

IS CONDITIONAL GENERATIVE MODELING ALL YOU NEED FOR DECISION MAKING?

a.k.a DECISION DIFFUSER

Tru

01 논문 요약

Abstract

- **Conditional Generative Modeling**을 **순차적 의사 결정 문제**에 적용해보았다.
- Generative Modeling이 표준 벤치마크에서 기존 오프라인 RL 접근 방식을 능가할 수 있는 Policy도 만들 수 있었다.
- Policy를 Conditional Diffusion model로 모델링함으로써 **Dynamic Programming**을 사용하지 않을 수 있었고 DP의 Side effect도 없앨 수 있었다. (contribution)
- Constraints와 Skill이라는 두 가지 다른 조건 변수를 추가하여 Conditional Diffusion model로서 정책 모델링의 이점을 추가로 보여준다. (Contribution 2)

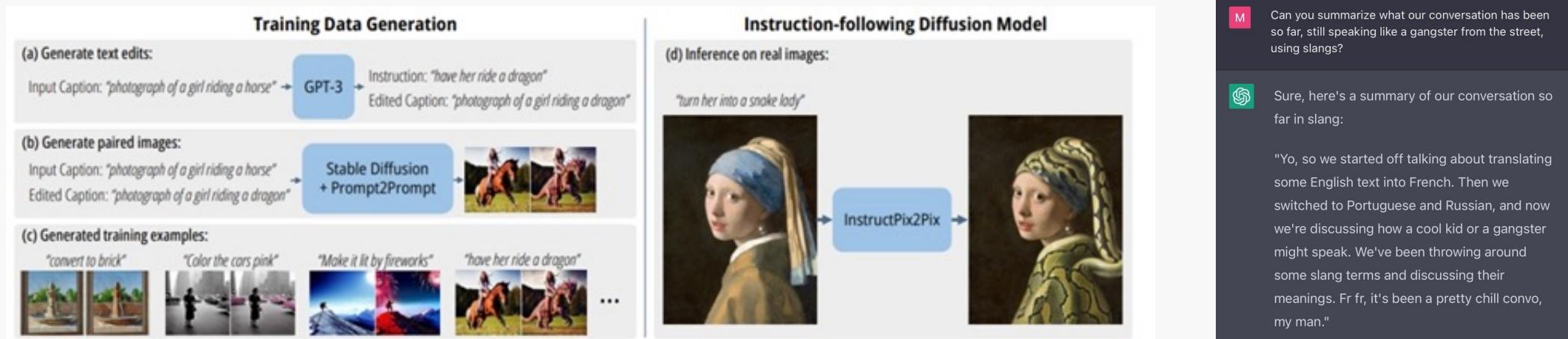
서론

ALL YOU NEED IS DIFFUSER IN REINFORCEMENT LEARNNING?

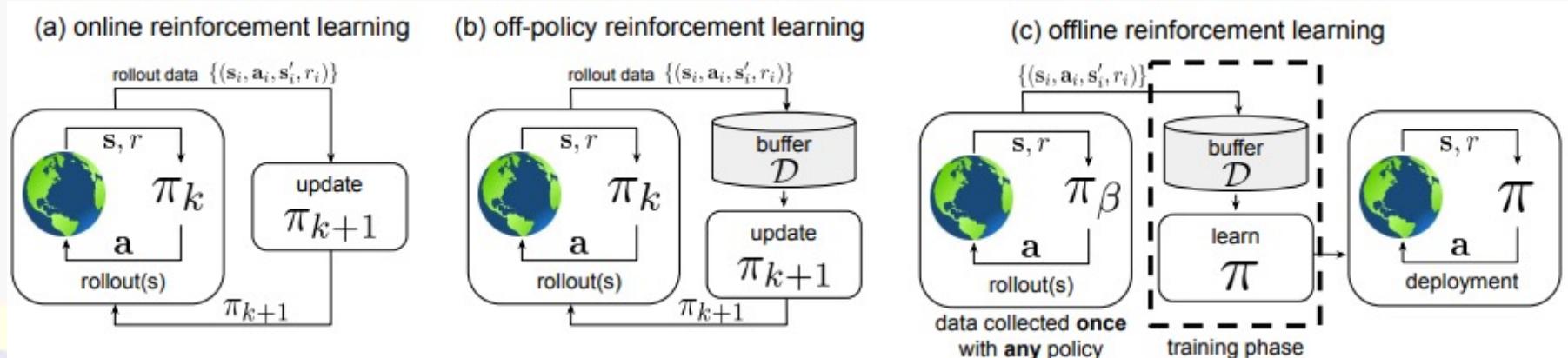
01 Decision Diffuser - 서론

서론

Over the last few years, there have been incredible advances in the field of **conditional generative modeling**



there exists a wide body of work that has looked into recovering high-performing policies from data logged by already operational systems (**offline RL**)

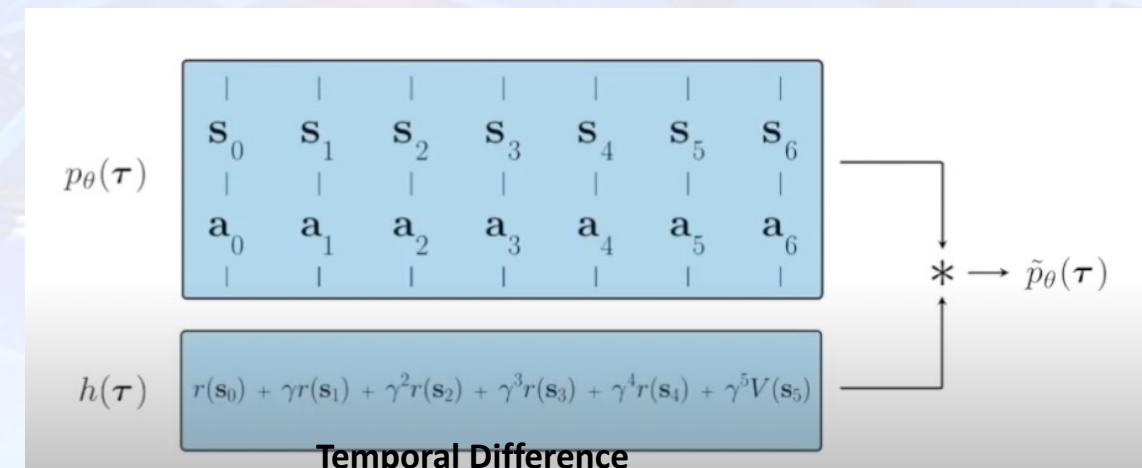
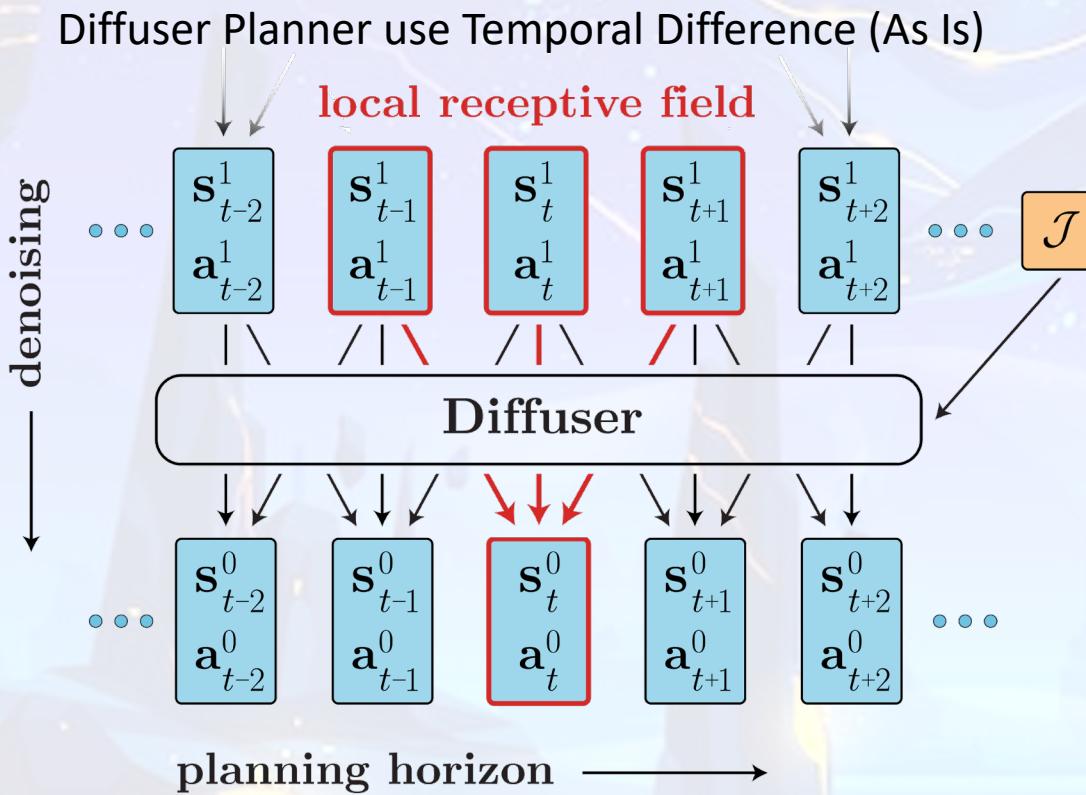


01 Decision Diffuser - 서론

서론

이 논문에서는 기존에 Diffuser Planner(Janner et al.)에서 사용했던 Guidance를 변경했다가 주 포인트입니다.

Classifier Guidance -> Classifier free guidance



01 Decision Diffuser - 서론

서론

We formulate the problem of **offline decision-making** from the perspective of **conditional generative modeling**.

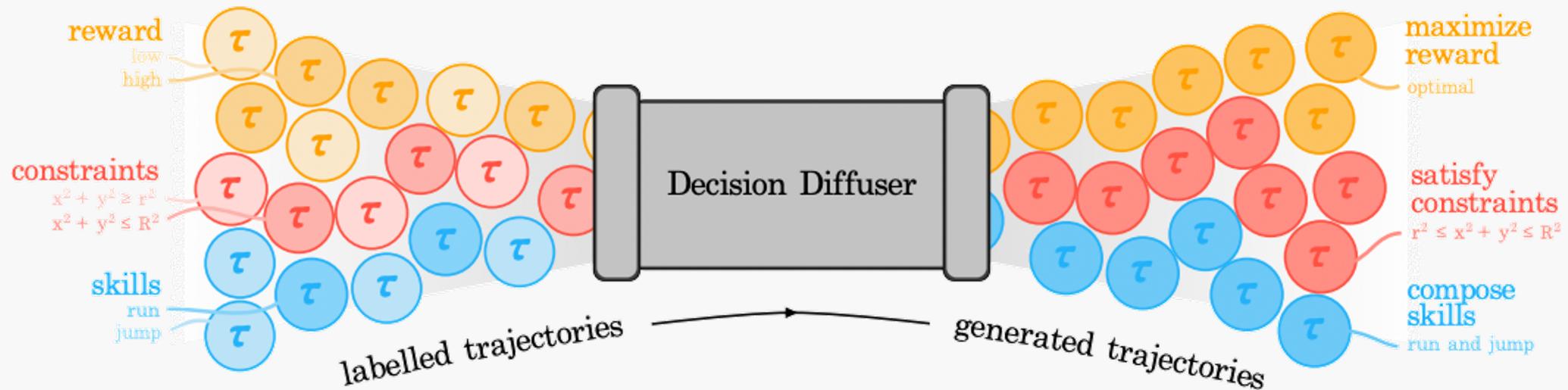
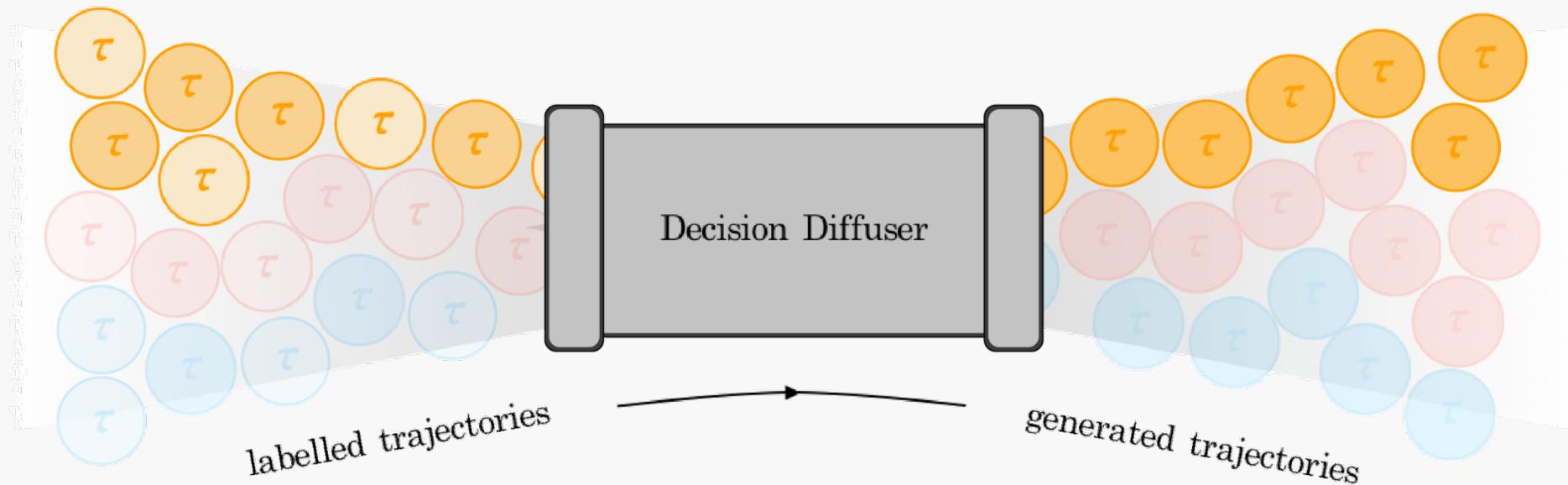


Figure 1: **Decision Making using Conditional Generative Modeling.** Framing decision making as a conditional generative modeling problem allows us to maximize rewards, satisfy constraints and compose skills.

01 Decision Diffuser - 서론

서론

given trajectories with **sub-optimal reward**, we can generate trajectories with **maximum reward**



1. Make Objective
2. Merge multiple Constraints
3. Compose Skills

03 Generative modeling with Decision Diffuser

모델 관련

Problem Definition $\max_{\theta} \mathbb{E}_{\tau \sim \mathcal{D}} [\log p_{\theta}(\mathbf{x}_0(\tau) | \mathbf{y}(\tau))]$

Noising process: $q(\mathbf{x}_{k+1}(\tau) | \mathbf{x}_k(\tau)),$

Denoising process: $p_{\theta}(\mathbf{x}_{k-1}(\tau) | \mathbf{x}_k(\tau), \mathbf{y}(\tau))$

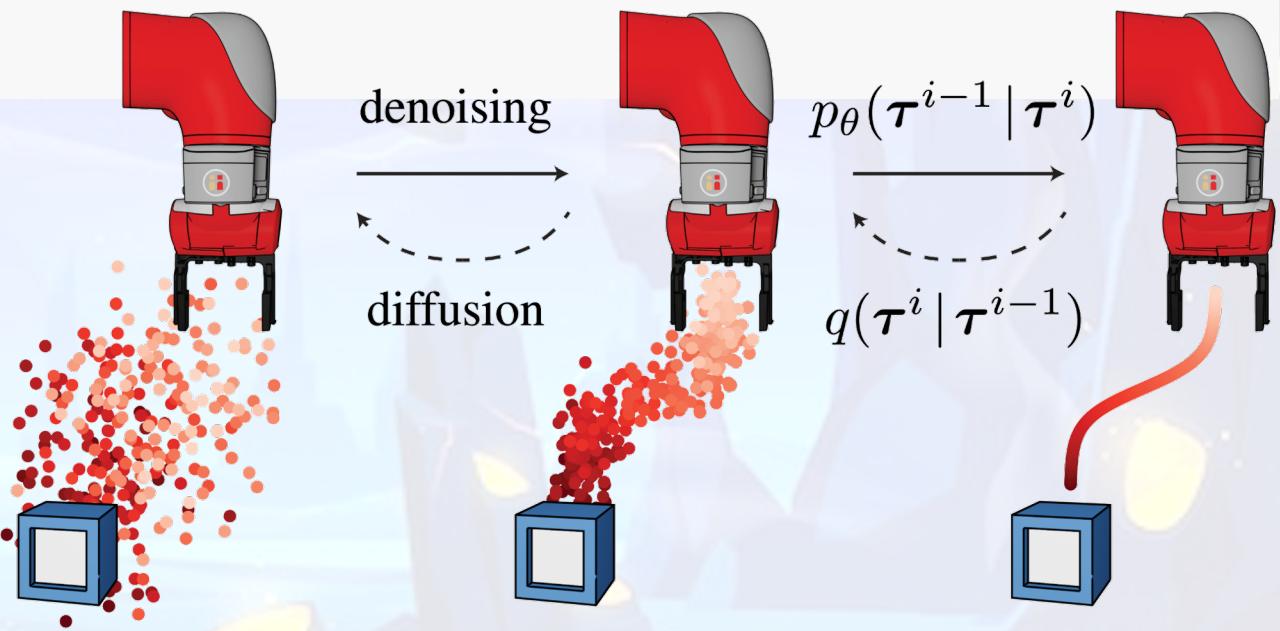
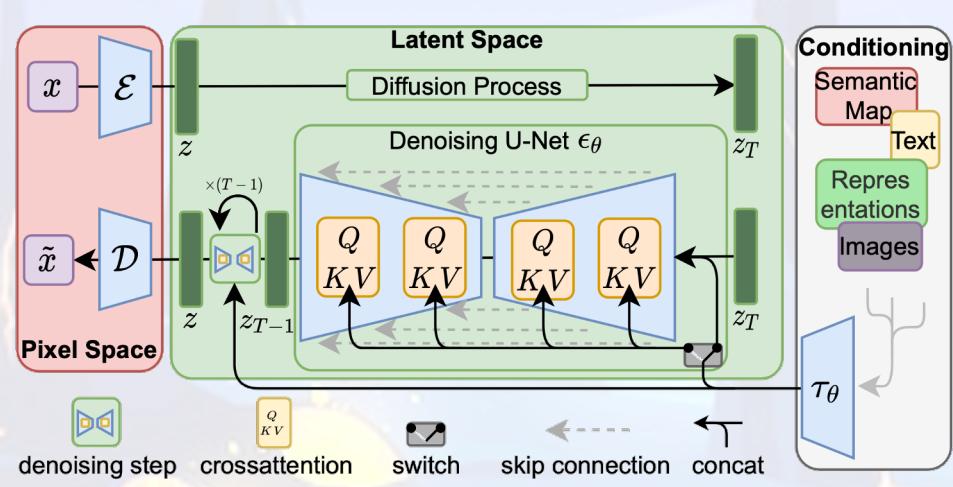
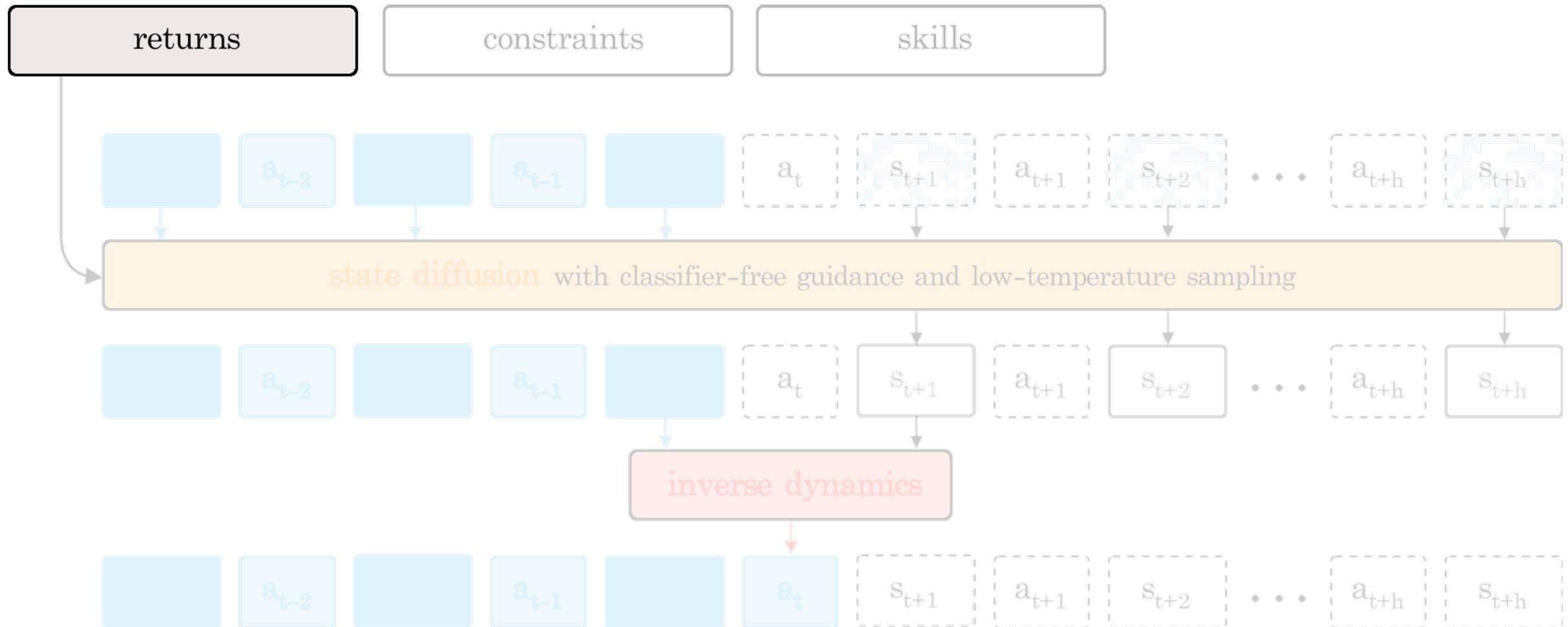


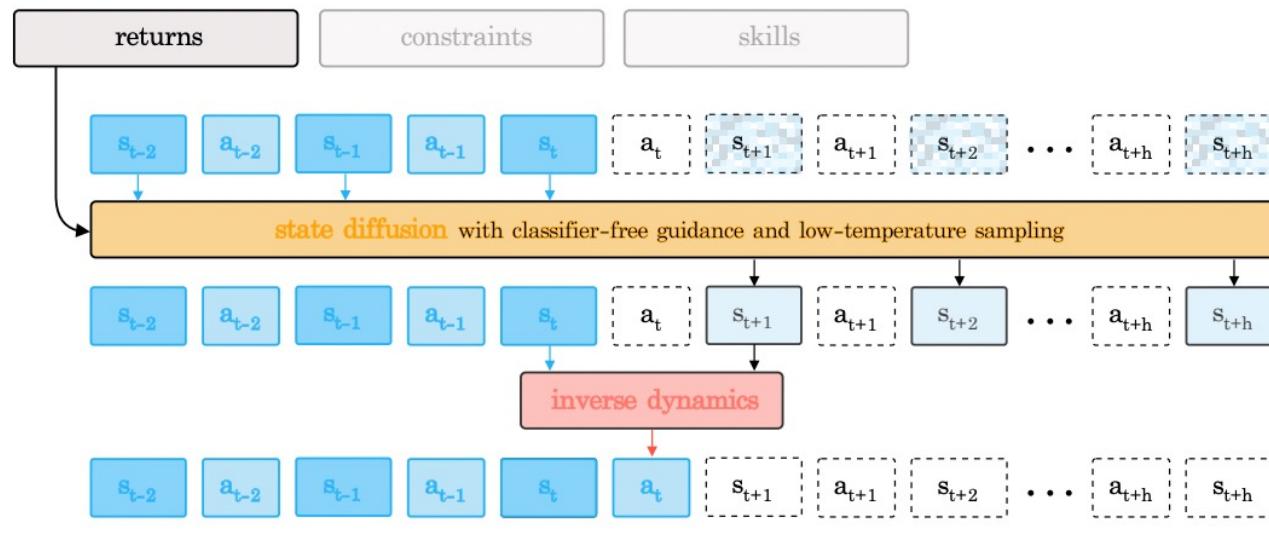
Figure. Diffusion model in Vision

02 Sampling

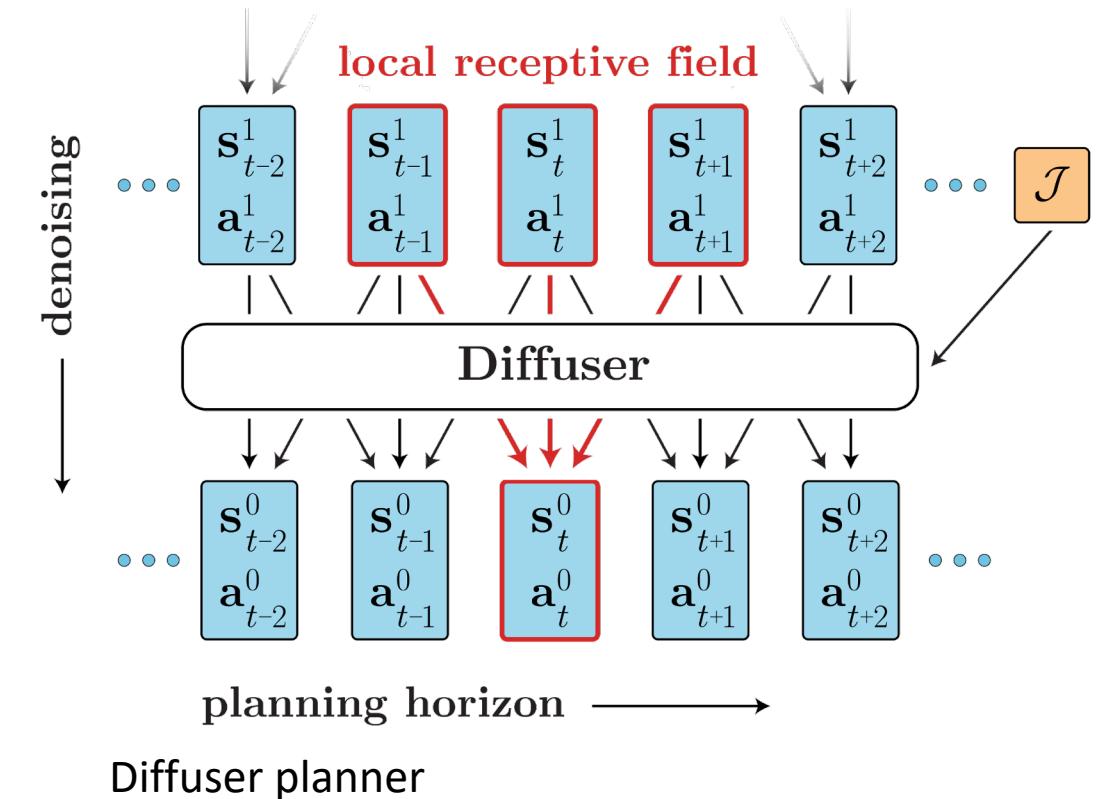
Sampling은 다음과 같이 할 수 있습니다.



02 Sampling



Decision diffuser



Diffuser planner

03 Experiments

Compare with Offline RL

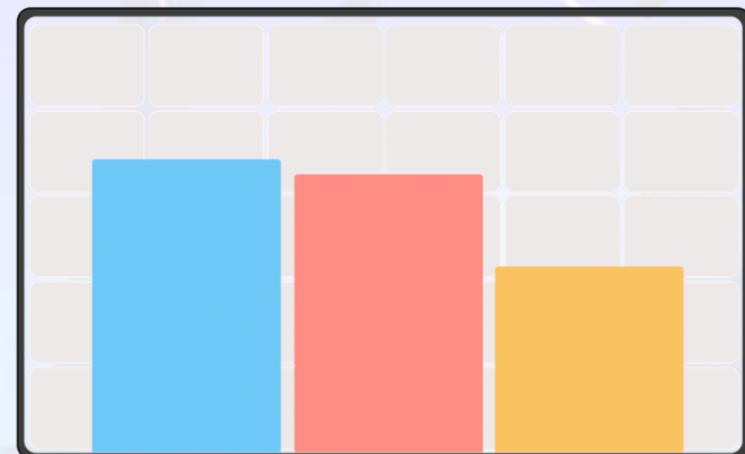
Decision Diffuser가 성능이 좋았다고 합니다.

Decision Diffuser

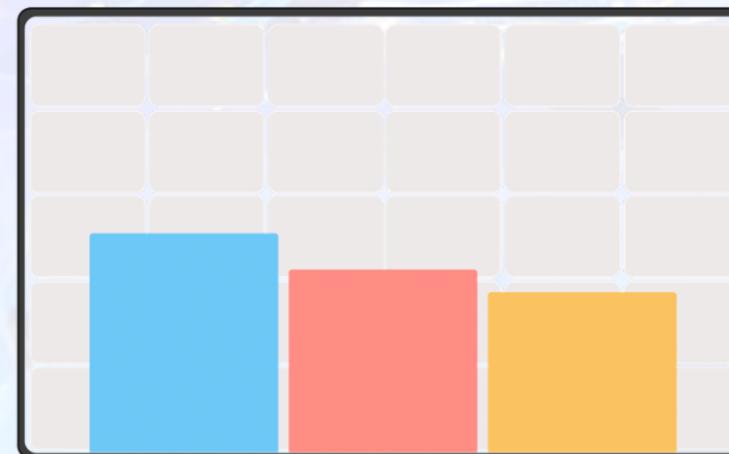
TD-learning

Behavior Cloning

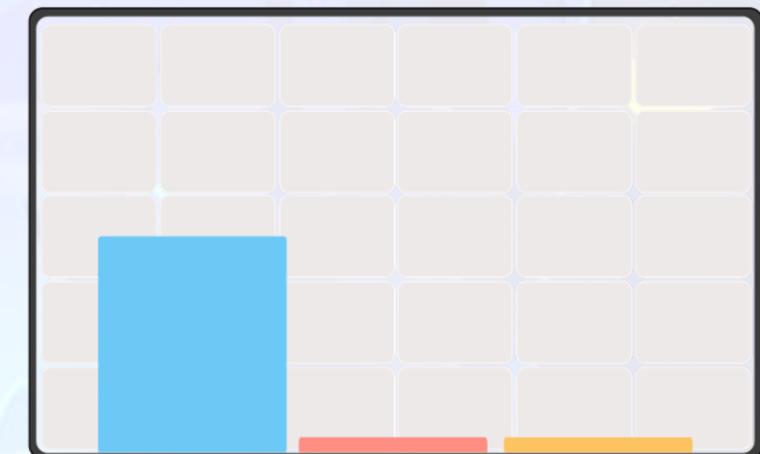
Performance



D4RL Locomotion



D4RL Kitchen



Kuka Block Stacking

03 Experiments

Performance

Dataset	Environment	BC	CQL	IQL	DT	TT	MOReL	Diffuser	DD
Med-Expert	HalfCheetah	55.2	91.6	86.7	86.8	95	53.3	79.8	90.6 ± 1.3
Med-Expert	Hopper	52.5	105.4	91.5	107.6	110.0	108.7	107.2	111.8 ± 1.8
Med-Expert	Walker2d	107.5	108.8	109.6	108.1	101.9	95.6	108.4	108.8 ± 1.7
Medium	HalfCheetah	42.6	44.0	47.4	42.6	46.9	42.1	44.2	49.1 ± 1.0
Medium	Hopper	52.9	58.5	66.3	67.6	61.1	95.4	58.5	79.3 ± 3.6
Medium	Walker2d	75.3	72.5	78.3	74.0	79	77.8	79.7	82.5 ± 1.4
Med-Replay	HalfCheetah	36.6	45.5	44.2	36.6	41.9	40.2	42.2	39.3 ± 4.1
Med-Replay	Hopper	18.1	95	94.7	82.7	91.5	93.6	96.8	100 ± 0.7
Med-Replay	Walker2d	26.0	77.2	73.9	66.6	82.6	49.8	61.2	75 ± 4.3
Average		51.9	77.6	77	74.7	78.9	72.9	75.3	81.8
Mixed	Kitchen	51.5	52.4	51	-	-	-	-	65 ± 2.8
Partial	Kitchen	38	50.1	46.3	-	-	-	-	57 ± 2.5
Average		44.8	51.2	48.7	-	-	-	-	61

Table 1: **Offline Reinforcement Learning Performance.** We show that Decision Diffuser (DD) either matches or outperforms current offline RL approaches on D4RL tasks in terms of normalized average returns (Fu et al., 2020). We report the mean and the standard error over 5 random seeds.

04 IMPORTANCE OF LOW TEMPERATURE SAMPLING

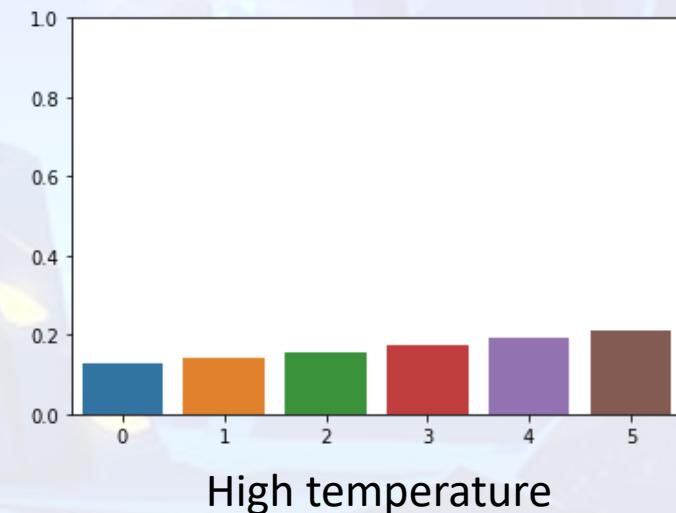
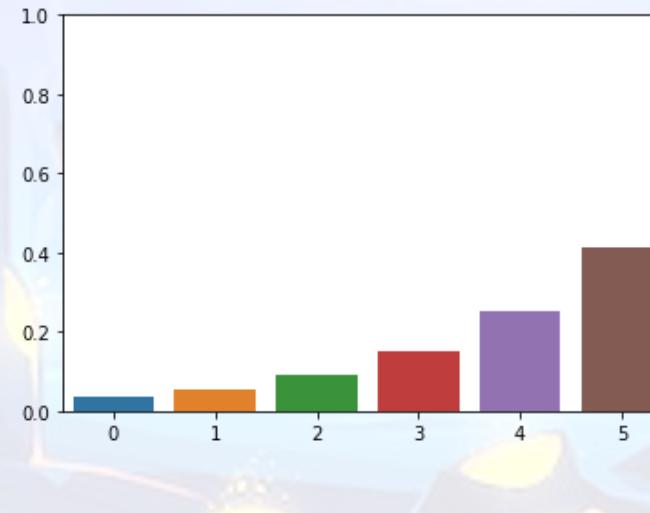
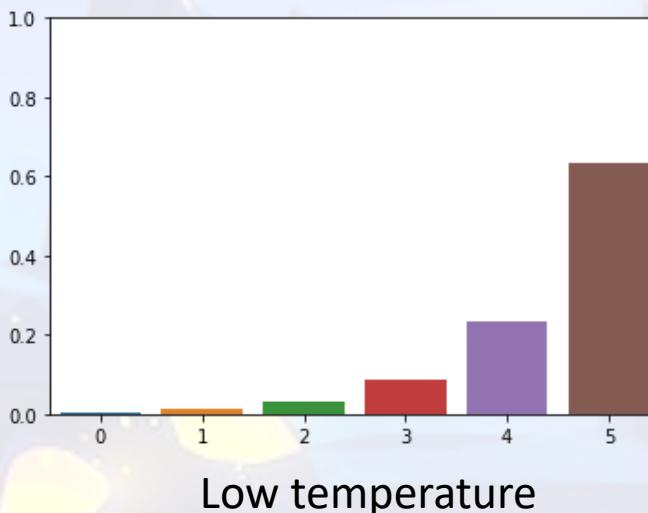
적당한 Temperature를 설정하는 것이 결과에 큰 영향을 줬다고 합니다.

Sampling : $x_{k-1} \sim N(\mu_{k-1}, \alpha\Sigma_{k-1})$

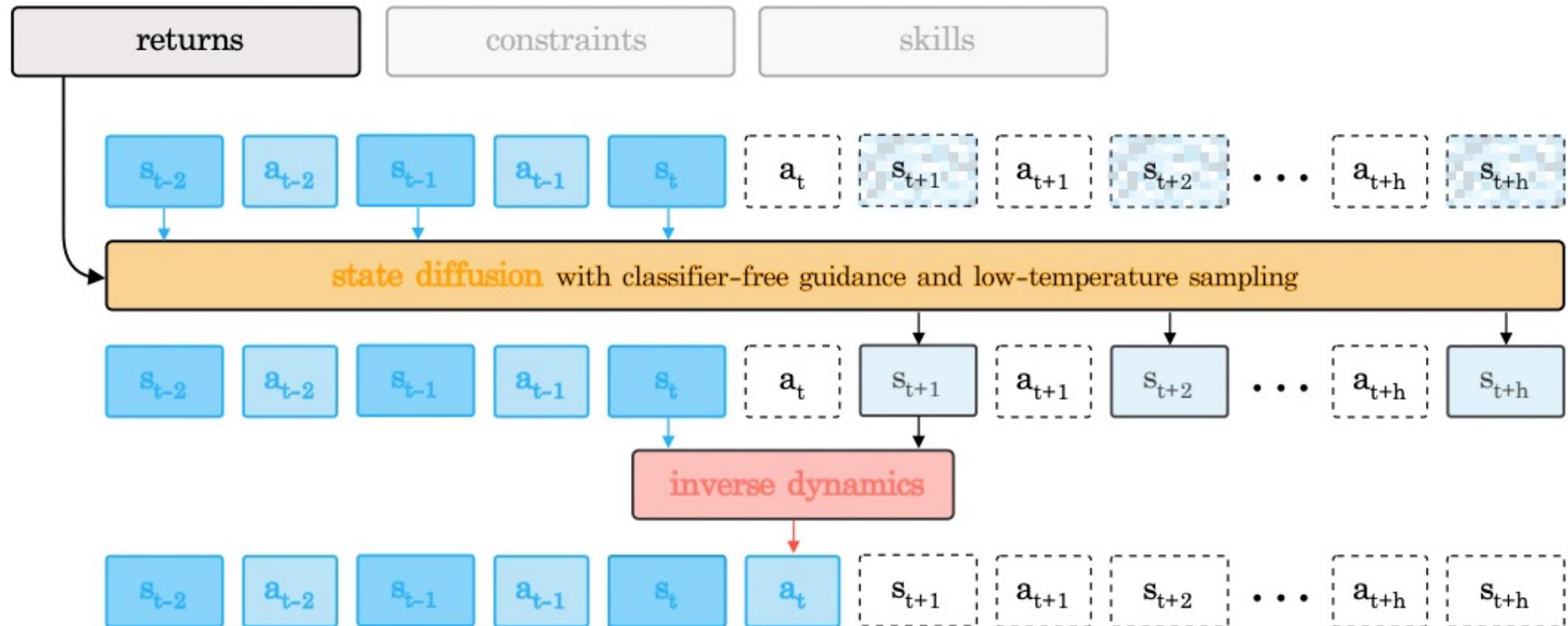
Results

Decision Diffuser	Hopper-Medium-Expert
$\alpha = 0$	104.3 ± 0.7
$\alpha = 0.5$	111.8 ± 1.6
$\alpha = 1.0$	107.1 ± 3.5

$$\text{softmax}(x)_i = \frac{e^{\frac{y_i}{T}}}{\sum_j^N e^{\frac{y_j}{T}}}$$

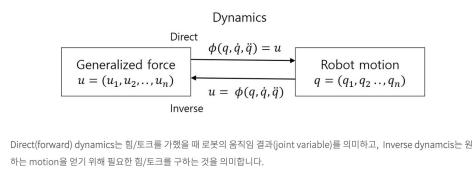


04 INVERSE DYNAMICS



Inverse dynamics model : $a_t := f_\phi(s_t, s_{\{t+1\}})$

State Diffusion using Inverse dynamics



04 INVERSE DYNAMICS

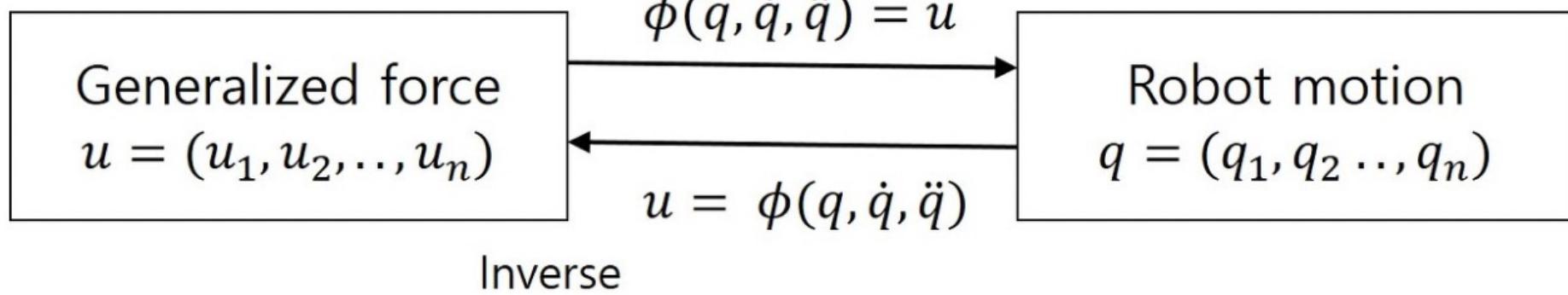
returns

constraints

skills

Dynamics

Direct



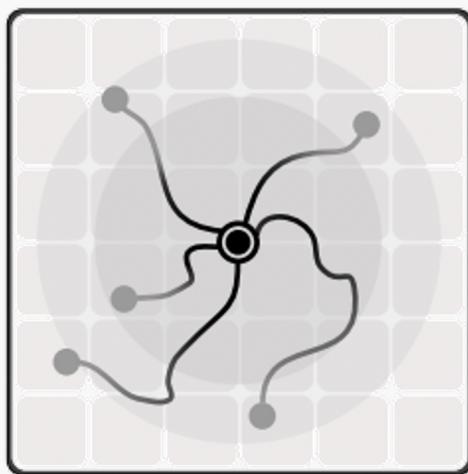
Inverse

Direct(forward) dynamics는 힘/토크를 가했을 때 로봇의 움직임 결과(joint variable)를 의미하고, Inverse dynamics는 원하는 motion을 얻기 위해 필요한 힘/토크를 구하는 것을 의미합니다.

04 Decision Diffuser - 서론

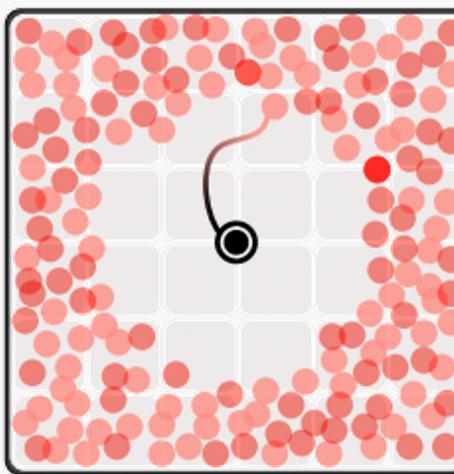
서론

Environment



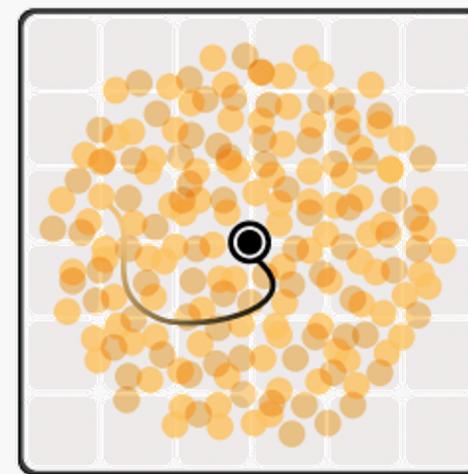
$$(x, y)$$

Training Dataset

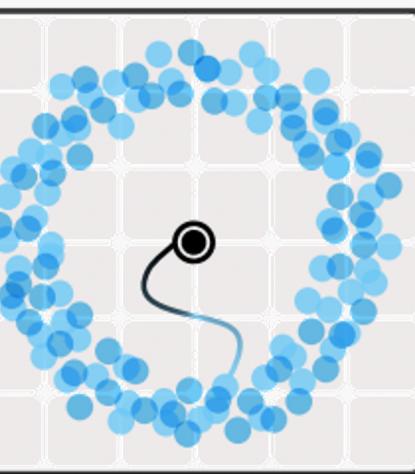


$$x^2 + y^2 \geq r^2$$

Generation



$$x^2 + y^2 \leq R^2$$



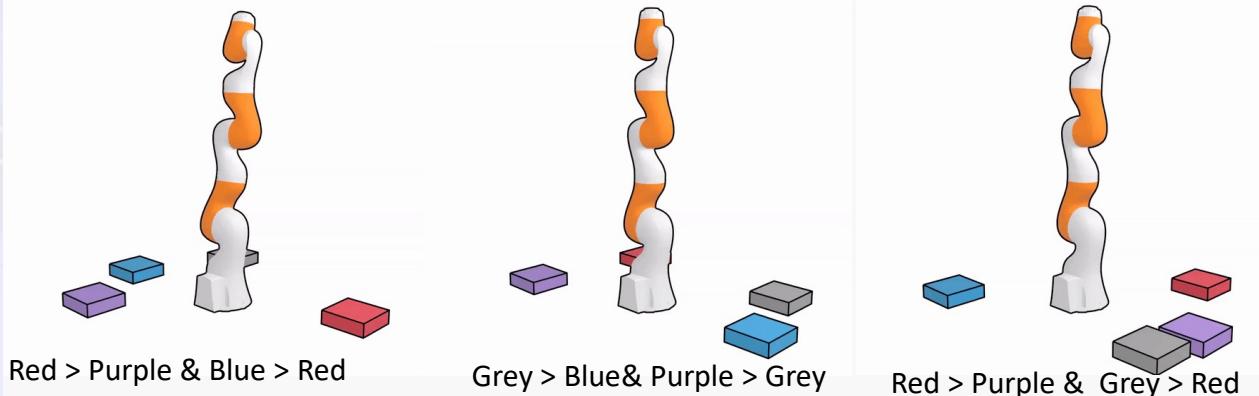
$$r^2 \leq x^2 + y^2 \leq R^2$$

Figure 2: **Illustrative example.** We visualize the 2d robot navigation environment and the constraints satisfied by the trajectories in the dataset derived from the environment. We show the ability of the conditional diffusion model to generate trajectories that satisfy the combined constraints.

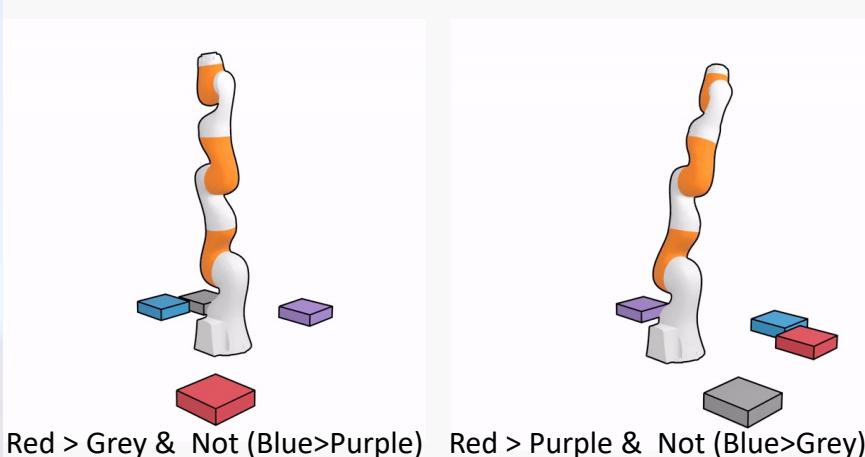
04 Skill Composition

COMPOSING

Combining Stacking Constraints



'NOT' constraints in Stacking and Rearrangement



$$\epsilon_{\theta}(x_k(\tau), y(\tau), k) := \epsilon_{\theta}(x_k(\tau), \mathbb{I}(\tau \in C_i), k)$$

One hot encoding을 condition으로 줍니다.

04 Skill Composition

Decision diffuser에서는 여러가지 Skill을 결합할수 있습니다.

SKILL COMPOSITION

Pace



$$\hat{\epsilon} := \epsilon_{\theta}(\mathbf{x}_k(\tau), \emptyset, k) + \omega \left(\sum_{i \neq j} (\epsilon_{\theta}(\mathbf{x}_k(\tau), \mathbf{y}^i(\tau), k) - \epsilon_{\theta}(\mathbf{x}_k(\tau), \emptyset, k)) \right. \\ \left. - (\epsilon_{\theta}(\mathbf{x}_k(\tau), \mathbf{y}^j(\tau), k) - \epsilon_{\theta}(\mathbf{x}_k(\tau), \emptyset, k)) \right)$$

Trott



Bound





Trott



Pace

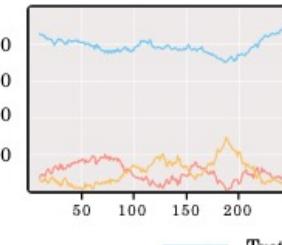


Trott + Pace

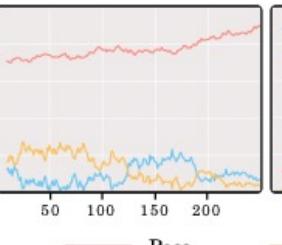


Trott + Pace

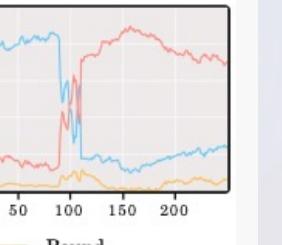
Only Trott



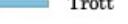
Only Pace



Trott + Pace



Trott



Pace



Bound



04 Skill Composition

COMPOSING

Composing Skills, Constraints가 만들어지는 방법

$$\begin{aligned}
 q(\mathbf{x}_k(\tau) | \{\mathbf{y}^i(\tau)\}_{i=1}^n) &\propto q(\mathbf{x}_k(\tau)) \prod_{i=1}^n \frac{q(\mathbf{x}_k(\tau) | \mathbf{y}^i(\tau))}{q(\mathbf{x}_k(\tau))} \quad (\text{Bayes Rule}) \\
 \Rightarrow \log q(\mathbf{x}_k(\tau) | \{\mathbf{y}^i(\tau)\}_{i=1}^n) &\propto \log q(\mathbf{x}_k(\tau)) + \sum_{i=1}^n (\log q(\mathbf{x}_k(\tau) | \mathbf{y}^i(\tau)) - \log q(\mathbf{x}_k(\tau))) \\
 \Rightarrow \nabla_{\mathbf{x}_k(\tau)} \log q(\mathbf{x}_k(\tau) | \{\mathbf{y}^i(\tau)\}_{i=1}^n) &= \nabla_{\mathbf{x}_k(\tau)} \log q(\mathbf{x}_k(\tau)) \\
 &\quad + \sum_{i=1}^n (\nabla_{\mathbf{x}_k(\tau)} \log q(\mathbf{x}_k(\tau) | \mathbf{y}^i(\tau)) - \nabla_{\mathbf{x}_k(\tau)} \log q(\mathbf{x}_k(\tau))) \\
 \Rightarrow \epsilon_\theta(\mathbf{x}_k(\tau), \{\mathbf{y}^i(\tau)\}_{i=1}^n, k) &= \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k) + \sum_{i=1}^n (\epsilon_\theta(\mathbf{x}_k(\tau), \mathbf{y}^i(\tau), k) - \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k))
 \end{aligned}$$

결과적으로 나오는 식

$$\begin{aligned}
 \hat{\epsilon} &:= \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k) + \omega(\epsilon_\theta(\mathbf{x}_k(\tau), \{\mathbf{y}^i(\tau)\}_{i=1}^n, k) - \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k)) \\
 &= \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k) + \omega \sum_{i=1}^n (\epsilon_\theta(\mathbf{x}_k(\tau), \mathbf{y}^i(\tau), k) - \epsilon_\theta(\mathbf{x}_k(\tau), \emptyset, k))
 \end{aligned}$$

04 Guidance

Q-Function guided diffusion vs Classifier-free guided diffusion

Algorithm 1 Guided Diffusion Planning

```

1: Require Diffuser  $\mu_\theta$ , guide  $\mathcal{J}$ , scale  $\alpha$ , covariances  $\Sigma^i$ 
2: while not done do
3:   Observe state  $s$ ; initialize plan  $\tau^N \sim \mathcal{N}(\mathbf{0}, I)$ 
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  $\tau_{\mathbf{a}_0}^0$ 

```

Algorithm 1 Conditional Planning with the Decision Diffuser

```

1: Input: Noise model  $\epsilon_\theta$ , inverse dynamics  $f_\phi$ , guidance scale  $\omega$ , history length  $C$ , condition  $\mathbf{y}$ 
2: Initialize  $h \leftarrow \text{Queue}(\text{length} = C)$ ,  $t \leftarrow 0$  // 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:      $\mathbf{x}_k(\tau)[:\text{length}(h)] \leftarrow h$  // Constrain plan to be consistent with history
7:      $\hat{\epsilon} \leftarrow \epsilon_\theta(\mathbf{x}_k(\tau), k) + \omega(\epsilon_\theta(\mathbf{x}_k(\tau), \mathbf{y}, k) - \epsilon_\theta(\mathbf{x}_k(\tau), k))$  // Classifier-free guidance
8:      $(\mu_{k-1}, \Sigma_{k-1}) \leftarrow \text{Denoise}(\mathbf{x}_k(\tau), \hat{\epsilon})$ 
9:      $\mathbf{x}_{k-1} \sim \mathcal{N}(\mu_{k-1}, \alpha \Sigma_{k-1})$ 
10:   end for
11:   Extract  $(s_t, s_{t+1})$  from  $x_0(\tau)$ 
12:   Execute  $a_t = f_\phi(s_t, s_{t+1})$ ;  $t \leftarrow t + 1$ 
13: end while

```

1. Train에서는 Classifier guided가 유리하지만, Testtime에서 Classifier free가 더 성능이 좋았다고 합니다.
2. Offline RL에서 Q-function은 Q-value를 높게 평가하는 경향이 있습니다. 대신에 Offline RL 세팅을 사용하는 diffuser에서 Q-function은 이런 error가 생길 수 있습니다. Classifier는 Q-function을 사용하지 않기에 이 문제가 발생하지 않습니다.

Limitations of Decision Diffuser

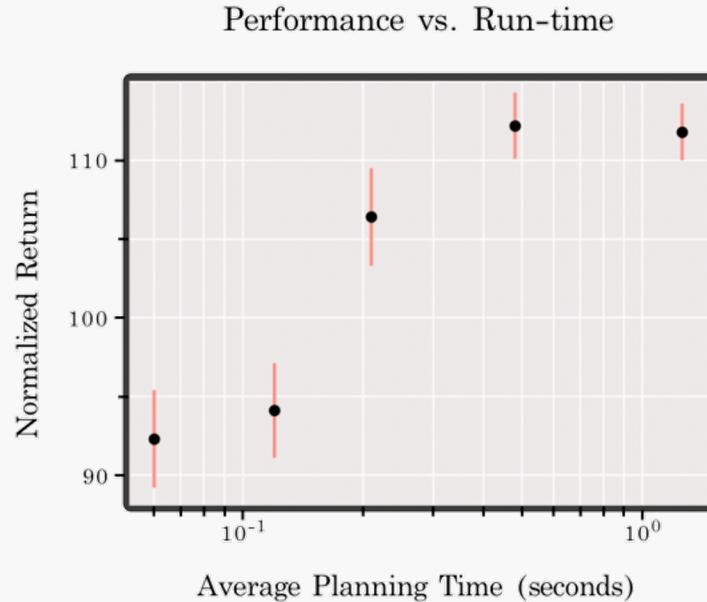


Figure A1: Performance vs planning time. We visualize the trade-off between performance, measured by normalized average return achieved in the environment, and planning time, measured in wall-clock time after warm-starting the reverse diffusion process.

run forward diffusion for a fixed number of steps (100 \rightarrow 40), and finally run the same number of reverse diffusion steps from the partially noised state sequence to generate another state sequence

5.4 Warm-Starting Diffusion for Faster Planning

A limitation of Diffuser is that individual plans are slow to generate (due to iterative generation). Naïvely, as we execute plans open loop, a new plan must be regenerated at each step of execution. To improve execution speed of Diffuser, we may further reuse previously generated plans to warm-start generations of subsequent plans.

To warm-start planning, we may run a limited number of forward diffusion steps from a previously generated plan and then run a corresponding number of denoising steps from this partially noised trajectory to regenerate an updated plan. In Figure 7, we illustrate the trade-off between performance and runtime budget as we vary the underlying number of denoising steps used to regenerate each a new plan from 2 to 100. We find that we may reduce the planning budget of our approach markedly with only modest drop in performance.

- Janner

Algorithm 1 Guided Diffusion Planning

```

1: Require Diffuser  $\mu_\theta$ , guide  $\mathcal{J}$ , scale  $\alpha$ , covariances  $\Sigma^i$ 
2: while not done do
3:   Observe state  $s$ ; initialize plan  $\tau^N \sim \mathcal{N}(0, I)$ 
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  $\tau_{s_0}^0$ 

```

05 한계점 (Limitations)

Limitations of Decision Diffuser

No partial observability (Offline RL)

Partial Observable environment에서는 offline RL에서 self-delusion 문제가 발생할 수 있습니다.

Inability to explore the environment and update itself in online setting (DT, DF 공통)

Online RL에서 사용하기에는 아직은 한계가 있습니다.

Experiments on only state-based environments

본 논문에서는 State 기반의 환경에서만 실행했습니다.

Only AND and NOT compositions are supported

Constraints \sqcap Skill을 OR 조건으로 결합하는 것은 아직 못했다고 합니다.

Performance degradation in environments with stochastic dynamics

Stochastic Dynamics를 갖는 환경에서는 성능이 떨어진다고 합니다.

Performance in limited data regime : Diffusion model처럼 Overfitting 문제가 발생합니다.

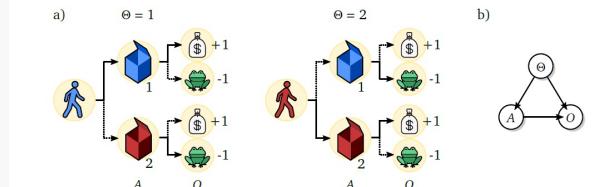


Figure 1 | The prize-or-frog problem. The objective is to choose the box containing the prize. Two problem instances exist: either the prize is in box 1 ($\Theta = 1$) or box 2 ($\Theta = 2$). There is also an expert who knows which box to open provided they are told the value of Θ . Panel (a) depicts the two configurations. Solid transitions are deterministic (taken with probability one) and dotted ones are possible transitions that aren't taken. Panel (b) shows the causal Bayesian network of the problem.



Latent model