=1}^T f(\boldsymbol{\tau}_t|\boldsymbol{\tau}_{t-1}), \quad \text{where } f(\boldsymbol{\tau}_t|\boldsymbol{\tau}_{t-1}) = \mathcal{N}(\boldsymbol{\tau}_t; \sqrt{1 - \beta_t}\boldsymbol{\tau}_{t-1}, \beta_t\mathbf{I}). \quad (6)$$

One can tractably reverse this process when conditioned on $\boldsymbol{\tau}_0$, which allows for the construction of a sum of the typical variational lower bounds for learning the backward process' density function (Sohl-Dickstein et al., 2015b). Since the backwards density also follows a Gaussian, it suffices to predict μ_θ and Σ_θ which parameterize the backwards distribution:

$$p_\theta(\boldsymbol{\tau}_{t-1} | \boldsymbol{\tau}_t) = \mathcal{N}(\boldsymbol{\tau}_{t-1}; \boldsymbol{\mu}_\theta(\boldsymbol{\tau}_t, t), \boldsymbol{\Sigma}_\theta(\boldsymbol{\tau}_t, t)). \quad (7)$$

In practice, Σ_θ is often fixed to constants, but can also be learned through reparameterization. Following (Ho et al., 2020) we consider learning only μ_θ , which can be computed just as a function of $\boldsymbol{\tau}_t$ and $\epsilon_\theta(\boldsymbol{\tau}_t, t)$. One can derive that $\boldsymbol{\tau}_t = \sqrt{\bar{\alpha}_t}\boldsymbol{\tau}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$ for $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, through a successive reparameterization of (6) until arriving at $f(\boldsymbol{\tau}_t|\boldsymbol{\tau}_0)$. Therefore to sample from $p(\boldsymbol{\tau})$, we need only to learn ϵ_θ , which is done by regressing to the ground truth ϵ given by the tractable backwards density. Assuming we have ϵ_θ , we can then follow a Markov chain of updates that eventually converges to the original data distribution, in a procedure reminiscent of Stochastic Gradient Langevin Dynamics (Welling & Teh):

$$\boldsymbol{\tau}_{t-1} = \frac{1}{\sqrt{1 - \beta_t}} \left(\boldsymbol{\tau}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\boldsymbol{\tau}_t, t) \right) + \sigma_t \mathbf{z}, \quad \text{where } \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (8)$$

To learn ϵ_θ , we can minimize the following variational lower bound on the negative log-likelihood:

$$\begin{aligned} L_{\text{CE}} &= -\mathbb{E}_{q(\boldsymbol{\tau}_0)} \log p_\theta(\boldsymbol{\tau}_0) \\ &= -\mathbb{E}_{q(\boldsymbol{\tau}_0)} \log \left(\int p_\theta(\boldsymbol{\tau}_{0:T}) d\boldsymbol{\tau}_{1:T} \right) \\ &= -\mathbb{E}_{q(\boldsymbol{\tau}_0)} \log \left(\int q(\boldsymbol{\tau}_{1:T}|\boldsymbol{\tau}_0) \frac{p_\theta(\boldsymbol{\tau}_{0:T})}{q(\boldsymbol{\tau}_{1:T}|\boldsymbol{\tau}_0)} d\boldsymbol{\tau}_{1:T} \right) \\ &= -\mathbb{E}_{q(\boldsymbol{\tau}_0)} \log \left(\mathbb{E}_{q(\boldsymbol{\tau}_{1:T}|\boldsymbol{\tau}_0)} \frac{p_\theta(\boldsymbol{\tau}_{0:T})}{q(\boldsymbol{\tau}_{1:T}|\boldsymbol{\tau}_0)} \right) \\ &\leq -\mathbb{E}_{q(\boldsymbol{\tau}_{0:T})} \log \frac{p_\theta(\boldsymbol{\tau}_{0:T})}{q(\boldsymbol{\tau}_{1:T}|\boldsymbol{\tau}_0)} \\ &= \mathbb{E}_{q(\boldsymbol{\tau}_{0:T})} \left[\log \frac{q(\boldsymbol{\tau}_{1:T}|\boldsymbol{\tau}_0)}{p_\theta(\boldsymbol{\tau}_{0:T})} \right] = L_{\text{VLB}}. \end{aligned} \quad (9)$$

$$L_{\text{VLB}} = L_T + L_{T-1} + \dots + L_0$$

where $L_T = D_{\text{KL}}(q(\boldsymbol{\tau}_T|\boldsymbol{\tau}_0) \parallel p_\theta(\boldsymbol{\tau}_T))$

$$L_t = D_{\text{KL}}(q(\boldsymbol{\tau}_t|\boldsymbol{\tau}_{t+1}, \boldsymbol{\tau}_0) \parallel p_\theta(\boldsymbol{\tau}_t|\boldsymbol{\tau}_{t+1}))$$

for $1 \leq t \leq T-1$ and

$$L_0 = -\log p_\theta(\boldsymbol{\tau}_0|\boldsymbol{\tau}_1).$$

Which enables us to find a tractable parameterization for training, as the KL between two Gaussians is analytically computable.

$$\begin{aligned}
 L_t &= \mathbb{E}_{\tau_0, \epsilon} \left[\frac{1}{2 \|\Sigma_\theta(\tau_t, t)\|_2^2} \|\tilde{\mu}_t(\tau_t, \tau_0) - \mu_\theta(\tau_t, t)\|^2 \right] \\
 &= \mathbb{E}_{\tau_0, \epsilon} \left[\frac{1}{2 \|\Sigma_\theta\|_2^2} \left\| \frac{1}{\sqrt{\alpha_t}} \left(\tau_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right) - \frac{1}{\sqrt{\alpha_t}} \left(\tau_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\tau_t, t) \right) \right\|^2 \right] \\
 &= \mathbb{E}_{\tau_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \bar{\alpha}_t) \|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(\tau_t, t)\|^2 \right] \\
 &= \mathbb{E}_{\tau_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \bar{\alpha}_t) \|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\alpha_t} \tau_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2 \right].
 \end{aligned} \tag{10}$$

After removing the coefficient at the beginning of this objective following (Ho et al., 2020), we arrive at the objective used in the practical algorithm 1:

$$\mathbb{E}_{\tau_0, \epsilon} [\|\epsilon_t - \epsilon_\theta(\sqrt{\alpha_t} \tau_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t)\|^2]. \tag{11}$$

Furthermore, thanks to the connection between noise conditioned score networks and diffusion models (Song et al., 2020b; Ho et al., 2020), we are able to state that $\epsilon_\theta \propto -\nabla \log p(\tau)$:

$$\begin{aligned}
 s_\theta(\tau_t, t) &\approx \nabla_{\tau_t} \log p(\tau_t) \\
 &= \mathbb{E}_{q(\tau_0)} [\nabla_{\tau_t} p(\tau_t | \tau_0)] \\
 &= \mathbb{E}_{q(\tau_0)} \left[- \frac{\epsilon_\theta(\tau_t, t)}{\sqrt{1 - \bar{\alpha}_t}} \right] \\
 &= - \frac{\epsilon_\theta(\tau_t, t)}{\sqrt{1 - \bar{\alpha}_t}}.
 \end{aligned} \tag{12}$$

Therefore, by using a variant of ϵ_θ conditioned on language to denoise our latent plans, we can effectively model $-\nabla_{\tau} \mathcal{P}_\beta(\tau | \mathcal{R})$ with our diffusion model, iteratively guiding our generated trajectory towards the optimal trajectories conditioned on language.

D. Proof of 3.1

Proof. The proof is fairly straightforward, and can be shown by translating our definition of suboptimality into the framework utilized by (Nachum et al., 2019). We are then able to leverage their first theorem bounding suboptimality by the Total Variation (TV) between transition distributions to show our result, as TV is bounded by the Lipschitz constant multiplied by the domain of the function.

(Nachum et al., 2019) first define a low level policy generator Ψ which maps from $S \times \tilde{A}$ to Π . Using the high level policy to sample a goal $g_t \sim \pi_{hi}(g|s_t)$, they use Ψ to translate this to a policy $\pi_t = \Psi(s_t, g_t)$, which samples actions $a_{t+k} \sim \pi_t(a|s_{t+k}, k)$ from $k \in [0, c - 1]$. The process is repeated from s_{t+c} . Furthermore, they define an inverse goal generator $\varphi(s, a)$, which infers the goal g that would cause Ψ to yield an action $\tilde{a} = \Psi(s, g)$. The following can then be shown:

Theorem D.1. *If there exists $\varphi : S \times A \rightarrow \tilde{A}$ such that,*

$$\sup_{s \in S, a \in A} D_{TV}(P(s'|s, a) || P(s'|s, \Psi(s, \varphi(s, a)))) \leq \epsilon, \quad (13)$$

then $\text{SubOpt}'(\Psi) \leq C\epsilon$, where $C = \frac{2\gamma}{(1-\gamma)^2} R_{max}$.

Note that their SubOpt' is different from ours; whilst we defined in terms of the encoder \mathcal{E} and action generator ϕ , they define it in terms of Ψ . Note, however, that the two are equivalent when the temporal stride $c = 1$, as Ψ becomes $\pi_{lo} = \phi \circ \mathcal{E}$. It is essential to note that when using a goal conditioned imitation learning objective, as we do in this paper, π_{lo} becomes equivalent to an inverse dynamics model $\text{IDM}(s, \mathcal{E}(s)) = a$ and that $\varphi(s, a)$ becomes equivalent to $\mathcal{E}(s')$. This is the key to our proof, as the second term in the total variation of D.1 reduces to

$$\begin{aligned} & P(s'|s, \Psi(s, \varphi(s, a))) \\ &= P(s'|s, \Psi(s, \mathcal{E}(s'))) \\ &= P(s'|s, a + \epsilon). \end{aligned} \quad (14)$$

Since we have that the transition dynamics are Lipschitz:

$$\begin{aligned} & \int_{\mathcal{A}, \mathcal{S}} |P(s'|s, a) - P(s'|s, a + \epsilon)| d\nu \\ & \leq \int_{\mathcal{A}, \mathcal{S}} K_f |a - (a + \epsilon)| d\nu \\ &= K_f \epsilon \int_{\mathcal{A}, \mathcal{S}} d\nu \\ &= K_f \epsilon \text{dom}(P(s'|s, a)) \end{aligned} \quad (15)$$

Which we can then plug into 13 to obtain the desired $C = \frac{2\gamma}{(1-\gamma)^2} R_{max} K_f \text{dom}(P(s'|s, a))$. \square

E. Inference Pipeline

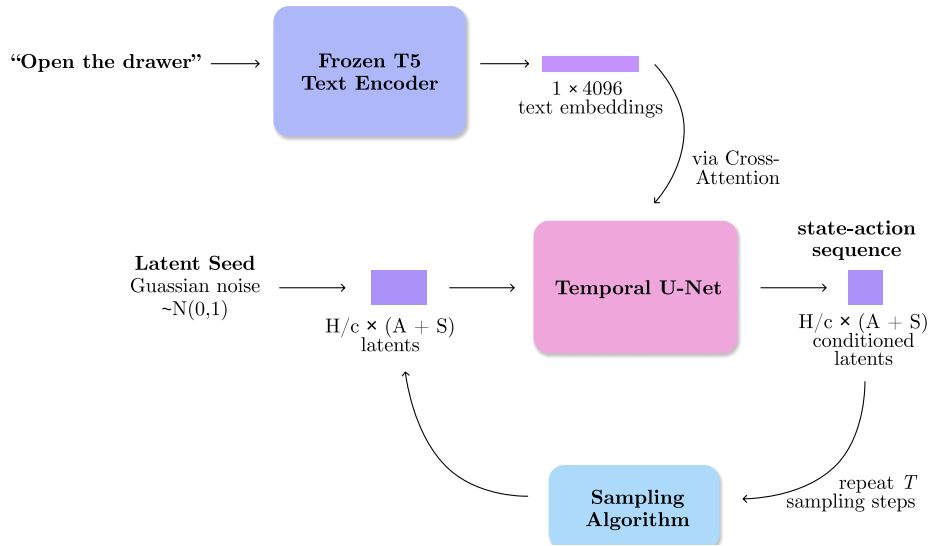


Figure 4: An overview of our inference pipeline. We pass sampled noise to our denoising autoencoder along with an encoded language prompt. The autoencoder is instantiated as a temporal U-Net. By repeating this process iteratively, we are able to generate high-fidelity latent plans conditioned on language. We are also able to decode the latent states into pixel space to analyze and interpret the plans generated by the denoising autoencoder.

F. Diffuser-2D Diffusion

In Figure 5 and Figure 6, we give an example of the denoising process of one of our ablations, the Diffuser-2D model. This model utilizes the 2D autoencoder of (Rombach et al., 2022) with (Janner et al., 2022).

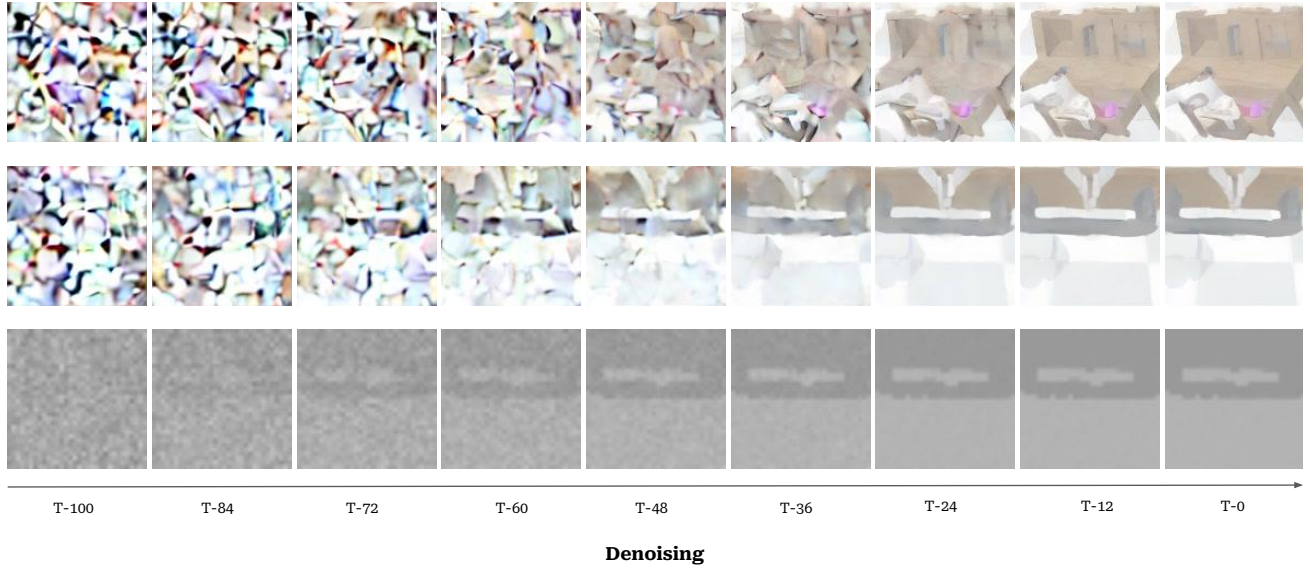


Figure 5: An overview of our Denoising process.

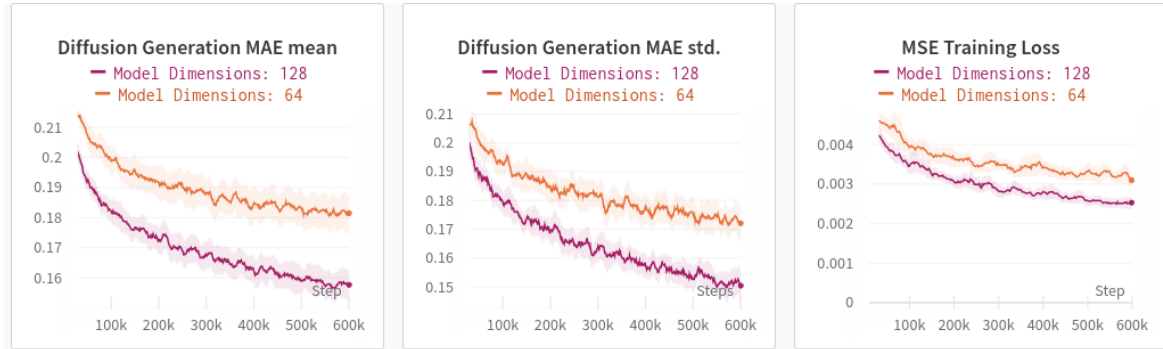


Figure 6: Diffusion Generation.

G. TSNE Comparison

We visualize the TSNE of GT-VAE trajectory and DM trajectory in Figure 7.

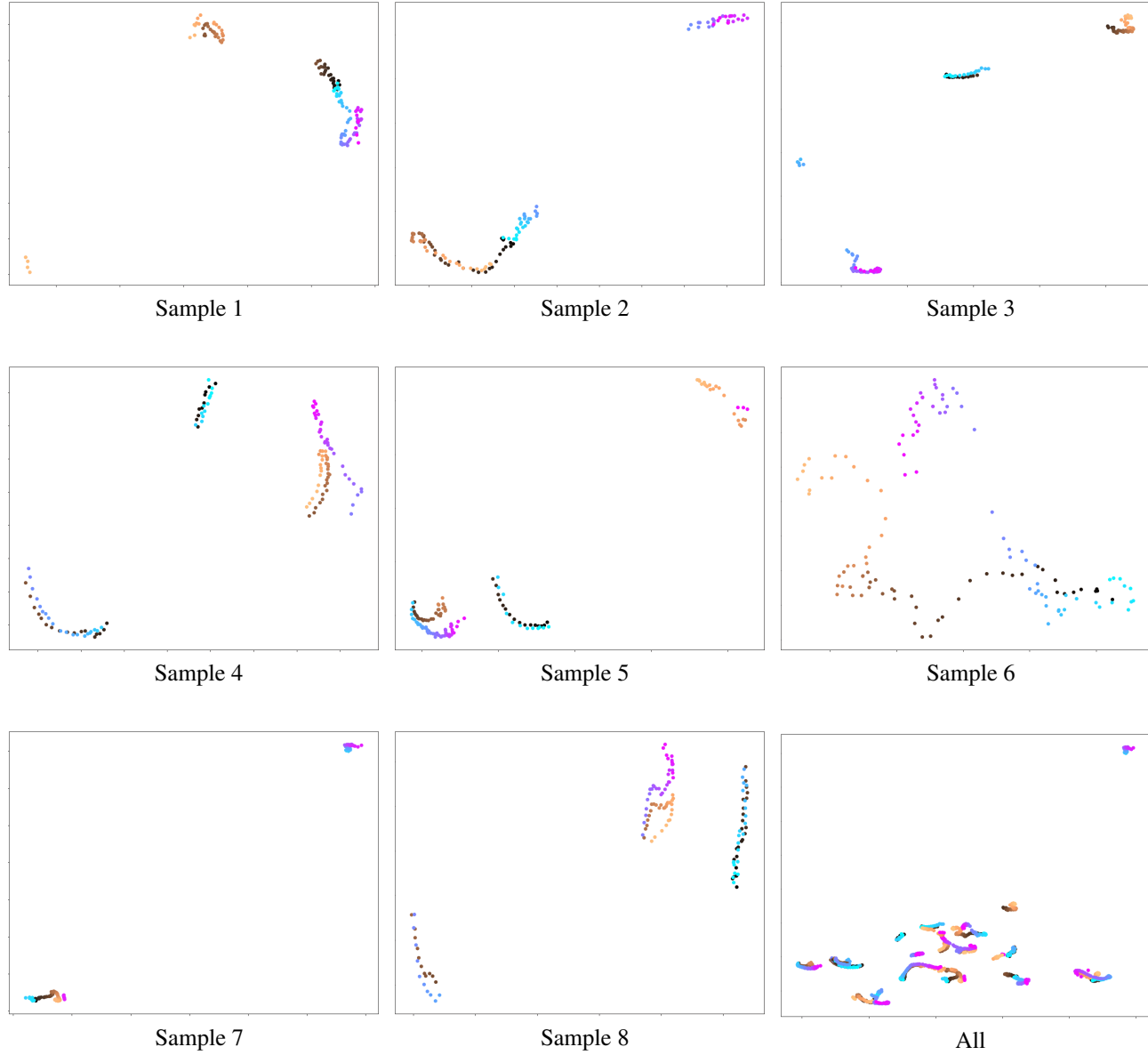


Figure 7: TSNE visualization of GT-VAE trajectory vs. Diffuser-1D trajectory, where the purple and light blue color range is the ground truth VAE, and the copper color range is Diffuser-1D. All states are normalized, and all trajectories are taken from the task “lift pink block table”.

H. Latent Plan TSNE

We give a TSNE embedding of the latent plan space of HULC in Figure 8. The latent plan space is the communication layer between the high level policy and low level policy of the HULC model, which corresponds to the intermediate layer between the lower level and lowest level policy in our method.

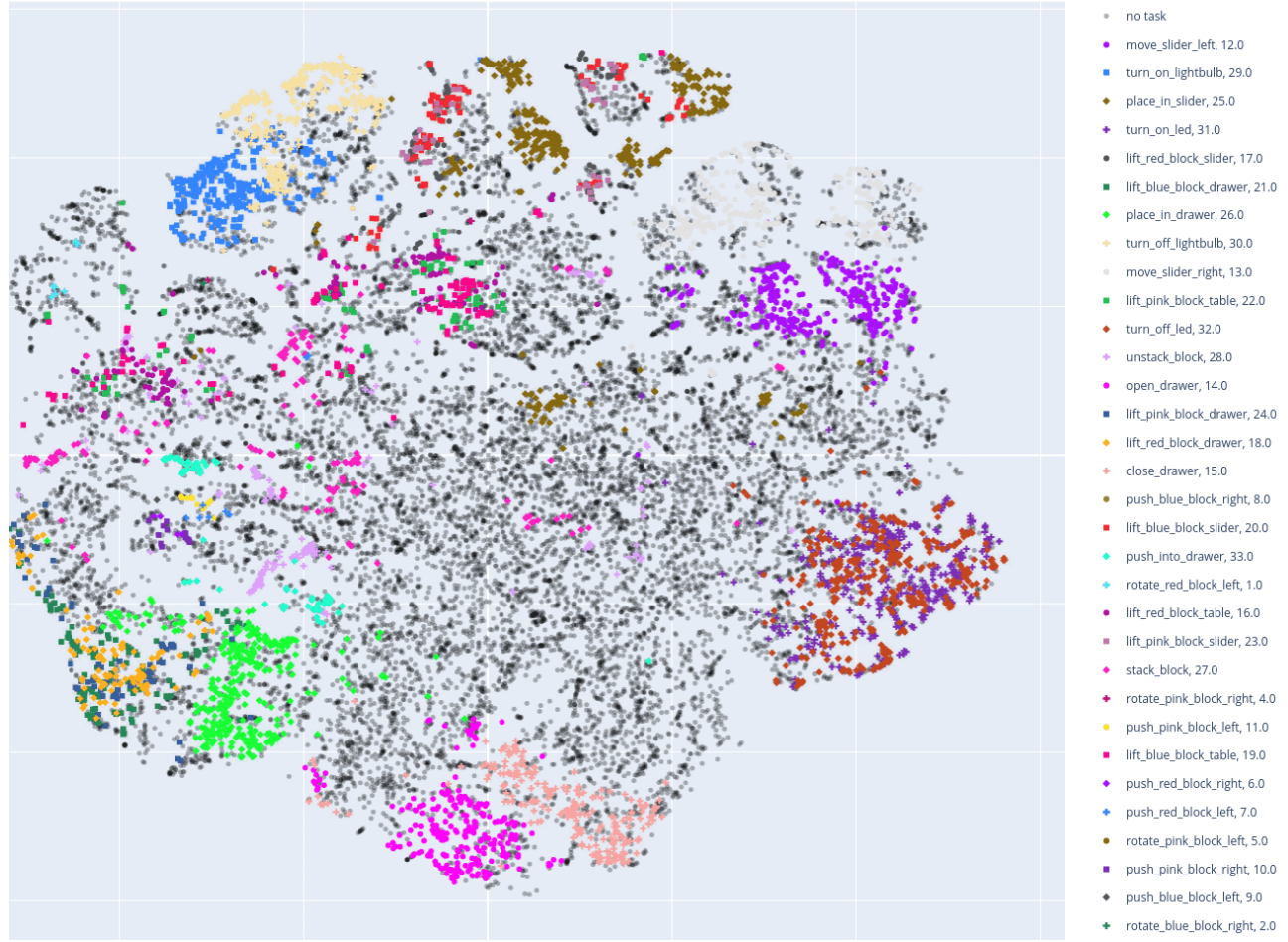


Figure 8: TSNE of Latent Plan.

I. Task Distribution

We visualize the distribution of all the tasks considered in our experiments in Figure 9. Note the long-tailedness of this distribution, and how it skews evaluation scores upwards if one can solve the relatively easier tasks that occur most frequently.

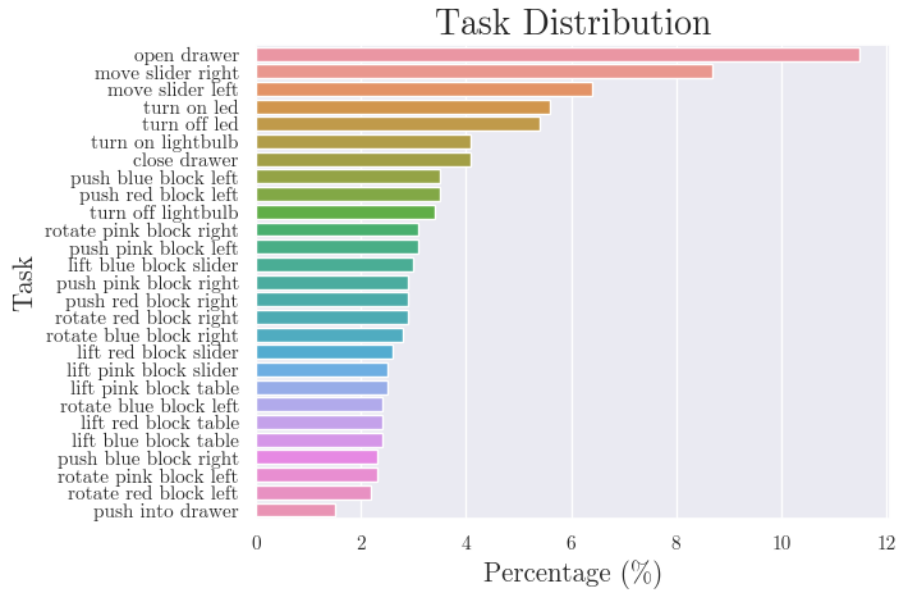


Figure 9: The Evaluation Task Distribution.