

Diffusion Model for Control and Planning Tutorial

Chaoyi Pan 2024.2.26

Tutorial outline:

1. Recap: What is a Diffusion Model / What Problem Does It Solve?
2. Motivation: Why Do We Need a Diffuser in Control and Planning?
3. Practice: How to Use a Diffuser in Control and Planning?
4. Literatures: Recent Research in a Diffuser for Control and Planning
5. Summary & Limitations: What We Can Do and What We Cannot Do

Recap: Diffusion Model

A diffusion model is a generative model capable of generating samples from a given distribution. It serves as a powerful tool for distribution matching, extensively utilized in image generation, text generation, and other creative tasks.



Figure 1: diffusion examples

At the core of the diffusion model is the score function, representing the noise direction used to update samples to match the target distribution. Through learning this score function, the diffusion model can adeptly generate samples from the specified distribution.

$$x_{i+1} \leftarrow x_i + c\nabla \log p(x_i) + \sqrt{2c\epsilon}, \quad i = 0, 1, \dots, K$$

The diffusion model have several advantages over other generative models, including but not limited to the following:

- **Multimodal:** It effectively handles multimodal distributions, a challenge often encountered when directly predicting distributions.
- **Scalable:** This approach scales well with high-dimensional distribution matching problems, making it versatile for various applications.
- **Stable:** Grounded in solid math and a standard multi-stage diffusion training procedure, the model ensures stability during training.

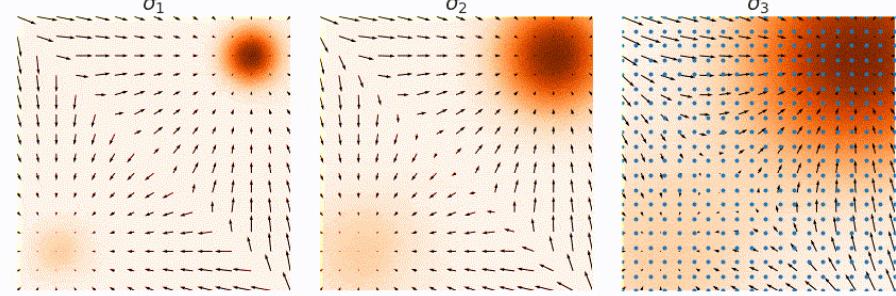


Figure 2: Annealed Langevin dynamics combine a sequence of Langevin chains with gradually decreasing noise scales.

- **Non-autoregressive:** Its capability to predict entire trajectory sequences in a non-autoregressive manner enables efficient handling of non-autoregressive and multimodal distribution matching challenges.

Motivation: Why Do We Need a Diffuser in Control and Planning?

Generative Models in Control and Planning

Before diffusion models became popular, there were other generative models used in control and planning. Below are a few examples of how other generative models have been applied in imitation learning:

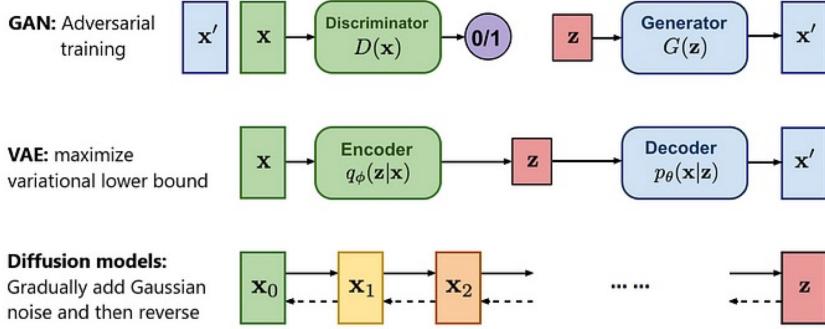
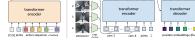


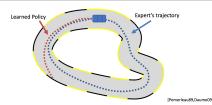
Figure 3: Other generative model overview

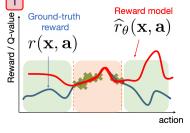
Generative model	Application	Idea	Limitation
GAN	Generative Adversarial Imitation Learning (GAIL): Learning a discriminator and training a policy to fool the discriminator	<pre> Algorithm 1 Generative adversarial imitation learning 1: Input: Expert trajectories $\tau_1 \sim \pi_\text{exp}$, initial policy and discriminator parameters θ_1, θ_2 2: $t = 1$ 3: Sample trajectories $\tau_t \sim \pi_\theta$ 4: Update discriminator θ_2 to minimize loss from π_θ w.r.t. π_{exp}, with the gradient $\mathbb{E}_{\tau}[\nabla_{\theta_2} \log(D_{\theta_2}(s, a)) + \mathbb{E}_{\tau}[\nabla_a \log(1 - D_{\theta_2}(s, a))]] \quad (D)$ 5: Take policy step from θ_1 to θ_{1+1}, using the TRPO rule with cost function $\log(J_{\text{imitate}}(\pi_\theta))$. Note: $J_{\text{imitate}}(\pi_\theta) = \mathbb{E}_{\tau}[\nabla_a \log(Q(s, a)) - \nabla_a \log(\pi_\theta(s, a))]$ when $Q(s, a) = \mathbb{E}_\pi[\log(Q_{\text{expert}}(s, a)) s = s, a = a]$ 6: <i>end for</i> </pre>	Difficult to handle multimodal distributions and unstable training
VAE	Action Chunking with Transformers (ACT) (ALOHA): Learning a latent space (encoder: expert action sequence + observation -> latent) and using the latent space for planning (decoder: latent + more observation -> action sequence prediction)		Challenging to train

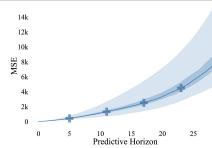
In these scenarios, the objective is to match the distribution of the dataset, similar to the goal of diffusion models. Compared with other generative models like GANs and VAEs, diffusion models are better at handling multimodal data and offer more stable training processes.

What to Learn with the Generative Model?

From a planning and reinforcement learning perspective, there are numerous scenarios where matching the dataset's distribution is crucial, such as:

Scenario	Challenge	Solution	Illustrations
Imitation Learning	$\min \ \ p_{\theta}(\tau) - p_{\text{data}}(\tau) \ \ $ Match the demonstrations' trajectory distribution (high-dimensional+multi-modality) with limited data. Common methods like GAIL use adversarial training to match the distribution.	Diffusion models excel at matching the distribution of the expert's actions with high capacity and expressiveness.	 <small>(From NeurIPS 2019)</small>

Scenario	Challenge	Solution	Illustrations
Offline Reinforcement Learning	$\max J_\theta(\tau) \geq J_{\text{data}}(\tau)$ s.t. $\ p_\theta(\tau) - p_{\text{data}}(\tau)\ < \epsilon$ Perform better than demonstrations with a large number of demonstrations. Here, it's essential to ensure the policy's action distribution is close to the dataset while improving performance. Common methods like CQL penalize out-of-distribution (OOD) samples, making the method overly conservative.	Diffusion models can match the dataset's action and regularize the policy's action distribution effectively.	

Scenario	Challenge	Solution	Illustrations
Model-based Reinforcement Learning	Match the dynamic model and (sometimes) the policy's action distribution. This involves first learning the model and then using it to plan in an autoregressive manner. This method suffers from compounding errors.	Diffusion models are adept at handling non-autoregressive and multimodal distribution matching by predicting the entire trajectory sequence at once.	

Diffusion models can be utilized to learn the policy, the planner, or the model itself. Each of these applications can be viewed as a distribution matching problem. The next section will delve into the choices and implications of using diffusion models in these contexts.

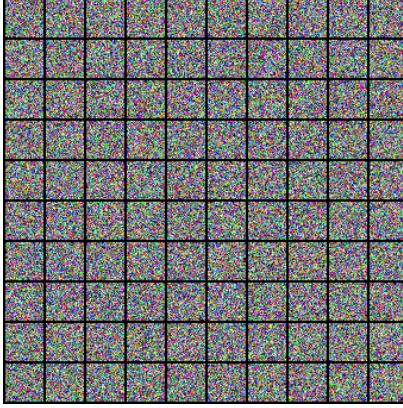
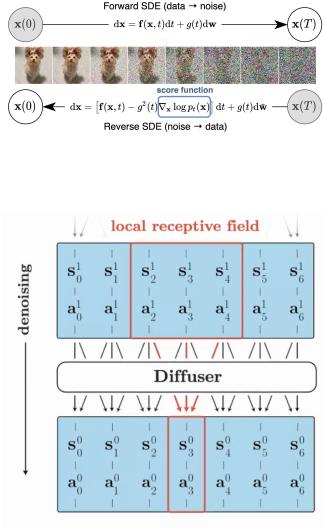
Practice: How to Use the Diffuser?

What to Diffuse?

In practice, there are several ways to incorporate the diffusion model into control and planning. The most common method is known as the ‘diffuser’:

The diffuser operates by concatenating the state and action, allowing us to diffuse the state-action sequence. This process is akin to diffusing a single-channel image. The training of the model follows a similar approach to image generation. Initially, noise is added to the state-action sequence. Subsequently, the model is trained to predict the score function or noise vector.

Here a local field method is implemented with a temporal convolutional network (TCN) to impose local consistency on the state-action sequence.

Task	Thing's to Diffuse	How to Diffuse
Image Generation		
Planning		

How to Impose Constraints/Objectives?

However, when training a model with given trajectory data, this model can only replicate the same actions as in the data, which is not what we want. We want the model to generalize to new tasks and constraints. To achieve this, we need to make the model conditional on the task and constraints. There are a few ways to do this:

Guidance Function

The guidance function directly shifts the distribution/cost or learned value, etc., or trains a discriminator (classifier) to obtain the guidance function. There are two common ways to get the guidance function:

1. Predefined the guidance function: This approach is easy to implement but might lead to out-of-distribution (OOD) samples, breaking the learned distribution.

$$\tilde{p}_\theta(\tau) \propto p_\theta(\tau) h(\tau) \tau^{i-1} = \mathcal{N}(\mu + \alpha \Sigma \nabla \mathcal{J}(\mu), \Sigma^i)$$

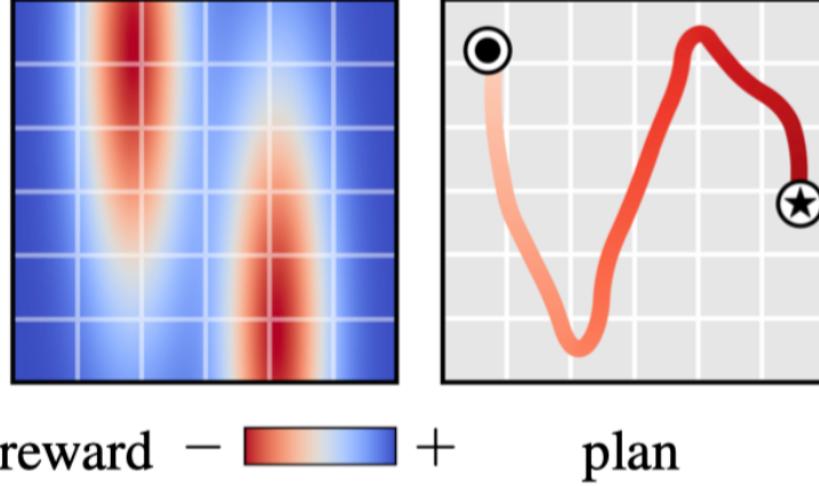


Figure 4: defined guidance function

1. Learned guidance function: This approach is more common in training image generation models, which involves training a classifier to obtain the guidance function. However, this method requires training a classifier in an adversarial manner, which might lead to unstable training.

$$\begin{aligned}\nabla \log p(x_t | y) &= \nabla \log \left(\frac{p(x_t) p(y | x_t)}{p(y)} \right) \\ &= \nabla \log p(x_t) + \nabla \log p(y | x_t) - \nabla \log p(y) \\ &= \underbrace{\nabla \log p(x_t)}_{\text{unconditional score}} + \underbrace{\nabla \log p(y | x_t)}_{\text{adversarial gradient}}\end{aligned}$$

Classifier-Free Method

By performing the following transformation, we can obtain the guidance function without training a classifier. The two terms below are known as the unconditional score and the conditional score, respectively. In practice, we can omit the input of the conditional score to achieve the unconditional score.

$$\begin{aligned}\nabla \log p(x_t | y) &= \nabla \log p(x_t) + \gamma (\nabla \log p(x_t | y) - \nabla \log p(x_t)) \\ &= \nabla \log p(x_t) + \gamma \nabla \log p(x_t | y) - \gamma \nabla \log p(x_t) \\ &= \underbrace{\gamma \nabla \log p(x_t | y)}_{\text{conditional score}} + \underbrace{(1 - \gamma) \nabla \log p(x_t)}_{\text{unconditional score}}\end{aligned}$$

Guidance Function Method	Classifier-Free Method
Algorithm 1 Guided Diffusion Planning <pre> 1: Require Diffuser μ_θ, guide \mathcal{J}, scale α, covariances Σ^i 2: while not done do 3: Observe state s; initialize plan $\tau^N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 4: for $i = N, \dots, 1$ do 5: // parameters of reverse transition 6: $\mu \leftarrow \mu_\theta(\tau^i)$ 7: // guide using gradients of return 8: $\tau^{i-1} \sim \mathcal{N}(\mu + \alpha \Sigma \nabla \mathcal{J}(\mu), \Sigma^i)$ 9: // constrain first state of plan 10: $\tau_{s_0}^{i-1} \leftarrow s$ 11: Execute first action of plan $\tau_{\mathbf{a}_0}^0$ </pre>	Algorithm 1 Conditional Planning with the Decision Diffuser <pre> 1: Initialize Noise model ϵ_θ, inverse dynamics f_θ, guidance scale ω, history length C, condition y 2: Initialize $d \leftarrow \text{Queue}(\text{length} = C)$, $t \leftarrow 0$ // Maintain a history of length C 3: while not done do 4: Observe state s; $h.\text{insert}(s)$; Initialize $\mathbf{x}_K(\tau) \sim \mathcal{N}(0, \alpha I)$ 5: for $k = K \dots 1$ do 6: $\mathbf{x}_k(\tau) \leftarrow \mathbf{x}_k(\tau)[k] = \mathbf{x}_k(\tau[k]) - \epsilon_\theta(\mathbf{x}_k(\tau), k)$ // Constrain plan to be consistent with history 7: $i \leftarrow \text{argmax}_{i \in \{0, \dots, K\}} \omega_i \epsilon_\theta(\mathbf{x}_k(\tau), y, k) - \epsilon_\theta(\mathbf{x}_k(\tau), i)$ // Classifier-free guidance 8: $(\mu_{k-1}, \Sigma_{k-1}) \leftarrow \text{Denoise}(\mathbf{x}_k(\tau), i)$ 9: $\mathbf{x}_{k-1} \sim \mathcal{N}(\mu_{k-1}, \alpha \Sigma_{k-1})$ 10: end for 11: Estimate (s_t, a_{t+1}) from $x_0(t)$ 12: Execute $a_t = f_\theta(s_t, s_{t+1})$; $t \leftarrow t + 1$ 13: end while </pre>

Inpainting

If the control problem involves specific state constraints (such as the initial/target state or constraints), we can simply fix the state and fill in the missing parts of the distribution. This approach is extremely useful in goal-reaching and navigation tasks.

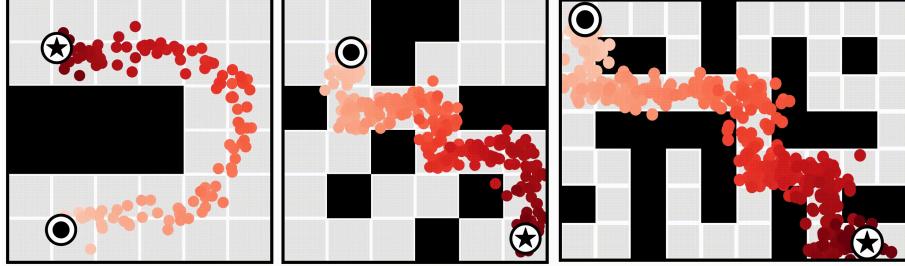


Figure 5: inpainting

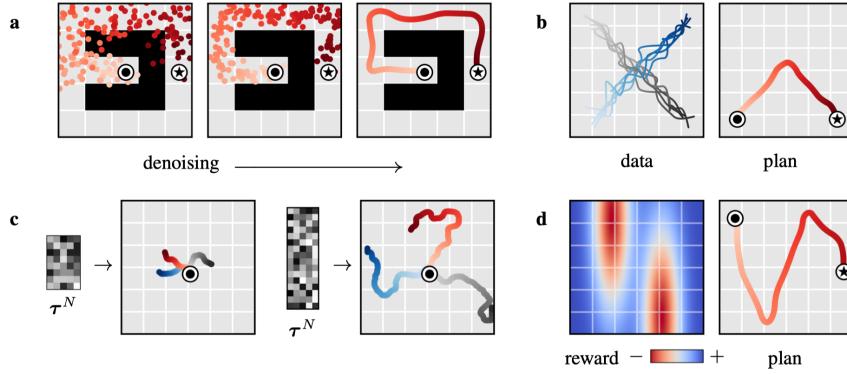


Figure 6: different ways to impose constraints/objectives

Literatures: Research Progress in Diffusion Models for Control and Planning

A detailed summary of each method can be found here.

$$\underbrace{\nabla_x \log P}_{\text{how to get score function}} \quad (\quad \underbrace{x}_{\text{what to diffuse}} \mid \underbrace{y}_{\text{how to impose constraints/objectives}} \quad)$$

The heart of the diffusion model is understanding how to obtain the score function. Based on the methods for obtaining the score function, what to diffuse, and how to impose constraints/objectives, we can categorize the recent research in diffusion models for control and planning into the following three axes:

Axis 1: How to Get the Score Function

- Data-driven: Learning from data by manually adding noise
- Hybrid: Learning from another optimization process
- Model-based: Calculating analytically from the model

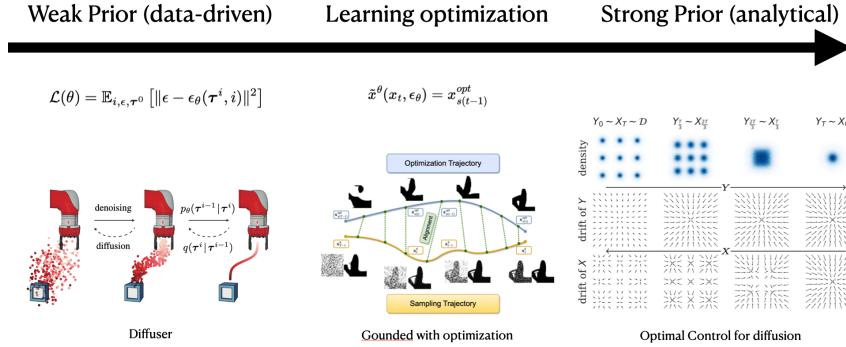


Figure 7: Different ways to get the score function

The data-driven method is the most common approach to obtaining the score function, which involves adding noise to the data and then training the model to predict the score function. The hybrid method learns the score function from the intermediate results of another optimization process, typically used in specific optimization problems. Finally, if you can calculate the score function analytically, then you can use Langevin dynamics to estimate the final distribution.

Axis 2: What to Diffuse

- Action: Directly diffuse for the next action
- State: Learn the model
- State-sequence: Diffuse for the next state sequence, or sometimes state-action sequence, or action sequence (for position control)

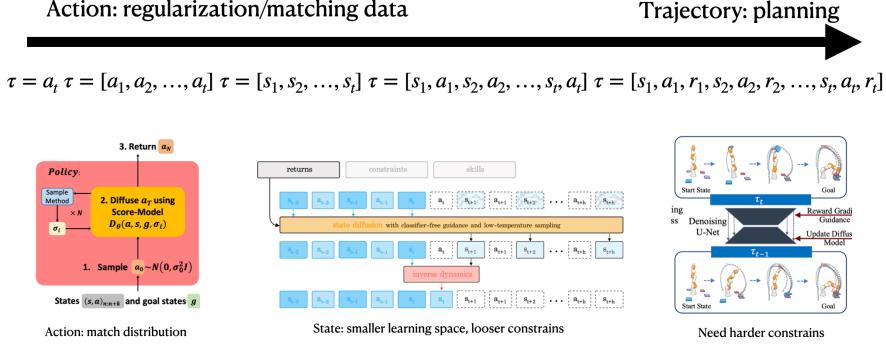


Figure 8: Different things to diffuse

What to diffuse depends on the goals. For instance, to regularize the action distribution to the dataset, diffusing the action is suitable. To learn a dynamically feasible optimal trajectory, diffusing the state-sequence is appropriate. Sometimes, diffusing the action sequence makes execution easier, while other times inverse dynamics are needed to obtain the action sequence.

Axis 3: How to Impose Constraints/Objectives

- Guidance function: Predefined or learned
- Classifier-free: Use the unconditional score and conditional score
- Inpainting: Fix the state and fill in the missing parts of the distribution

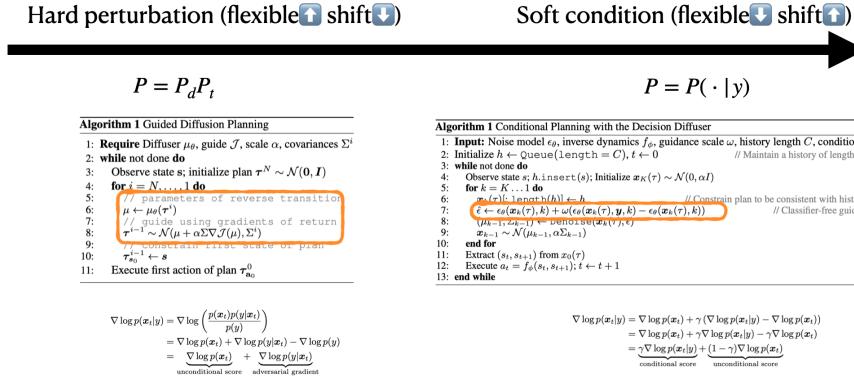


Figure 9: Different ways to impose constraints/objectives

The guidance function is the easiest way to impose constraints/objectives, which involves multiplying the guidance function with the model's distribution. The classifier-free method uses the unconditional and conditional scores to impose constraints/objectives, as seen in methods like **decision diffuser**, **adaptive**

diffuser, etc. Inpainting involves fixing the state and filling in the missing parts of the distribution, useful in goal-reaching or navigation tasks. This approach is complementary to the guidance function and classifier-free method.

One interesting work in this line called **safe diffuser** is to solve the quadratic programming (QP) problem at each step to add hard constraints. This method is useful in safety-critical tasks, where they prove that the diffusion model can satisfy the constraints if it converges.

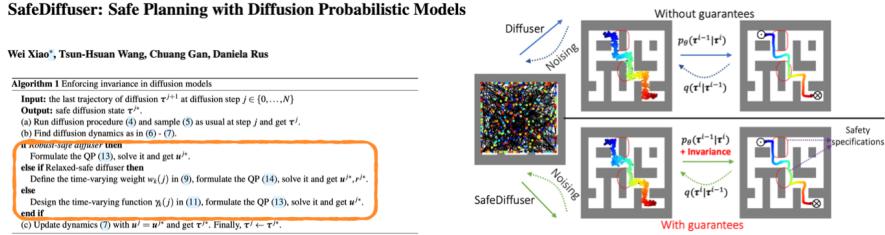


Figure 10: safe diffuser solving QP each step to add hard constraints

Summary & Limitations: What are the Challenges?

From the above discussion, we can see that the diffusion model can be used in control and planning to match the distribution of the dataset, which is widely used in imitation learning, offline reinforcement learning, and model-based reinforcement learning. The diffusion model can be used to learn the policy, the planner, or the model, which can also be viewed as a distribution matching problem.

Why Does Diffusion Work?

Compared with learning the explicit policy directly or learning the energy-based model, the diffusion model can handle multimodal distribution and higher-dimensional distribution matching by iteratively predicting the score function. This greatly smooths the distribution matching process and makes the training more stable and scalable.

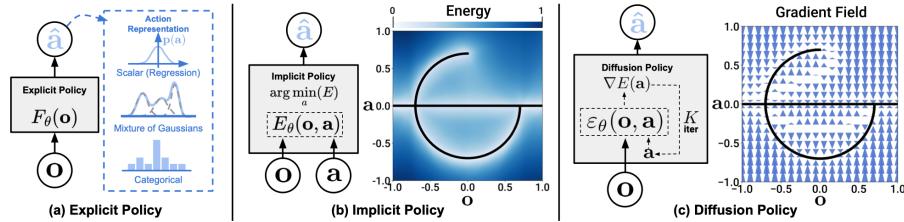


Figure 11: Compare distribution with other models

Limitations

However, the diffusion model also has its limitations, which include:

1. Computational cost: The diffusion model requires a longer time to train (a few GPU days compared with tens of minutes) and inference (iterative sample steps compared with one forward pass). This makes high-frequency control and planning difficult to use with the diffusion model.
2. Handling shifting distribution: In online RL, the distribution of the policy will keep changing. Adapting the diffusion model to the new distribution requires a large amount of data and a long time to train. This limits the diffusion model to be trained in a fixed rather than a dynamic dataset.
3. High variance: Depending on the initial guess and random sampling, the variance of the diffusion model is high, which limits its application in high-precision or safety-critical tasks.
4. Constraint satisfaction: The diffusion model does not guarantee to satisfy the constraints, especially when tested in a constraint different from the training set. This limits its application in adapting to new constraints and tasks.