

WHAT MAKES A GOOD DIFFUSION PLANNER FOR DECISION MAKING?

Anonymous authors

Paper under double-blind review

ABSTRACT

Diffusion models have recently shown significant potential in solving decision-making problems, particularly in generating behavior plans – also known as diffusion planning. While numerous studies have demonstrated the impressive performance of diffusion planning, the mechanisms behind the key components of a good diffusion planner remain unclear and the design choices are highly inconsistent in existing studies. In this work, we address this issue through systematic empirical experiments on diffusion planning in an offline reinforcement learning (RL) setting, providing practical insights into the essential components of diffusion planning. We trained and evaluated over 6,000 diffusion models, identifying the critical components such as guided sampling, network architecture, action generation and planning strategy. We revealed that some design choices opposite to the common practice in previous work in diffusion planning actually lead to better performance, e.g., unconditional sampling with selection can be better than guided sampling and Transformer outperforms U-Net as denoising network. Based on these insights, we suggest a simple yet strong diffusion planning baseline that achieves state-of-the-art results on standard offline RL benchmarks.

1 INTRODUCTION

Decision making by learning from offline data has been a fundamental approach in robotics and artificial intelligence (Bellman, 1957). It enables agents to acquire complex behaviors by observing and mimicking expert demonstrations, circumventing the need for explicit programming or exhaustive exploration. However, this paradigm faces significant challenges, particularly when dealing with long-horizon planning and high-dimensional action spaces. The complexity of modeling sequential dependencies and capturing the intricacies of action distributions makes it difficult to scale traditional methods (Deisenroth & Rasmussen, 2011) to more complex tasks (Parmas et al., 2018).

Recently, diffusion models have achieved remarkable success in image and video generation, demonstrating their ability to handle complex distribution and long-range dependencies (Ho et al., 2020; Dhariwal & Nichol, 2021). Inspired by these works, several recent studies have applied diffusion models to planning sequential decisions, especially with continuous state and action spaces such as robotic manipulation tasks (Janner et al., 2022; Ajay et al., 2022; Lu et al., 2023; Li et al., 2023). The diffusion models are used to approximate the sequence of states and actions from current time step to future – and by exploiting the diffusion models’ conditional generation capacity such as diffusion guidance (Ho et al., 2020; Ho & Salimans, 2021), the model can make plans (i.e. state trajectory) with desired properties such as reward maximization (i.e. offline reinforcement learning (Levine et al., 2020)).

Despite achieving impressive performance across a diverse array of tasks, there has been limited exploration into the fundamental components and mechanisms that constitute an effective diffusion planning model for decision making. Previous research exhibits a lack of consistency and coherence in design choices. It remains uncertain whether sub-optimal design choices might hinder the full potential of diffusion models within decision-making domains. Specifically, existing approaches have not adequately addressed essential facets such as the choice of diffusion guidance algorithm, network architecture, and whether the plan should contain states or state-action pairs. This naturally raises the following fundamental question:

What makes a good diffusion planner for decision making, especially offline RL?

We seek to answer the question by conducting a comprehensive empirical investigation into key design choices in diffusion models for decision-making, in particular for state-based robotics tasks. Our work contributes to the field of decision making and diffusion models in several aspects.

- **Comprehensive Experiments:** We conducted an extensive empirical study to explore what constitutes an effective diffusion planner. By training and evaluating over 6,000 models, we analyzed key components critical to decision making in diffusion planning, including guided sampling algorithms, network architectures, action generation methods, and planning strategies.
- **Insights and Tips:** We ran detailed experiments and data analysis to understand the role of each key component in constituting a good diffusion planner. In particular, we discovered that certain design choices, contrary to common practice in diffusion planning actually lead to better performance. Our work offers intuitive explanations and practical tips about the choices and provides insights about the strengths and limitations of diffusion planning.
- **A simple yet strong baseline:** Building on the insights from our study, we suggest a simple yet highly competitive baseline, named *Diffusion Veteran* (DV). This model achieves state-of-the-art performance in planning tasks in standard offline RL benchmarks.

2 BACKGROUND AND RELATED WORK

Offline Reinforcement Learning (Fujimoto et al., 2019; Levine et al., 2020; Fu et al., 2020) is a subfield of reinforcement learning (RL) where the agent learns from a fixed dataset of past experiences. This dataset typically consists of state-action-reward-next-state tuples, which encapsulate the agent’s interactions with the environment. The challenge in offline RL is for the agent to derive an effective policy from this static dataset without further exploration or interaction with the environment. Two major challenges arise in this context. First, the state and action spaces may be high-dimensional and involve long-range dependencies, making it difficult to model effectively (Levine et al., 2020). Second, the learned policy must be optimal, even though the behavior policy that generated the offline data may be sub-optimal or different from the desired policy (Fujimoto et al., 2019).

Recently, diffusion models have emerged as a powerful framework for tasks such as image and video generation due to their ability to model complex distributions (Croitoru et al., 2023), which could mitigate the first problem. Moreover, diffusion guidance techniques (Ho et al., 2020; Ho & Salimans, 2021) allow the model to generate samples that adhere to the desired properties. The second challenge in offline RL, learning an optimal policy, can be addressed by diffusion guidance techniques to produce behavior that maximizes rewards. Building on this insight, a growing body of research has explored the use of diffusion models to generate behavior trajectories, denoted as τ .

Diffusion planning (Ajay et al., 2022; Janner et al., 2022; Liang et al., 2023; Dai et al., 2023; Yang et al., 2023; Li et al., 2023; Yang et al., 2023; Chen et al., 2024; Dong et al., 2024c) considers that at the time step t , a trajectory τ consists of the current and subsequent H steps of state-action pairs or states:

$$\tau = \begin{bmatrix} s_t & s_{t+1} & \cdots & s_{t+H-1} \\ a_t & a_{t+1} & \cdots & a_{t+H-1} \end{bmatrix}, \text{ or } \tau = [s_t \quad s_{t-1} \quad \cdots \quad s_{t+H-1}]. \quad (2.1)$$

There is a guidance function to model the reward, such as the immediate reward r_t or the state value function $v(s_t) = \mathbb{E} \left[\sum_{h=0}^{\text{end}} \gamma^h r_{t+h} \right]$, where γ is the discount factor (Sutton & Barto, 1998). In classifier guidance (CG) (Ho et al., 2020), a guidance network is learned simultaneously with the diffusion model, whose input is the generated trajectory and the output is accumulated rewards or value function. The gradient of the guidance network is used in the generation process of diffusion model to maximize the rewards. Examples of diffusion planning with CG are (Janner et al., 2022; Liang et al., 2023; Zhang et al., 2022) In classifier-free guidance (CFG) (Ho & Salimans, 2021), it takes the desired reward or value function as an additional argument feed into the diffusion process. Instances are (Ajay et al., 2022; Li et al., 2023; Yang et al., 2023). However, despite some literature reviews such as Zhu et al. (2023), **the field lack a systematical study to elucidate the design space of diffusion planning in offline RL with substantial experimental results.**

Diffusion policy (Pearce et al., 2023; Wang et al., 2023b; Hansen-Estruch et al., 2023; Chen et al., 2023) is another kind of popular usage of diffusion model in decision making. The trajectory only

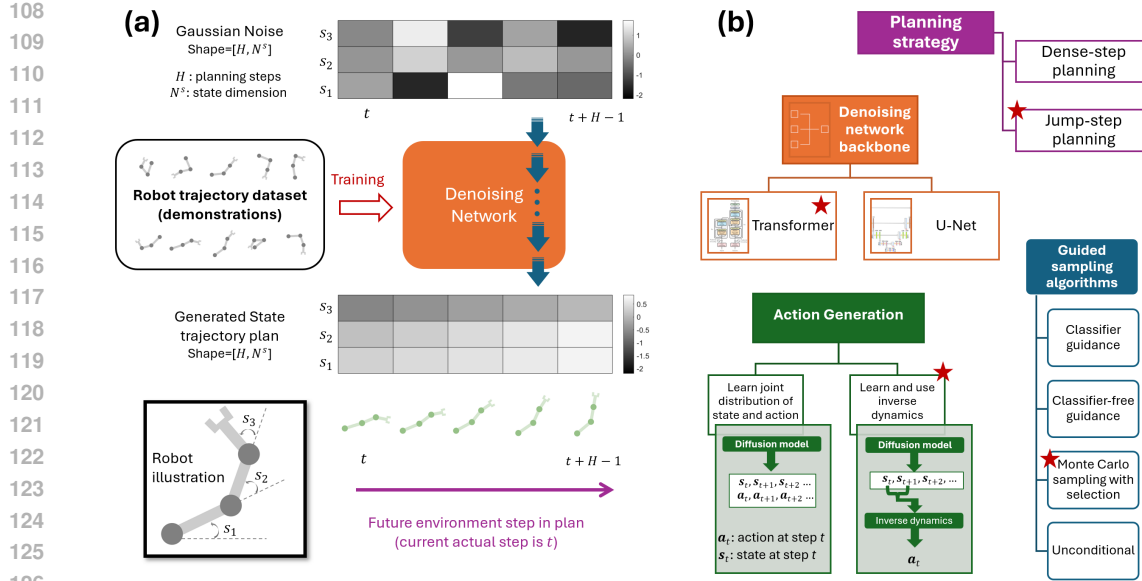


Figure 1: **Diffusion planning framework for decision making.** (a) The generation of a sequence plan using the denoising process of a diffusion model. A 3-joints robot arm is used as an illustrative example. (b) Important components and candidates in the framework. Each color corresponds to one component in the framework. A star indicates the preferred choice in experiments.

includes $\tau = a_t$, without lookahead planning. The model is trained by combining the loss of imitation learning and model-free RL as in classic offline RL methods (Kumar et al., 2020; Fujimoto & Gu, 2021). Diffusion policy methods hope to improve the performance of by leveraging the capacity of diffusion model to model complex distribution of actions (policy function). A recent study (Dong et al., 2024b) investigated the design space of diffusion policy, proposed that diffusion policies such as DQL (Wang et al., 2023b) can be a computationally efficient and powerful candidate for decision-making tasks.

3 STUDY DESIGN

3.1 KEY COMPONENTS AND MECHANISMS OF DIFFUSION PLANNER

Recent pioneering work in diffusion planning (Janner et al., 2022; Ajay et al., 2022; Chen et al., 2024) has demonstrated the potential of this approach in offline RL. However, the design choices in these studies vary significantly, and it remains unclear whether there is an optimal configuration for different domains. Our aim is to conduct a systematic analysis supported by comprehensive experimental results. To achieve this, we begin by listing key design components (excluding common deep learning hyperparameters such as learning rates) that have varied in previous studies. See Fig. 1(b) for an overview.

Guided sampling algorithms: Classifier guidance (CG) (Ho et al., 2020), Classifier-free guidance (CFG) (Ho & Salimans, 2021), Monte Carlo sampling with selection (sample N *unconditional* trajectories and select the best, the criteria of which is given by a critic function learned simultaneously with diffusion model). Most previous diffusion planners used CG (Janner et al., 2022; Wang et al., 2023a; Chen et al., 2024) or CFG (Ajay et al., 2022; Li et al., 2023; Yang et al., 2023) for offline RL.

Denoising network backbone: U-Net (Ronneberger et al., 2015); Transformer (Vaswani et al., 2017). U-Net was used in most previous diffusion planners for state-based offline RL (Janner et al., 2022; Ajay et al., 2022; Wang et al., 2023a; Li et al., 2023; Chen et al., 2024)

Action Generation: Learn joint distribution of state and action and directly execute the generated action at the current step (used in, e.g. Janner et al. (2022); Liang et al. (2023)); Learn and use inverse dynamics to compute action from state plan (used in e.g., Ajay et al. (2022); Wang et al. (2023a)).

Planning strategy: Dense-step planning means the planned trajectory τ (Eq. 2.1) corresponds to contiguous H steps in the environment (this is a conventional setting in diffusion planning (Janner et al., 2022; Ajay et al., 2022; Lu et al., 2023)) ; Jump-step planning models $H \times m$ environment steps, where $m \in \mathbb{N}^+$ is the planning stride; Hierarchical planning (studied by Li et al. (2023); Chen et al. (2024)).

Details of the implementation are deferred to Appendix A and B.

3.2 EXPERIMENT PROCEDURE

Given the multitude of components involved, it is challenging to draw scientific conclusions directly from the collective results. Therefore, we structured our study using the following procedure:

- (1) Conduct a comprehensive search on the key components (Sect. 3.1) by combining grid search and manual tuning to obtain the best results.
- (2) Evaluate the effect of each component using the control variable method; that is, modify only one component of the best model at a time and compare it with the original.
- (3) After identifying which components are important and understanding how they affect performance, perform a deeper analysis to derive useful insights.

3.3 BENCHMARK

We conducted experiments on the D4RL dataset (Fu et al., 2020), one of the most widely used benchmarks for offline RL and imitation learning. The dataset covers a variety of task domains, including maze navigation, robot locomotion, robot arm manipulation, and vehicle driving, among others. For our experiments, we selected three sets of behavior planning tasks that were most commonly studied in prior works in offline RL and diffusion planning (Janner et al., 2021; Ajay et al., 2022; Janner et al., 2022; Liang et al., 2023; Li et al., 2023; Lu et al., 2023; Chen et al., 2024). These tasks (Fig. 2) encompass both planning and control challenges, providing a comprehensive evaluation in various problem settings. The performance metric considered in this work is the standard RL objective: the average total rewards in an online testing episode.

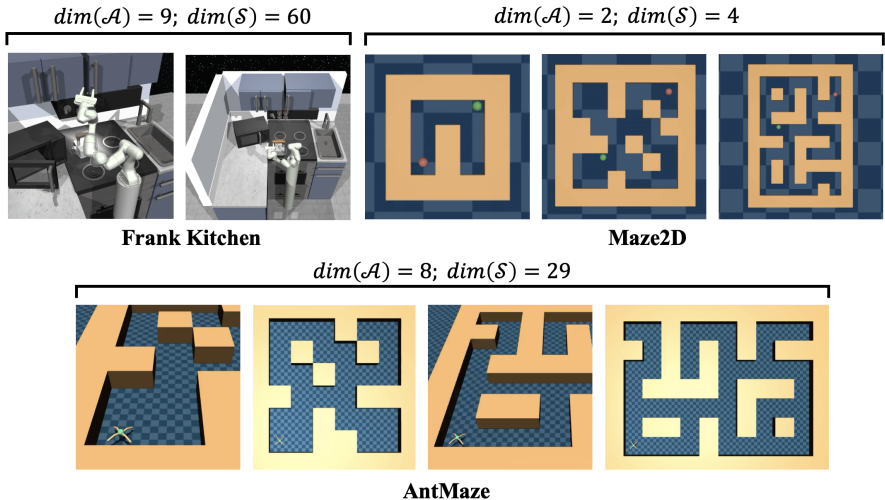


Figure 2: **Rendering of the benchmarking tasks considered in this study, where $dim(S)$ and $dim(A)$ denote the dimension of the state and action spaces.**

Maze2D involve navigating a 2D maze, requiring the agent to find an optimal path to a goal. These tasks are used to test planning capabilities in environments where spatial reasoning is critical.

AntMaze presents a navigation challenge with a simulated ant robot. The agent controls a multi-legged robot to navigate through a 2D maze, combining both locomotion and planning.

Franka Kitchen simulates a robot arm performing a variety of manipulation tasks in a kitchen environment to achieve task goals across multiple stages.

Category	Env	Kitchen			Antmaze				Maze2D				
		Mixed	Partial	avg.	L.-div.	L.-play	M.-div.	M.-play	avg.	L.	M.	Umaze	avg.
Non-diffusion	BC	47.5	33.8	40.7	0.0	0.0	0.0	0.0	0.0	5	30.3	3.8	13.0
	BCQ	8.1	18.9	13.5	2.2	6.7	0.0	0.0	2.2	6.2	8.3	12.8	9.1
	CQL	51.0	49.8	50.4	61.2	53.7	15.8	14.9	36.4	12.5	5.0	5.7	7.7
	IQL	51.0	46.3	48.7	47.5	39.6	70.0	71.2	57.1	58.6	34.9	47.4	47.0
Diffusion Policies	SfBC	45.4	47.9	46.7	45.5	59.3	82.0	81.3	67.0	74.4	73.8	73.9	74.0
	DQL	62.6	60.5	61.6	56.6	46.4	78.6	76.6	64.6	–	–	–	–
	DQL*	55.1	65.5	60.3	70.6	81.3	82.6	87.3	80.5	186.8	152.0	140.6	159.8
	IDQL	66.5	66.7	66.6	67.9	63.5	84.8	84.5	75.2	90.1	89.5	57.9	79.2
	IDQL*	66.5	66.7	66.6	40.0	48.7	83.3	67.3	59.8	–	–	–	–
	CEP	–	–	–	64.8	66.6	83.8	83.6	74.7	–	–	–	–
Diffusion Planners	Diffuser	52.5	55.7	54.1	27.3	17.3	2.0	6.7	13.3	123	121.5	113.9	119.5
	AdpDfsr	51.8	55.5	53.7	8.7	5.3	6.0	12.0	8.0	167.9	129.9	135.1	144.3
	DD	75.0	56.5	65.8	0.0	0.0	4.0	8.0	3.0	–	–	–	–
	HD	71.7	73.3	72.5	83.6	–	88.7	–	–	128.4	135.6	155.8	139.9
	DV (Ours)	73.6	94.0	83.8	80.0	76.4	87.4	89.0	83.2	203.6	150.7	136.6	163.6

Table 1: **Normalized performance of various offline-RL methods.** Our results (DV) are averaged over 500 episode seeds. The results of other methods are obtained from literature. We omit the variance over seeds for simplicity; however, it can be found in the detailed tables in Appendix C. The best average performance on each task set are marked in bold fonts. BC: vanilla imitation learning, BCQ: Fujimoto et al. (2019), CQL: Kumar et al. (2020), IQL: Kostrikov et al. (2021), SfBC: Chen et al. (2023), DQL: Wang et al. (2023b), IDQL: Hansen-Estruch et al. (2023), DQL* and IDQL*: replicated by Dong et al. (2024b), CEP: Lu et al. (2023), Diffuser: Janner et al. (2022), AdpDfsr: Liang et al. (2023), DD: Ajay et al. (2022), HD: Chen et al. (2024).

4 EXPERIMENTAL RESULTS

We trained and evaluated over 6,000 diffusion models by sweeping the key components discussed in Sect. 3.1 and other hyper-parameters (See Appendix B for details).

By summarizing the results from the experiments, we identified one kind of diffusion planning framework, called the Diffusion Veteran (DV). The pseudocode of DV can be found in Algorithm 1. As shown in Table 1, DV outperforms all previous diffusion planning and diffusion policy methods. We hope DV will serve as a simple yet strong baseline for future research in diffusion planning.

Algorithm 1: Diffusion Veteran (DV) Simplified Pseudocode

Input: Planning horizon H , Dataset \mathcal{D} , Discount factor γ , Candidate num N , Planning stride M .

Initialize: Diffusion Transformer Planner ϵ_θ , Diffusion Inverse dynamics ϵ_ω , Critic V_ϕ

1 Calculate accumulated discounted returns $R_t = \sum_{h=0}^{\text{end}} \gamma^h r_{t+h}$ for every step t .

2 **Function TRAINING:**

3 Sample $\mathbf{s}_{t,t+M,\dots,t+(H-1)M}$, $\mathbf{a}_{t,t+M,\dots,t+(H-1)M}$, R_t from \mathcal{D}

4 Train planner ϵ_θ using \mathbf{s}_t as condition and $\mathbf{s}_{t,t+M,\dots,t+(H-1)M}$ as target output

5 Train Inverse dynamics ϵ_ω using \mathbf{s}_t , \mathbf{s}_{t+M} as input, \mathbf{a}_t as target output

6 Train critic V_ϕ using $\mathbf{s}_{t,t+M,\dots,t+(H-1)M}$ as input, R_t as target output

7 **end**

8 **Function EXECUTION(\mathbf{s}):**

9 Randomly generate N plans using ϵ_θ , while fixing the first state as \mathbf{s} during sampling

10 Select the best plan using critic V_ϕ

11 Use the inverse dynamics ϵ_ω to generate action using \mathbf{s} and the next state in the best plan

12 **end**

With DV in place, we can then analyze the impact of each component in diffusion planning by looking into how each component influences its performance. Each of the following sub-sections will focus on one component that we have found to be crucial. In the end of this section, we will conclude our findings into practical tips.

270
271
272
273
274
275
276
277
278
279

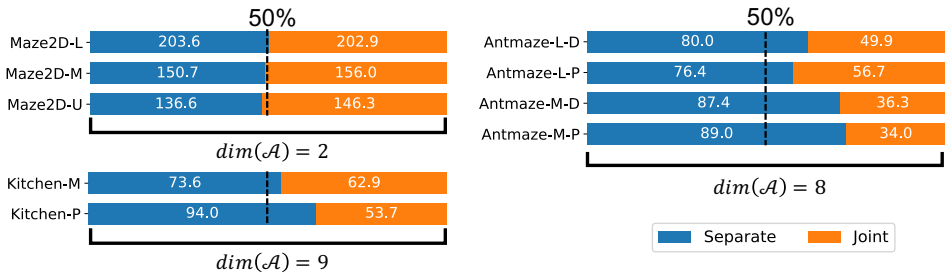


Figure 3: **Comparison of performance between two action generation strategies.** “Separate” learns and uses inverse dynamics to compute action from state plan. “Joint” means learning joint distribution of state and action and directly executing the generated action at the current step (see “action generation” in Fig. 1(b)). A straightforward conclusion drawn from the results is that “separate” is better than “Joint” when tackling higher-dimensional action spaces. The vertical dashed line indicates on-par performance.

4.1 ACTION GENERATION

The choice of action generation design (Sect. 3.1) remains a subject of ongoing debate within the field. On one side, the pioneering diffusion planner, Diffuser, along with subsequent studies (Janner et al., 2022; Liang et al., 2023; Chen et al., 2024), employs a diffusion model to generate the joint distribution of action and state trajectories (“joint”). In contrast, studies by Ajay et al. (2022); Wang et al. (2023a); Du et al. (2024) have adopted inverse dynamics to generate actions based on planned states (“separate”).

Our experimental findings favor the latter approach: Although both strategies perform comparably in simpler environments such as Maze2D, which lacks robotic control elements, the “separate” approach significantly outperforms the “joint” strategy in more complex settings like Kitchen and AntMaze, which feature robotic control and higher-dimensional action spaces.

This observed disparity may be attributed to the additional complexity introduced when modeling the joint distribution of sequential states and actions, compared to modeling only the states. This complexity becomes particularly pronounced in environments where state transitions involve more complex actions due to higher-dimensional action spaces.

We tested both diffusion models and vanilla MLP as the inverse dynamics, and found similar performance between them. We adhered to diffusion inverse dynamics (Appendix B.1).

4.2 PLANNING STRATEGY

310
311
312
313
314
315
316
317
318
319
320
321
322
323

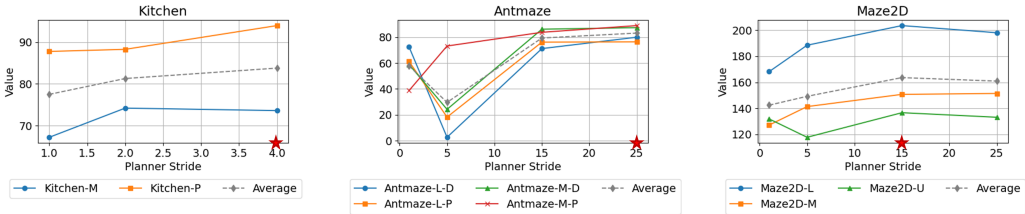


Figure 4: **Performance change of DV over planning stride.** It reduces to dense-step planning when Stride=1. The star indicates the choice of DV.

One crucial result we found is that jump-step planning (Sect. 3.1) is beneficial in almost all cases, despite the fact that most previous work used dense-step planning. This is observed in DV (Fig. 4) and generally in diffusion planners (see Appendix C for extensive results).

An obvious benefit from jump-step planning is that with the same planning steps, the model can look ahead farther. This may be crucial for planning tasks that require long-term credit assignment. The choice of stride should be related to the actual clock-time interval between two environment steps. Nonetheless, we suggest to try jump-steps and sweep the stride.

This observed phenomenon also implies that the diffusion planner should play the role of planning at a more abstract level or with a longer timescale. Interestingly, this is consistent with the neuroscientific fact that the intrinsic timescale of the prefrontal cortex (higher-level planning) is longer than that of the motor cortex (low-level control) (Murray et al., 2014; Runyan et al., 2017; Wang et al., 2018). A recent study (Chen et al., 2024) demonstrated impressive planning performance (Table 1, HD) using multi-timescale diffusion planning. Exploring the hierarchical paradigm of diffusion planning could be an interesting future direction.

4.3 DENOISING NETWORK BACKBONE

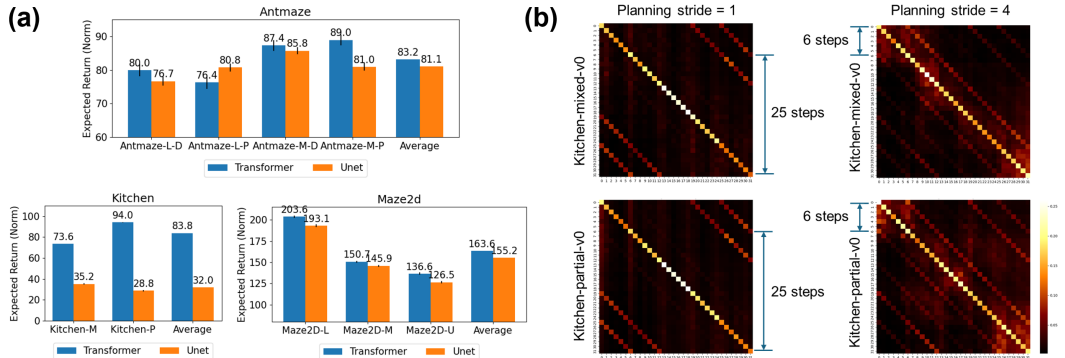


Figure 5: Using Transformer as the backbone of denoising network. (a) Performance comparison between Transformer and U-Net. The Transformer outperforms U-Net in 8 out of 9 sub-tasks and in all 3 main tasks. The amount of parameters in U-Net is comparable to that in Transformers. (b) Visualization of attention weights of the first layer in the Transformer network during the denoising process. More plots can be found in Appendix C.

Most diffusion planners on the D4RL dataset use 1-D U-Net for the denoising network. It is natural to question whether attention is all you need (Vaswani et al., 2017) for diffusion planning. Thus, we examined the benefit of replacing U-Net with the Transformer architecture as the backbone of the denoising model (Sect. 3.1) (see Appendix B for details about network structures). The experimental results clearly support the utilization of Transformer (Fig. 5(a)) in diffusion planning, consistent with the latest trend in image and video generation (Peebles & Xie, 2023; OpenAI, 2024).

We also conducted a case study by looking into the attention weights of the trained Transformer in the Kitchen environment (Fig. 5(b)), which reflect the temporal credit assignment (i.e., to how many steps later should be paid attention in the planning sequence). First, we see that the model pays more attention to the long-range element in the trajectory compared to the short-range ones. It suggests that the long-term dependency is crucial in this task, which breaks the local inductive bias of convolutional neural networks such as U-Net. Second, an interesting finding is that the characteristic attention length is consistent even with different planning stride (Sect. 4.2): 6 (attention step) $\times 4$ (stride) ≈ 25 (attention step) $\times 1$ (stride), as depicted in Fig. 5(b). It suggests that the Transformer finds the invariant correlations across the stride, contributing to the generalization performance. We have included more attention weights visualization in Appendix C, which also implies the importance of long-term dependency in trajectory modeling. In-depth study will be needed to fully understand the long-term dependency in the future.

4.4 IMPACT OF NETWORK SIZE

Since the experimental results are in favor of Transformer, one may wonder whether a “scaling law” (Kaplan et al., 2020) holds, in particular, whether performance scales up with model depth (Ye et al., 2024). The results presented in Fig. 6 pass two clear messages: First, 1-layer Transformer is not

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

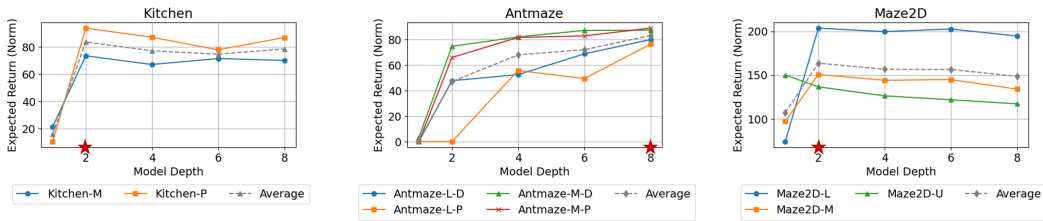


Figure 6: Performance change over depth of the Transformer network as diffusion planner. The star indicates the choice of DV.

enough, except for the simplest sub-task (Maze2D-U). Second, a deeper model is not always better. This may be due to an intrinsic difference between decision making and natural language processing and limitations of dataset size and quality, which requires further study to systematically address.

4.5 GUIDED SAMPLING ALGORITHMS

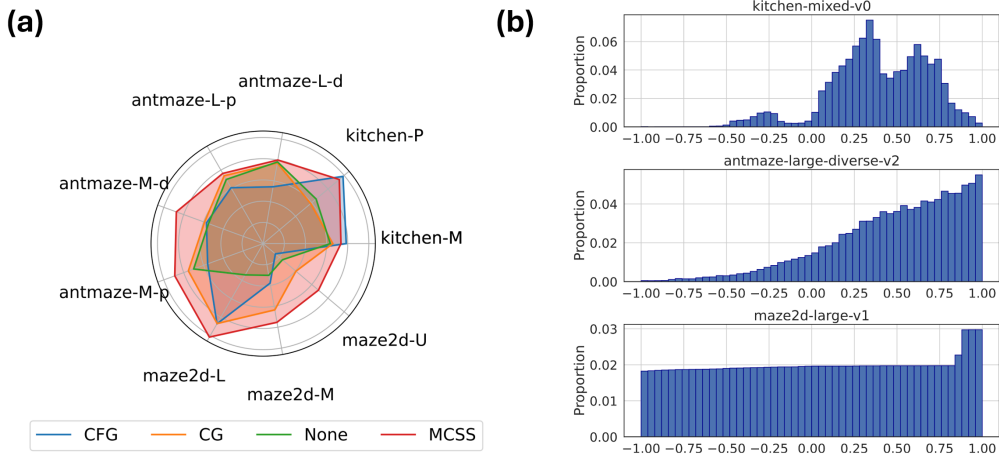


Figure 7: Analysis of guided sampling algorithm. (a) Performance comparison among different guided sampling algorithms for reward maximization. (b) Histogram of the value (accumulated discounted return in the future $(\sum_{h=0}^{\text{end}} \gamma^h r_{t+h})$, normalized to $[-1, 1]$) of the data points in each environment. For AntMaze, the failed trajectories are omitted since their values are all 0.

Another inconsistent design in previous work lies in the choice of guided sampling algorithm (Sect. 3.1), which enables the diffusion planner to generate plans that perform better than the average level of the dataset. Fig. 7(a) visualizes the corresponding empirical results (normalized) in our model. We can draw several conclusions from the results.

First, classifier guidance (CG) is comparable with classifier-free guidance (CFG), despite the fact that CFG is generally considered better than CG in image synthesis (Ho & Salimans, 2021). A potential reason is that the target value of CFG may need to be adjusted over time since the total rewards an agent can obtain in the future may vary depending on the task stage, but we can only use a fixed target value for CFG since there is no trivial solution.

Also, we observed that non-guidance can be better than guidance – Monte Carlo sampling with selection (MCSS) performs overall the best, except for Franka Kitchen where MCSS lags slightly behind CFG. This is an important finding since existing diffusion planners usually used CG or CFG (Chen et al., 2023; Wang et al., 2023b)). To understand the potential underlying reasons, we plotted the value distribution of data in each environment (Fig. 7(b)). It can be seen that in Maze2D and AntMaze, there is a substantial amount of optimal and near-optimal experiences, whereas in Kitchen

most samples are sub-optimal (note that here the optimality is with respect to condition of diffusion model). This may explain why CFG performs better than MCSS in Kitchen. Thus we can propose a hypothesis: No guidance (MCSS) can be better than guided generation (CG, CFG) if the dataset contains a substantial portion of expert demonstration.

4.6 COMPARISON TO DIFFUSION POLICY

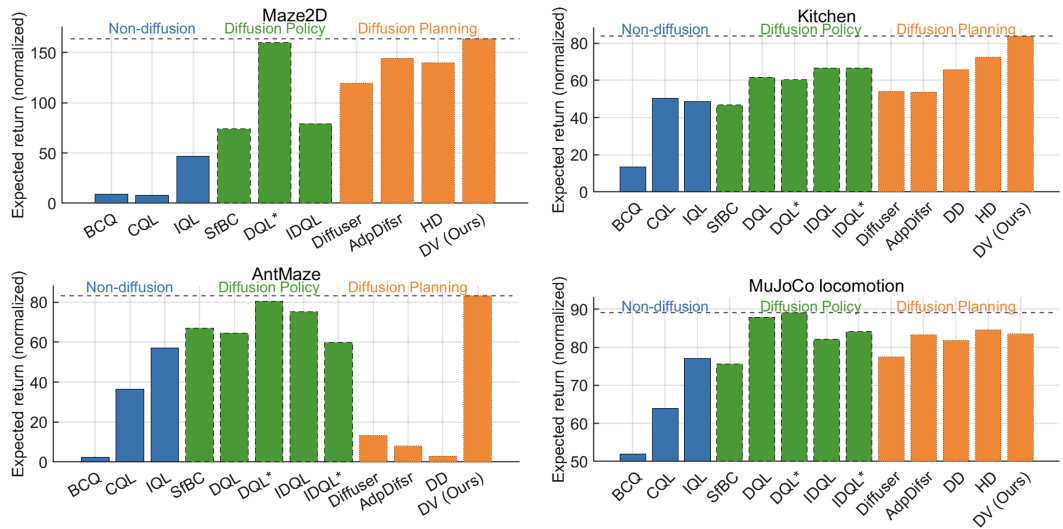


Figure 8: **Average performance of methods on different tasks.** The horizontal dashed line indicates the best performance over all methods. DV (Diffusion planning) shows the best performance for Kitchen, Maze2D, and AntMaze; while DQL (Diffusion policy) (Wang et al., 2023b) outperforms all diffusion planning methods in MuJoCo locomotion tasks. Refer to the caption of Table 1 for the reference of each method.

Diffusion planning and diffusion policy represent two key approaches within diffusion-based decision-making. After examining the core components of diffusion planners, we turn to a comparison of diffusion planning and diffusion policy across different environments. The experimental results are illustrated in Fig. 8. We observed that diffusion planning outperforms diffusion policy in AntMaze, Kitchen, and Maze2D, whereas diffusion policy excels in MuJoCo locomotion tasks. The first three environments require precise goal achievement, such as positioning an object exactly, necessitating long-term planning. This makes them well-suited to diffusion planning, which generates entire trajectories in one step. Furthermore, these environments feature sparse reward structures, posing challenges for model-free RL algorithms typically used in diffusion policies (Wang et al., 2023b). In contrast, the objective in MuJoCo is simply to control agents to run faster, a task that is less related to lookahead planning and does not require intricate planning. RL loss functions can help diffusion policy (Wang et al., 2023b) achieve better results in such scenarios.

4.7 PRACTICAL TIPS TO TAKE HOME

Takeaway 1: Diffusion planning is most effective for tasks requiring long-term credit assignment, while diffusion policies are better suited for locomotion tasks that do not demand long-term planning(Sect. 4.6)

Takeaway 2: It is recommended to generate state plans with diffusion planners and use an inverse dynamics model to compute the corresponding actions(Sect. 4.1).

Takeaway 3: Implementing jump-step planning can be highly beneficial; experimenting with different planning strides is encouraged(Sect. 4.2).

Takeaway 4: It is worth trying to use Transformer as the backbone of diffusion planner, especially in the tasks that require long-term lookahead planning (Sect. 4.3).

486 **Takeaway 5:** A single-layer Transformer is insufficient for effective planning (Sect. 4.4).

487
488 **Takeaway 6:** Larger models do not necessarily lead to better performance in diffusion planner for
489 offline RL (Sect. 4.4).

490 **Takeaway 7:** Non-guidance approaches, such as Monte Carlo unconditional sampling with selection,
491 can outperform classifier or classifier-free guidance when the dataset contains enough near-optimal
492 trajectories (Sect. 4.5).

493 494 495 496 5 DISCUSSIONS 497

498 **Synergy between diffusion planning and diffusion policy.** A significant avenue for future research
499 involves a deeper exploration of the distinctions between diffusion planning and diffusion policy.
500 Drawing on Daniel Kahneman’s seminal work *Thinking, Fast, and Slow* (Kahneman, 2011), human
501 cognitive processes are categorized into System 1 and System 2. Diffusion policies are analogous
502 to System 1 processes, as they operate rapidly and efficiently (Wang et al., 2023b), making them
503 well-suited for tasks such as locomotion (Fig. 8) that do not require extensive deliberation or long-
504 term planning. These policies manage routine decision making with the same efficiency as intuitive
505 responses in human cognition. Conversely, diffusion planning mirrors System 2 thinking, character-
506 ized by its slower, more deliberate, and effortful nature. This approach is particularly effective
507 for tasks that demand long-term credit assignment (Fig. 8), involving more computations to develop
508 effective plans. In RL terminology, diffusion planning can be broadly classified as model-based, while
509 diffusion policy aligns with model-free methodologies. Investigating the interplay between these two
510 systems presents a compelling intersection for both machine learning and cognitive neuroscience
511 (Gläscher et al., 2010; Duan et al., 2016; Botvinick et al., 2019). Studies from cognitive science
512 indicate that the brain may use a synergistic approach which arbitrates and selects the better system
513 according to the current situation, and the preference may change over time (Lee et al., 2014; Han
514 et al., 2024). We anticipate extensive future research focused on integrating the strengths of diffusion
515 planning and diffusion policies to enable both efficient and effective decision-making AI.

516 **Computational efficiency.** Despite the effectiveness of diffusion planners, their computational cost
517 is substantial. Our work is orthogonal to the optimization of computational cost (Dong et al., 2024a).
518 Nonetheless, future work may consider new schemes such as the consistency model (Song et al.,
519 2023) to improve computational efficiency.

520 **Interpretability and safety.** Our study focuses on a single performance metric (total return),
521 potentially overlooking qualitative aspects such as the interpretability and reliability of the diffusion
522 planner. Future work may consider issues such as explainability (Puiutta & Veith, 2020) and safety
523 (Xiao et al., 2023) of diffusion planning. Leveraging the experiences from computer vision domain
524 will be worth investigating.

525 **Sustainability.** Our work required significant computational resources, particularly in terms of
526 GPU energy consumption, as we trained and evaluated thousands of models across diverse tasks.
527 However, this investment in energy is not without purpose. We aim to provide a solid foundation
528 for future research. Subsequent work can build upon our findings, reducing the need for extensive
529 trial-and-error experimentation. In this way, our research contributes to energy efficiency in the
530 long term, as researchers can reference our results and apply proven methods rather than duplicating
531 resource-intensive exploratory efforts.

532 **Open problems and future directions.** In the current study, we have focused on standard Markov
533 decision process problems (Bellman, 1957) using a popular offline RL benchmark (Fu et al., 2020).
534 The planning and control are based on joint states and coordinates. Numerous untouched problems
535 exist, such as vision-based decision making (Du et al., 2024; Yang et al., 2024), goal-conditioned
536 reinforcement learning (Liu et al., 2022; Wang et al., 2023a), partially observable environments
537 (Schmidhuber, 1991), offline-to-online deployment (Matsushima et al., 2021), and the scalability of
538 diffusion planning models (Kaplan et al., 2020). Future efforts are anticipated to fully address these
539 limitations. However, even within the scope of the current work, we have found several interesting
phenomena and tips that are counter to common practices. Our work should be considered as a new
but solid starting point for behavior planning using decision models.

REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our results. To facilitate this, we include the source code of DV in supplementary material and will release the code upon acceptance. Detailed descriptions of our experimental setup, including model architectures, training procedures, and hyperparameter settings, are provided in Appendix A and Appendix B. We have included comprehensive information on the datasets used, along with any preprocessing steps, in Appendix B. For all key experiments, we have specified the evaluation protocols and metrics in Sect. 3 and provided extensive results in Appendix C. We have included the full list of hyperparameters and configurations in Appendix B.4.

REFERENCES

- Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B Tenenbaum, Tommi S Jaakkola, and Pulkit Agrawal. Is conditional generative modeling all you need for decision making? In *The Eleventh International Conference on Learning Representations*, 2022.
- Richard Bellman. A Markovian decision process. *Journal of Mathematics and Mechanics*, pp. 679–684, 1957.
- Matthew Botvinick, Sam Ritter, Jane X Wang, Zeb Kurth-Nelson, Charles Blundell, and Demis Hassabis. Reinforcement learning, fast and slow. *Trends in cognitive sciences*, 23(5):408–422, 2019.
- Chang Chen, Fei Deng, Kenji Kawaguchi, Caglar Gulcehre, and Sungjin Ahn. Simple hierarchical planning with diffusion. *arXiv preprint arXiv:2401.02644*, 2024.
- Huayu Chen, Cheng Lu, Chengyang Ying, Hang Su, and Jun Zhu. Offline reinforcement learning via high-fidelity generative behavior modeling. In *The Eleventh International Conference on Learning Representations*, 2023.
- Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(9): 10850–10869, 2023.
- Yilun Dai, Mengjiao Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *arXiv preprint arXiv:2302.00111*, 2023.
- Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the International Conference on Machine Learning*, pp. 465–472, 2011.
- Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- Zibin Dong, Jianye Hao, Yifu Yuan, Fei Ni, Yitian Wang, Pengyi Li, and Yan Zheng. Diffuserlite: Towards real-time diffusion planning. *arXiv preprint arXiv:2401.15443*, 2024a.
- Zibin Dong, Yifu Yuan, Jianye Hao, Fei Ni, Yi Ma, Pengyi Li, and Yan Zheng. Cleandiffuser: An easy-to-use modularized library for diffusion models in decision making. *arXiv preprint arXiv:2406.09509*, 2024b.
- Zibin Dong, Yifu Yuan, Jianye HAO, Fei Ni, Yao Mu, YAN ZHENG, Yujing Hu, Tangjie Lv, Changjie Fan, and Zhipeng Hu. Aligndiff: Aligning diverse human preferences via behavior-customisable diffusion model. In *The Twelfth International Conference on Learning Representations*, 2024c.
- Yilun Du, Sherry Yang, Bo Dai, Hanjun Dai, Ofir Nachum, Josh Tenenbaum, Dale Schuurmans, and Pieter Abbeel. Learning universal policies via text-guided video generation. *Advances in Neural Information Processing Systems*, 36, 2024.
- Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel. RL2: Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.

- 594 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep
595 data-driven reinforcement learning, 2020.
596
- 597 Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning.
598 *Advances in neural information processing systems*, 34:20132–20145, 2021.
599
- 600 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
601 exploration. In *International Conference on Machine Learning*, pp. 2052–2062. PMLR, 2019.
602
- 603 Jan Gläscher, Nathaniel Daw, Peter Dayan, and John P O’Doherty. States versus rewards: dissociable
604 neural prediction error signals underlying model-based and model-free reinforcement learning.
605 *Neuron*, 66(4):585–595, 2010.
606
- 607 Dongqi Han, Kenji Doya, Dongsheng Li, and Jun Tani. Synergizing habits and goals with variational
608 bayes. *Nature Communications*, 15(1):4461, 2024.
609
- 610 Philippe Hansen-Estruch, Ilya Kostrikov, Michael Janner, Jakub Grudzien Kuba, and Sergey Levine.
611 Idql: Implicit q-learning as an actor-critic method with diffusion policies. *arXiv preprint*
612 *arXiv:2304.10573*, 2023.
613
- 614 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on*
615 *Deep Generative Models and Downstream Applications*, 2021.
616
- 617 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
618 *neural information processing systems*, 33:6840–6851, 2020.
619
- 620 Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence
621 modeling problem. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
622
- 623 Michael Janner, Yilun Du, Joshua Tenenbaum, and Sergey Levine. Planning with diffusion for
624 flexible behavior synthesis. In *International Conference on Machine Learning*, pp. 9902–9915.
625 PMLR, 2022.
626
- 627 Daniel Kahneman. *Thinking, fast and slow*. macmillan, 2011.
628
- 629 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott
630 Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models.
631 *arXiv preprint arXiv:2001.08361*, 2020.
632
- 633 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
634 *arXiv:1412.6980*, 2014.
635
- 636 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit
637 q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
638
- 639 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
640 reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
641
- 642 Sang Wan Lee, Shinsuke Shimojo, and John P O’Doherty. Neural computations underlying arbitration
643 between model-based and model-free learning. *Neuron*, 81(3):687–699, 2014.
644
- 645 Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial,
646 review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
647
- 648 Wenhao Li, Xiangfeng Wang, Bo Jin, and Hongyuan Zha. Hierarchical diffusion for offline decision
649 making. In *International Conference on Machine Learning*, pp. 20035–20064. PMLR, 2023.
650
- 651 Zhixuan Liang, Yao Mu, Mingyu Ding, Fei Ni, Masayoshi Tomizuka, and Ping Luo. Adaptdiffuser:
652 Diffusion models as adaptive self-evolving planners. *arXiv preprint arXiv:2302.01877*, 2023.
653
- 654 Minghuan Liu, Menghui Zhu, and Weinan Zhang. Goal-conditioned reinforcement learning: Problems
655 and solutions. *arXiv preprint arXiv:2201.08299*, 2022.

- 648 Cheng Lu, Huayu Chen, Jianfei Chen, Hang Su, Chongxuan Li, and Jun Zhu. Contrastive energy
649 prediction for exact energy-guided diffusion sampling in offline reinforcement learning. In
650 *International Conference on Machine Learning*, pp. 22825–22855. PMLR, 2023.
- 651
- 652 Tatsuya Matsushima, Hiroki Furuta, Yutaka Matsuo, Ofir Nachum, and Shixiang Gu. Deployment-
653 efficient reinforcement learning via model-based offline optimization. In *International Conference*
654 *on Learning Representations*, 2021.
- 655 John D Murray, Alberto Bernacchia, David J Freedman, Ranulfo Romo, Jonathan D Wallis, Xinying
656 Cai, Camillo Padoa-Schioppa, Tatiana Pasternak, Hyojung Seo, Daeyeol Lee, et al. A hierarchy of
657 intrinsic timescales across primate cortex. *Nature Neuroscience*, 17(12):1661, 2014.
- 658
- 659 OpenAI. Sora. <https://openai.com/index/sora/>, 2024.
- 660 Paavo Parmas, Carl Edward Rasmussen, Jan Peters, and Kenji Doya. PIPPS: Flexible model-based
661 policy search robust to the curse of chaos. In *International Conference on Machine Learning*, pp.
662 4065–4074. PMLR, 2018.
- 663
- 664 Tim Pearce, Tabish Rashid, Anssi Kanervisto, Dave Bignell, Mingfei Sun, Raluca Georgescu,
665 Sergio Valcarcel Macua, Shan Zheng Tan, Ida Momennejad, Katja Hofmann, et al. Imitating
666 human behaviour with diffusion models. *arXiv preprint arXiv:2301.10677*, 2023.
- 667 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
668 *the IEEE/CVF International Conference on Computer Vision*, pp. 4195–4205, 2023.
- 669
- 670 Erika Puiutta and Eric MSP Veith. Explainable reinforcement learning: A survey. In *International*
671 *cross-domain conference for machine learning and knowledge extraction*, pp. 77–95. Springer,
672 2020.
- 673 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical
674 image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI*
675 *2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III*
676 *18*, pp. 234–241. Springer, 2015.
- 677
- 678 Caroline A Runyan, Eugenio Piasini, Stefano Panzeri, and Christopher D Harvey. Distinct timescales
679 of population coding across cortex. *Nature*, 548(7665):92, 2017.
- 680 Jürgen Schmidhuber. Reinforcement learning in Markovian and non-Markovian environments. In
681 *Advances in Neural Information Processing Systems*, pp. 500–506, 1991.
- 682
- 683 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*
684 *preprint arXiv:2010.02502*, 2020.
- 685
- 686 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *International*
687 *Conference on Machine Learning*, pp. 32211–32252. PMLR, 2023.
- 688
- 689 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT
690 press Cambridge, 1998.
- 691
- 692 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
693 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information*
694 *Processing Systems*, pp. 5998–6008, 2017.
- 695
- 696 Jane X Wang, Zeb Kurth-Nelson, Dharshan Kumaran, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo,
697 Demis Hassabis, and Matthew Botvinick. Prefrontal cortex as a meta-reinforcement learning
698 system. *Nature Neuroscience*, 21(6):860, 2018.
- 699
- 700 Wei Wang, Dongqi Han, Xufang Luo, Yifei Shen, Charles Ling, Boyu Wang, and Dongsheng Li.
701 Toward open-ended embodied tasks solving. In *Second Agent Learning in Open-Endedness*
Workshop, 2023a.
- 702
- 703 Zhendong Wang, Jonathan J Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy
704 class for offline reinforcement learning. In *The Eleventh International Conference on Learning*
Representations, 2023b.

702 Wei Xiao, Tsun-Hsuan Wang, Chuang Gan, and Daniela Rus. Safediffuser: Safe planning with
703 diffusion probabilistic models. *arXiv preprint arXiv:2306.00148*, 2023.
704

705 Cheng-Fu Yang, Haoyang Xu, Te-Lin Wu, Xiaofeng Gao, Kai-Wei Chang, and Feng Gao. Planning
706 as in-painting: A diffusion-based embodied task planning framework for environments under
707 uncertainty. *arXiv preprint arXiv:2312.01097*, 2023.

708 Sherry Yang, Yilun Du, Seyed Kamyar Seyed Ghasemipour, Jonathan Tompson, Leslie Pack Kael-
709 bling, Dale Schuurmans, and Pieter Abbeel. Learning interactive real-world simulators. In *The*
710 *Twelfth International Conference on Learning Representations*, 2024.
711

712 Tian Ye, Zicheng Xu, Yuanzhi Li, and Zeyuan Allen-Zhu. Physics of language models: Part 2.1,
713 grade-school math and the hidden reasoning process. *arXiv preprint arXiv:2407.20311*, 2024.

714 Mingyuan Zhang, Zhongang Cai, Liang Pan, Fangzhou Hong, Xinying Guo, Lei Yang, and Ziwei
715 Liu. Motiondiffuse: Text-driven human motion generation with diffusion model. *arXiv preprint*
716 *arXiv:2208.15001*, 2022.

717 Zhengbang Zhu, Hanye Zhao, Haoran He, Yichao Zhong, Shenyu Zhang, Yong Yu, and Weinan
718 Zhang. Diffusion models for reinforcement learning: A survey. *arXiv preprint arXiv:2311.01223*,
719 2023.
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A GUIDED SAMPLING ALGORITHMS

For decision making tasks, guided sampling algorithms are used to generate desired plans or actions. In this work, we compare three types of different guided sampling methods: classifier guidance (Dhariwal & Nichol, 2021) (CG), classifier-free guidance (CFG) (Ho & Salimans, 2021), and Monte Carlo sampling from selections (MCSS).

Classifier guidance: Classifier guidance (CG) is introduced to guide the unconditional diffusion models $q_t(\mathbf{x}_t)$ to generate data over condition c . The conditioned score function is formulated as:

$$\nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t|c) = \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t) + \nabla_{\mathbf{x}} \log q_t(c|\mathbf{x}_t)$$

where the second term is also known as a noised classifier that predicts the condition using noised data \mathbf{x}_t . During sampling, the gradient of the classifier is then applied to the predicted noise $\epsilon_{\theta}(\mathbf{x}_t, t)$:

$$\bar{\epsilon}_{\theta}(\mathbf{x}_t, t, c) = \epsilon_{\theta}(\mathbf{x}_t, t) - w\sigma_t \nabla_{\mathbf{x}} \log q_t(c|\mathbf{x}_t)$$

where w is a weighting factor that controls the strength of the classifier guidance. For CG sampling, we tuned w in range $[0.001, 10]$ on each task.

Classifier-free guidance: To avoid training classifiers, classifier-free guidance (CFG) is proposed. The main idea of CFG is to train a diffusion model that can be used for both conditional noise predictor $\epsilon_{\theta}(\mathbf{x}_t, t, c)$ and unconditional noise predictor $\epsilon_{\theta}(\mathbf{x}_t, t)$:

$$\bar{\epsilon}_{\theta}(\mathbf{x}_t, t, c) = \epsilon_{\theta}(\mathbf{x}_t, t) + w(\epsilon_{\theta}(\mathbf{x}_t, t, c) - \epsilon_{\theta}(\mathbf{x}_t, t))$$

where $\epsilon_{\theta}(\mathbf{x}_t, t) = \epsilon_{\theta}(\mathbf{x}_t, t, \emptyset)$. Noise prediction of $\epsilon_{\theta}(\mathbf{x}_t, t, \emptyset)$ and $\epsilon_{\theta}(\mathbf{x}_t, t, c)$ can be jointly learnt by randomly discarding conditioning with probability of p_{uncond} . For decision making tasks, we can train diffusion models using condition of discounted returns, and using classifier-free guidance for better plan sampling. We can normalize the discounted return in the dataset for training, and use condition of 1 as target return for CFG sampling during inference (Ajay et al., 2022). However, experiments show that fixing 1 as target may lead to unrealistic or unstable plans. Consequently, besides tuning the guidance strength $w \in [1.0, 6.0]$, we also tune for the best target return within range $[0.5, 1.5]$ for each task to test CFG’s best performance.

Monte Carlo sampling from selections: For Monte Carlo sampling from selections (MCSS), N selections are firstly sampled from an unconditional generative model as candidates. Then these candidates are evaluated with a learnt critic for the selection of the optimal one. One advantage for MCSS is that, it does not rely on any task-specific hyperparameters for inference, such as the guidance strength w in CFG and CG, and target return in CFG. However, it needs to sample $N - 1$ more candidates during each decision making step, and requires more computation.

B IMPLEMENTATION DETAILS

B.1 MODEL ARCHITECTURE

Planner: Our code is based on CleanDiffuser (Dong et al., 2024b). We examined U-Net (Ronneberger et al., 2015) and Transformer (Vaswani et al., 2017) as the neural network backbone for all the diffusion planners. Specifically, we keep consistent with the implementation of U-Net1D (Janner et al., 2022), with 5 kernel size, $(1, 2, 2, 2)$ for channel multiplication, 32 base channels on MuJoCo, Kitchen and Maze2D, and 64 base channels on AntMaze. For Transformer, we use DiT1D (Peebles & Xie, 2023; Dong et al., 2024c), with hidden dimension of 256, head dimension of 32, 2 DiT blocks on MuJoCo, Kitchen and Maze2D, and 8 DiT blocks on AntMaze. All the planner diffusion models are trained with the Adam (Kingma & Ba, 2014) optimizer with learning rate of $3e - 4$, batch size of 128, for 1M gradient steps. All the diffusion models in this work are trained to predict the noise. However, for U-Net1D experiments on Kitchen, the diffusion planner to predict the clean estimation, because it could achieve quite better performance.

Inverse dynamics: We used an MLP-based diffusion model as the inverse dynamic model, whose input is the current state and the planned next-state; and the output is the action to execute. It is implemented with a 3-layer MLP with additional 2-layer embedding layer and trained with 1M gradient steps. All the inverse dynamic models are trained with the Adam (Kingma & Ba, 2014)

optimizer with learning rate of $3e - 4$, batch size of 128. We found that the using diffusion model as inverse dynamics show similar performance with vanilla MLP inverse dynamics.

Critic: We also implemented two types of critic models for guided sampling. The first type has the architecture of the U-Net1D for the planner, with a linear output layer to produce the critic value. The second type is a 2 blocks vanilla transformer, with hidden dimension of 256 as the value function, with a linear projection head on the first token output. We trained all the critic model using the Adam (Kingma & Ba, 2014) optimizer with learning rate of $3e - 4$, batch size of 128. The model will be trained with 200K gradient steps, if it is a clean critic model ¹. Otherwise, it is trained for 1M gradient steps.

B.2 DIFFUSION SOLVER

We use DDIM (Song et al., 2020) of temperature 1.0 for planner diffusion sampling, and DDPM with temperature of 0.5 for inverse dynamic action sampling. The sampling temperature is introduced to reduce sampling randomness (Ajay et al., 2022).

B.3 DATASET PRE-PROCESSING

Diffusion policy baselines (Chen et al., 2023; Wang et al., 2023b; Hansen-Estruch et al., 2023) commonly learns policy, Q functions, and value functions using a temporal difference manner, on standard transitions (s_t, a_t, s_{t+1}, r_t) . Admittedly, diffusion planners often require careful dataset pre-processing, including horizon padding, planning strides, return calculation, and truncation-termination handling. An unsuccessful sequential dataset pre-processing may greatly reduce the planning ability of diffusion planners. Most planning tasks are usually sparse rewarded, and how optimality is defined, combined with different temporal credit assignment methods is also important. For MuJoCo and Kitchen, we use discount factor $\gamma = 0.997$. For Maze2D and AntMaze, we use IQL-maze (Kostrikov et al., 2021) reward shaping methods for temporal credit assignment in navigation planning tasks, where an -1 penalty is always applied to agent during every timestep.

B.4 FULL HYPER-PARAMETERS

Table 2 displays the hyperparameters and default choices in our work.

¹A clean critic model is only trained on clean input data, while a noised critic model is trained using noised data using the diffusion model’s noise schedule

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

Table 2: Configuration Settings

Settings	Default	Choices
Guidance Type	MCSS	[MCSS, CG, CFG, None]
State-Action Generation	Separate	[Joint, Separate]
Advantage Weighting	True	[True, False]
Inverse Dynamic	Diffusion	[Diffusion, Regular]
Time Credit Assignment	discount=0.997	[discount=0.997, IQL-maze]
Planner Net. Backbone	Transformer	[Transformer, UNet]
UNet Channels Mult	(1, 2, 2, 2)	(1, 2, 2, 2)
UNet Base Channels	32	[16, 32, 64]
Transformer Hidden	256	256
Transformer Block	2	[2, 4, 6, 8]
Planner Solver	DDIM	[DDIM, DDPM]
Planner Sampling Steps	20	20
Planner Training Steps	1000000	1000000
Planner Temperature	1	1
MCSS Candidates	50	[1, 20, 50]
Planning Horizon	32	[4, 32, 40]
Planning Stride	1	[1, 2, 4, 5, 15, 25]
Inverse Dynamics Net. Backbone	MLP	MLP
Inverse Dynamics Hidden	256	256
Inverse Dynamics Solver	DDPM	DDPM
Inverse Dynamics Sampling Steps	10	10
Inverse Dynamics Training Steps	1000000	1000000
Policy Temperature	0.5	0.5

C EXTENSIVE RESULTS

The following plots show the attention weights in different DDIM denoising steps in each task. We can see that a long-term dependency is generally existing in the Transformer.

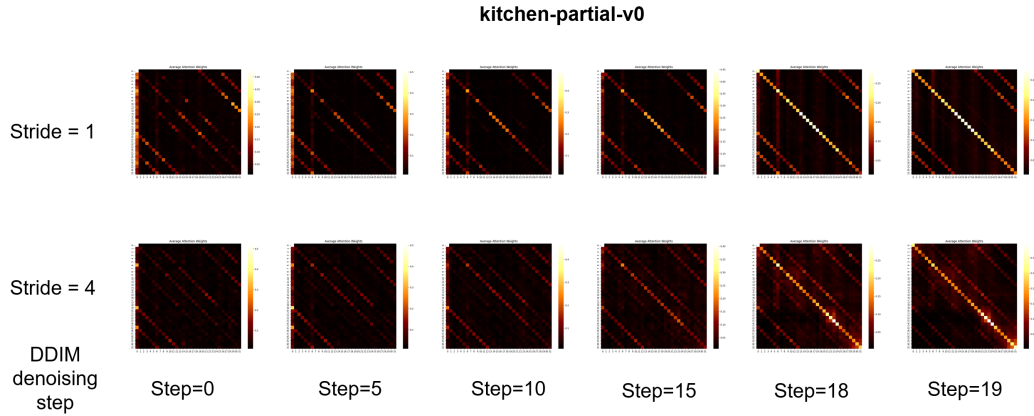


Figure 9: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the Kitchen-Partial-v0 dataset.

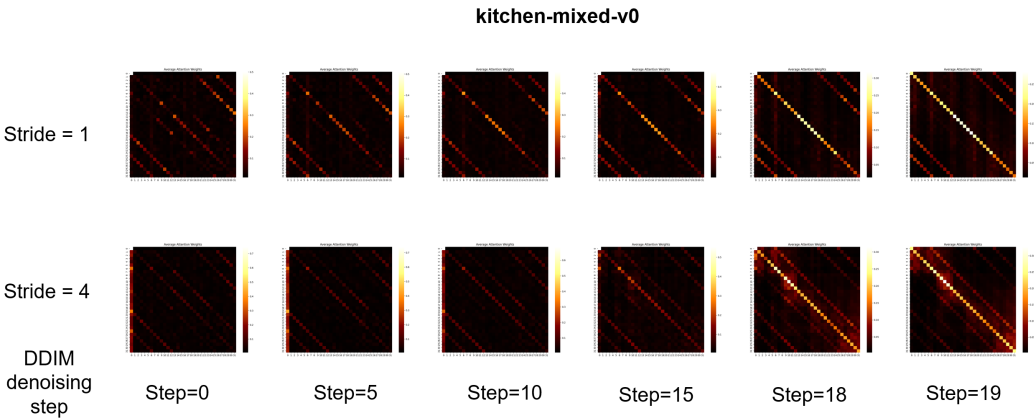


Figure 10: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the Kitchen-Mixed-v0 dataset.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

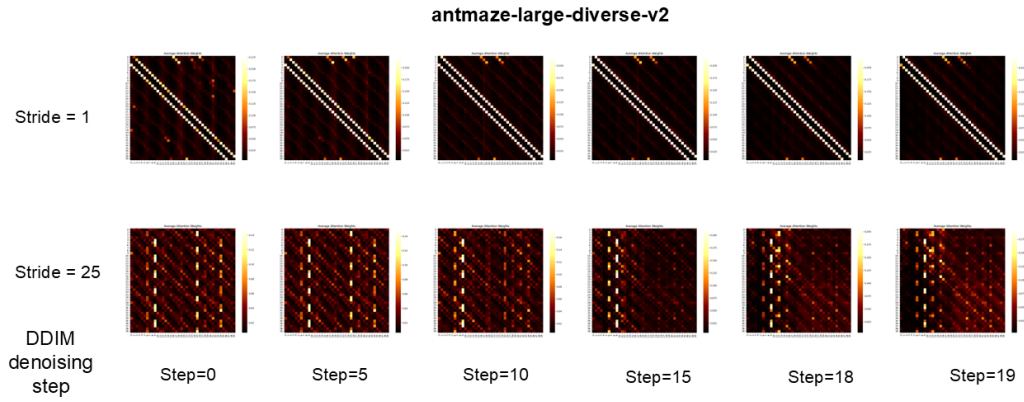


Figure 11: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the AntMaze-Large-Diverse-v2 dataset.

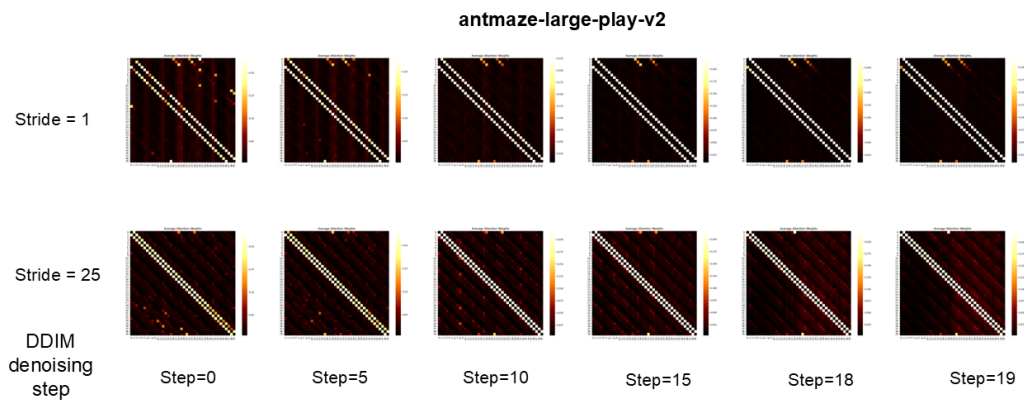


Figure 12: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the AntMaze-Large-Play-v2 dataset.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

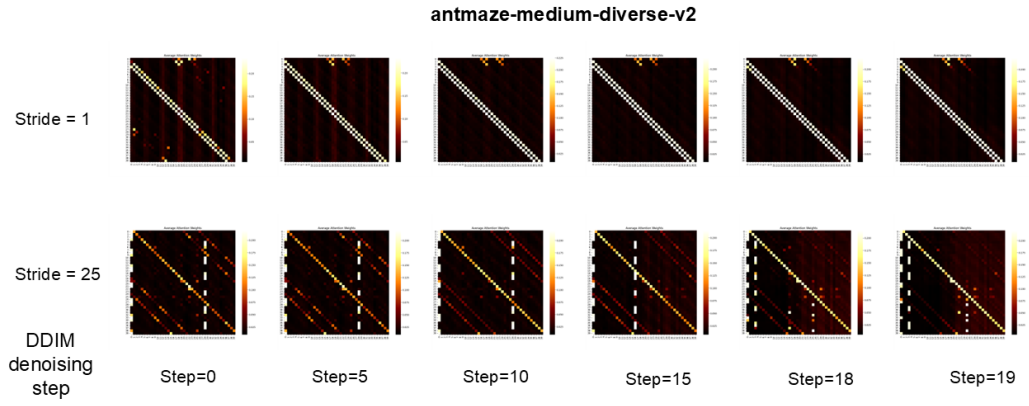


Figure 13: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the AntMaze-Medium-Diverse-v2 dataset.

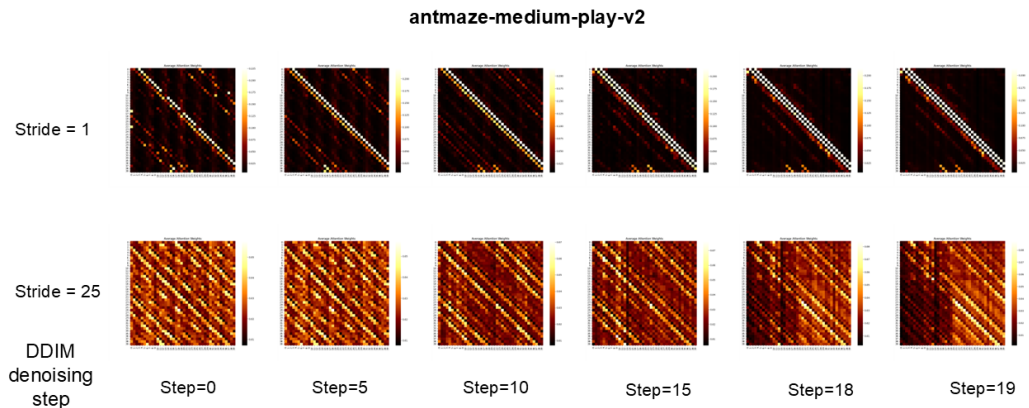


Figure 14: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the AntMaze-Medium-Play-v2 dataset.

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

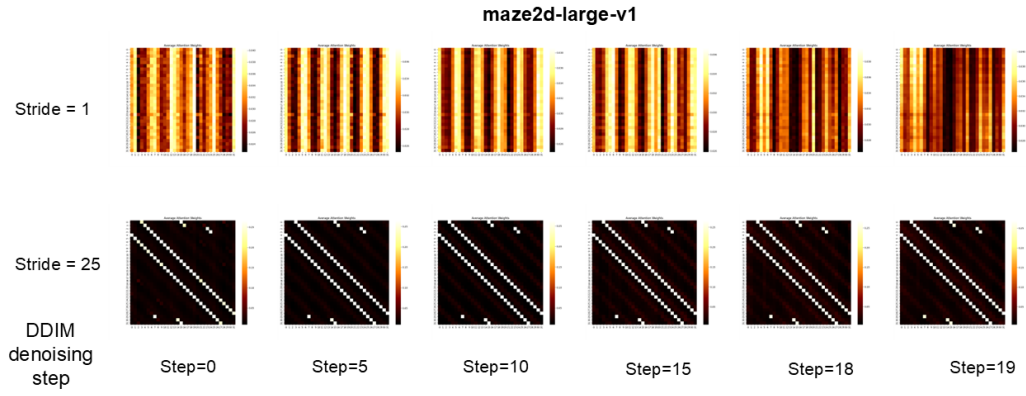


Figure 15: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the Kitchen-Mixed-v0 dataset.

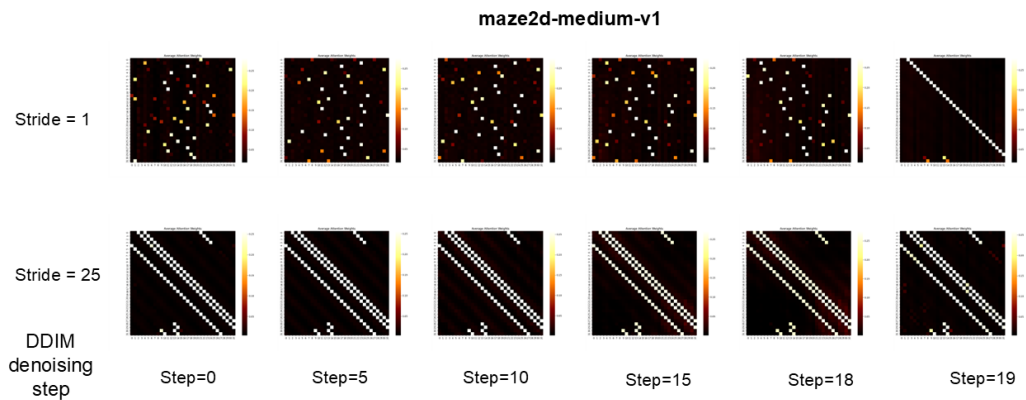


Figure 16: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the Maze2D-Large-v1 dataset.

1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172
 1173
 1174
 1175
 1176
 1177
 1178
 1179
 1180
 1181
 1182
 1183
 1184
 1185
 1186
 1187

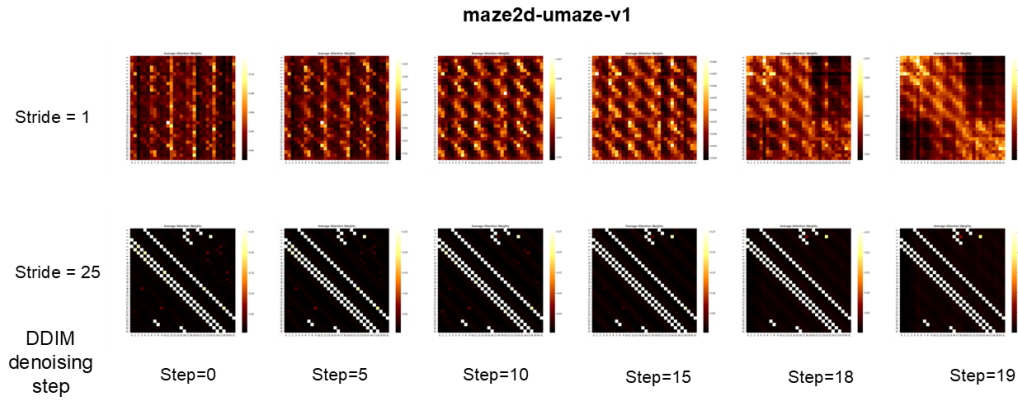


Figure 17: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the Maze2D-medium-v1 dataset.

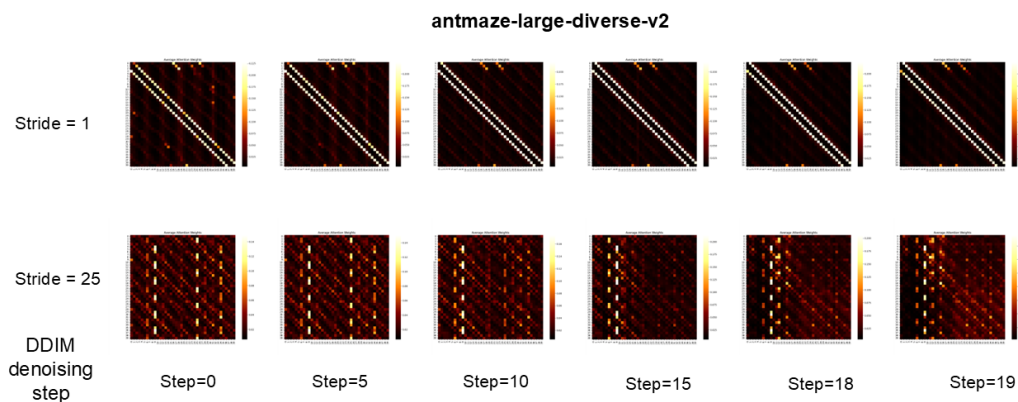


Figure 18: Attention weights (averaged on multi-heads) of the first Transformer layer of DV on the Maze2D-UMaze-v1 dataset.

Table 3: Normalized performance of various offline-RL methods. Data are Mean \pm Standard Error (if available).

Dataset	Environment	IL														
		None-diffusion Policies					Diffusion Policies									
		BC	BCQ	CQL	IQL	SBBC	DQL	DQL*	IDQL	IDQL*	CEP	Diffuser	AdaptDiffuser	DD	HD	DV
Medium-Expert	HalfCheetah	35.8	64.7	62.4	86.7	92.6 \pm 0.5	96.8 \pm 0.3	95.5 \pm 0.1	95.9	91.3 \pm 0.6	93.5 \pm 0.3	88.9 \pm 0.3	89.6 \pm 0.8	90.6 \pm 1.3	92.5 \pm 0.3	92.7 \pm 0.3
Medium-Replay	HalfCheetah	38.4	38.2	46.2	44.2	37.1 \pm 1.7	47.8 \pm 0.3	47.9 \pm 0.0	45.9	46.5 \pm 0.3	47.6 \pm 1.4	37.7 \pm 0.5	38.3 \pm 0.9	39.3 \pm 4.1	115.3 \pm 1.1	45.8 \pm 0.1
Medium	HalfCheetah	36.1	40.7	44.4	47.4	45.9 \pm 2.2	51.1 \pm 0.5	52.3 \pm 0.2	51.0	51.5 \pm 0.1	54.1 \pm 0.4	42.8 \pm 0.3	44.2 \pm 0.6	49.1 \pm 1.0	107.1 \pm 1.1	50.4 \pm 0.0
Medium-Expert	Hopper	111.9	110.9	98.7	91.5	108.6 \pm 2.1	111.1 \pm 1.3	111.1 \pm 0.4	108.6	110.1 \pm 0.7	108.0 \pm 2.5	103.3 \pm 1.3	111.6 \pm 2.0	111.8 \pm 1.8	46.7 \pm 0.2	110.0 \pm 0.5
Medium-Replay	Hopper	11.3	33.1	48.6	94.7	86.2 \pm 9.1	101.3 \pm 0.6	101.6 \pm 0.0	92.1	99.4 \pm 0.1	96.9 \pm 2.6	93.6 \pm 0.4	92.2 \pm 1.5	100.0 \pm 0.7	99.3 \pm 0.3	91.9 \pm 0.0
Medium	Hopper	29.0	54.5	58.0	66.3	57.1 \pm 4.1	90.5 \pm 4.6	96.5 \pm 1.3	65.4	70.1 \pm 2.0	98.0 \pm 2.6	74.3 \pm 1.4	96.6 \pm 2.7	79.3 \pm 3.6	84.0 \pm 0.6	83.6 \pm 1.2
Medium-Expert	Walker2d	6.4	57.5	111.0	109.6	109.8 \pm 0.2	110.1 \pm 0.3	111.6 \pm 0.0	112.7	110.6 \pm 0.0	110.7 \pm 0.6	106.9 \pm 0.2	108.2 \pm 0.8	108.8 \pm 1.7	38.1 \pm 0.7	109.2 \pm 0.0
Medium-Replay	Walker2d	11.8	15.0	26.7	73.9	65.1 \pm 5.6	95.5 \pm 1.5	98.2 \pm 0.1	85.1	89.1 \pm 2.4	84.4 \pm 4.1	70.6 \pm 1.6	84.7 \pm 3.1	75.0 \pm 4.3	94.7 \pm 0.7	85.0 \pm 0.5
Medium	Walker2d	6.6	53.1	79.2	78.3	77.9 \pm 2.5	87.0 \pm 0.9	86.8 \pm 0.2	82.5	88.1 \pm 0.4	86.0 \pm 0.7	79.6 \pm 0.6	84.4 \pm 2.6	82.5 \pm 1.4	84.1 \pm 2.2	82.8 \pm 0.1
	Average	31.9	52.0	63.9	77.0	75.6	87.9	89.1	82.1	84.1	86.6	77.5	83.3	81.8	84.6	83.5
Mixed	Kitchen	47.5	8.1	51.0	51.0	45.4 \pm 1.6	62.6 \pm 5.1	55.1 \pm 1.58	66.5	66.5 \pm 4.1	-	52.5 \pm 2.5	51.8 \pm 0.8	75.0 \pm 0.0	71.7 \pm 2.7	73.6 \pm 0.1
Partial	Kitchen	33.8	18.9	49.8	46.3	47.9 \pm 4.1	60.5 \pm 6.9	65.5 \pm 1.38	66.7	66.7 \pm 2.5	-	55.7 \pm 1.3	55.5 \pm 0.4	56.5 \pm 5.8	73.3 \pm 1.4	94.0 \pm 0.3
	Average	40.7	13.5	50.4	48.7	46.7	61.6	60.3	66.6	66.6	-	54.1	53.7	65.8	72.5	83.8
Antmaze-Large	Diverse	0.0	2.2	61.2	47.5	45.5 \pm 6.6	56.6 \pm 7.6	70.6 \pm 3.7	67.9	40.0 \pm 11.4	64.8 \pm 5.5	27.3 \pm 2.4	8.7 \pm 2.5	0.0 \pm 0.0	83.6 \pm 5.8	80.0 \pm 1.8
Antmaze-Large	Play	0.0	6.7	53.7	39.6	59.3 \pm 14.3	46.4 \pm 8.3	81.3 \pm 3.1	63.5	48.7 \pm 4.7	66.6 \pm 9.8	17.3 \pm 1.9	5.3 \pm 3.4	0.0 \pm 0.0	-	76.4 \pm 2.0
Antmaze-Medium	Diverse	0.0	0.0	15.8	70.0	82.0 \pm 3.1	78.6 \pm 10.3	82.6 \pm 3.0	84.8	83.3 \pm 5.0	83.8 \pm 3.5	2.0 \pm 1.6	6.0 \pm 3.3	4.0 \pm 2.8	88.7 \pm 8.1	87.4 \pm 1.6
Antmaze-Medium	Play	0.0	0.0	14.9	71.2	81.3 \pm 2.6	76.6 \pm 10.8	87.3 \pm 2.7	84.5	67.3 \pm 5.7	83.6 \pm 4.4	6.7 \pm 5.7	12.0 \pm 7.5	8.0 \pm 4.3	-	89.0 \pm 1.6
	Average	0.0	2.2	36.4	57.1	67.0	64.6	80.5	75.2	59.8	74.7	13.3	8.0	3.0	-	83.2
Maze2D	Large	5	6.2	12.5	58.6	74.4 \pm 1.7	-	186.8 \pm 1.7	90.1	-	-	123	167.9 \pm 5.0	-	128.4 \pm 3.6	203.6 \pm 1.4
Maze2D	Medium	30.3	8.3	5.0	34.9	73.8 \pm 2.9	-	152.0 \pm 0.8	89.5	-	-	121.5	129.9 \pm 4.6	-	135.6 \pm 3.0	150.7 \pm 1.0
Maze2D	Umaze	3.8	12.8	5.7	47.4	73.9 \pm 6.6	-	140.6 \pm 1.0	57.9	-	-	113.9	135.1 \pm 5.8	-	155.8 \pm 2.5	136.6 \pm 1.3
	Average	13.0	9.1	7.7	47.0	74.0	-	159.8	79.2	-	-	119.5	144.3	-	139.9	163.6

Table 4: The effect of guided sampling algorithms for DV, with different planning strides. In "Stride x/y", x is for Kitchen, and y is for Maze2D & AntMaze. Data are Mean \pm Standard Error over 500 episode seeds.

		MCSS					CFG				
Dataset	Environment	Stride 1	Stride 2/5	Stride 4/15	Stride 8/25	Stride 1	Stride 2/5	Stride 4/15	Stride 8/25		
Mixed	Kitchen	67.2 \pm 0.3	74.2 \pm 0.1	73.6 \pm 0.1	-	72.7 \pm 0.2	78.6 \pm 0.1	72.4 \pm 0.1	-		
Partial	Kitchen	87.8 \pm 0.6	88.3 \pm 0.4	94.0 \pm 0.3	-	95.9 \pm 0.4	98.6 \pm 0.3	95.2 \pm 0.6	-		
Average		77.5	81.3	83.8	-	84.3	88.6	83.8	-		
Antmaze-Large	Diverse	72.4 \pm 1.4	2.6 \pm 0.1	71.2 \pm 2.0	80.0 \pm 1.8	54.4 \pm 0.6	46.5 \pm 0.1	6.4 \pm 0.1	67.4 \pm 0.3		
Antmaze-Large	Play	61.2 \pm 1.5	18.2 \pm 0.5	76.2 \pm 0.9	76.4 \pm 2.0	60.8 \pm 0.8	0.4 \pm 0.0	59.4 \pm 0.3	56.4 \pm 0.5		
Antmaze-Medium	Diverse	58.8 \pm 1.5	24.2 \pm 0.8	86.2 \pm 1.6	87.4 \pm 1.6	57.1 \pm 0.6	5.0 \pm 0.1	33.0 \pm 0.8	49.0 \pm 0.7		
Antmaze-Medium	Play	38.8 \pm 1.5	73.2 \pm 1.5	83.8 \pm 1.9	89.0 \pm 1.6	55.9 \pm 0.8	51.3 \pm 0.7	28.7 \pm 0.4	29.4 \pm 0.5		
Average		57.8	29.6	79.4	83.2	57.1	25.8	31.9	50.6		
Maze2D	Large	168.2 \pm 2.0	188.5 \pm 1.6	203.6 \pm 1.4	198.1 \pm 1.5	15.1 \pm 1.8	29.6 \pm 1.9	174.6 \pm 2.3	42.0 \pm 2.4		
Maze2D	Medium	127.1 \pm 1.4	141.3 \pm 1.2	150.7 \pm 1.0	151.4 \pm 0.9	31.8 \pm 1.7	97.8 \pm 2.2	75.7 \pm 2.1	153.4 \pm 4.2		
Maze2D	Umaze	131.8 \pm 1.3	117.7 \pm 1.6	136.6 \pm 1.3	133.1 \pm 1.2	21.6 \pm 2.4	114.1 \pm 2.7	30.2 \pm 2.4	58.7 \pm 2.5		
Average		142.4	149.2	163.6	160.9	22.8	80.5	93.5	84.7		
CG											
Dataset	Environment	Stride 1	Stride 2/5	Stride 4/15	Stride 8/25	Stride 1	Stride 2/5	Stride 4/15	Stride 8/25	None	
Mixed	Kitchen	64.0 \pm 0.4	65.9 \pm 0.4	61.3 \pm 0.4	-	63.5 \pm 0.4	63.5 \pm 0.4	54.2 \pm 0.3	-		
Partial	Kitchen	58.9 \pm 0.5	58.7 \pm 0.6	58.4 \pm 0.5	-	65.4 \pm 0.5	56.0 \pm 0.5	49.7 \pm 0.5	-		
Average		61.5	62.3	59.9	-	64.5	59.8	52.0	-		
Antmaze-Large	Diverse	79.1 \pm 1.3	35.3 \pm 1.5	20.6 \pm 1.1	77.2 \pm 1.3	73.1 \pm 1.4	16.8 \pm 1.6	11.6 \pm 1.4	78.2 \pm 1.8		
Antmaze-Large	Play	64.9 \pm 1.5	12.2 \pm 0.9	8.6 \pm 0.8	73.8 \pm 1.4	58.6 \pm 1.5	30.8 \pm 2.0	0.2 \pm 0.1	69.8 \pm 2.0		
Antmaze-Medium	Diverse	71.6 \pm 1.4	57.3 \pm 1.5	67.3 \pm 1.4	59.0 \pm 1.5	67.2 \pm 1.4	57.0 \pm 2.2	61.6 \pm 2.1	55.0 \pm 2.2		
Antmaze-Medium	Play	50.6 \pm 1.5	61.0 \pm 1.5	78.3 \pm 1.3	75.3 \pm 1.4	39.6 \pm 1.5	60.8 \pm 2.1	67.0 \pm 2.1	70.0 \pm 2.0		
Average		66.6	41.5	43.7	71.3	59.6	41.4	35.1	68.3		
Maze2D	Large	148.6 \pm 2.9	142.1 \pm 2.9	131.4 \pm 2.6	174.1 \pm 2.6	68.1 \pm 3.0	45.2 \pm 2.8	42.6 \pm 2.7	47.5 \pm 2.6		
Maze2D	Medium	110.1 \pm 2.2	113.3 \pm 2.2	123.2 \pm 2.1	126.9 \pm 2.2	60.8 \pm 2.3	61.9 \pm 2.3	41.5 \pm 2.2	47.3 \pm 2.3		
Maze2D	Umaze	80.4 \pm 2.4	64.3 \pm 2.4	67.6 \pm 2.4	81.0 \pm 2.4	47.6 \pm 2.4	38.9 \pm 2.4	27.5 \pm 2.4	35.3 \pm 2.4		
Average		113.0	106.6	107.4	127.3	58.8	48.7	37.2	43.4		

Table 5: Changing denoising network backbone for DV, with different planning strides. In ‘‘Stride x/y’’, x is for Kitchen, and y is for Maze2D & AntMaze. Data are Mean \pm Standard Error over 500 episode seeds.

Dataset	Environment	Transformer (DITID)					U-Net (U-NetID)				
		Stride 1	Stride 2/5	Stride 4/15	Stride 8/25	Stride 1	Stride 2/5	Stride 4/15	Stride 8/25		
Mixed	Kitchen	67.2 \pm 0.3	74.2 \pm 0.1	73.6 \pm 0.1	–	–	11.6 \pm 0.5	35.2 \pm 0.5	38.4 \pm 0.8	–	
Partial	Kitchen	87.8 \pm 0.6	88.3 \pm 0.4	94.0 \pm 0.3	–	–	12.4 \pm 0.5	28.8 \pm 0.6	10.8 \pm 0.5	–	
Average		77.5	81.3	83.8			12.0	32.0	24.6		
Antmaze-Large	Diverse	72.4 \pm 1.4	2.6 \pm 0.1	71.2 \pm 2.0	80.0 \pm 1.8	53.7 \pm 1.5	59.9 \pm 1.5	76.7 \pm 1.3	68.0 \pm 1.4		
Antmaze-Large	Play	61.2 \pm 1.5	18.2 \pm 0.5	76.2 \pm 0.9	76.4 \pm 2.0	53.4 \pm 1.5	48.3 \pm 1.5	80.8 \pm 1.2	72.7 \pm 1.4		
Antmaze-Medium	Diverse	58.8 \pm 1.5	24.2 \pm 0.8	86.2 \pm 1.6	87.4 \pm 1.6	27.9 \pm 1.4	55.5 \pm 1.5	85.8 \pm 1.1	81.1 \pm 1.2		
Antmaze-Medium	Play	38.8 \pm 1.5	73.2 \pm 1.5	83.8 \pm 1.9	89.0 \pm 1.6	28.5 \pm 1.4	66.2 \pm 1.4	81.0 \pm 1.2	81.9 \pm 1.2		
Average		57.8	29.6	79.4	83.2	40.9	57.5	81.1	75.9		
Maze2D	Large	168.2 \pm 2.0	188.5 \pm 1.6	203.6 \pm 1.4	198.1 \pm 1.5	172.5 \pm 1.9	168.7 \pm 2.2	193.1 \pm 1.5	189.7 \pm 1.5		
Maze2D	Medium	127.1 \pm 1.4	141.3 \pm 1.2	150.7 \pm 1.0	151.4 \pm 0.9	140.0 \pm 1.1	138.7 \pm 1.1	145.9 \pm 1.1	146.6 \pm 1.0		
Maze2D	Umaze	131.8 \pm 1.3	117.7 \pm 1.6	136.6 \pm 1.3	133.1 \pm 1.2	126.4 \pm 1.4	120.0 \pm 1.6	126.5 \pm 1.4	121.4 \pm 1.4		
Average		142.4	149.2	163.6	160.9	146.3	142.5	155.2	152.6		

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Table 6: The impact of action generation method for DV, with different planning strides. In “Stride x/y”, x is for Kitchen, and y is for Maze2D & AntMaze. Data are Mean \pm Standard Error over 500 episode seeds.

Dataset	Environment	Separate					Joint				
		Stride 1	Stride 2/5	Stride 4/15	Stride 8/25	Stride 1	Stride 2/5	Stride 4/15	Stride 8/25		
Mixed	Kitchen	67.2 \pm 0.3	74.2 \pm 0.1	73.6 \pm 0.1	–	60.8 \pm 0.4	62.9 \pm 0.4	56.6 \pm 0.6	–		
Partial	Kitchen	87.8 \pm 0.6	88.3 \pm 0.4	94.0 \pm 0.3	–	43.8 \pm 0.7	53.7 \pm 0.8	41.2 \pm 0.7	–		
Average		77.5	81.3	83.8		52.3	58.3	48.9			
Antmaze-Large	Diverse	72.4 \pm 1.4	2.6 \pm 0.1	71.2 \pm 2.0	80.0 \pm 1.8	49.9 \pm 1.5	0.6 \pm 0.2	0.2 \pm 0.1	7.3 \pm 0.8		
Antmaze-Large	Play	61.2 \pm 1.5	18.2 \pm 0.5	76.2 \pm 0.9	76.4 \pm 2.0	56.7 \pm 1.5	0.4 \pm 0.2	0.1 \pm 0.1	3.7 \pm 0.5		
Antmaze-Medium	Diverse	58.8 \pm 1.5	24.2 \pm 0.8	86.2 \pm 1.6	87.4 \pm 1.6	36.3 \pm 1.5	2.5 \pm 0.4	0.3 \pm 0.1	6.3 \pm 0.7		
Antmaze-Medium	Play	38.8 \pm 1.5	73.2 \pm 1.5	83.8 \pm 1.9	89.0 \pm 1.6	34.0 \pm 1.4	2.7 \pm 0.5	3.7 \pm 0.6	27.8 \pm 1.4		
Average		57.8	29.6	79.4	83.2	44.2	1.6	1.1	11.3		
Maze2D	Large	168.2 \pm 2.0	188.5 \pm 1.6	203.6 \pm 1.4	198.1 \pm 1.5	187.0 \pm 1.6	202.6 \pm 1.4	202.9 \pm 1.5	198.5 \pm 1.6		
Maze2D	Medium	127.1 \pm 1.4	141.3 \pm 1.2	150.7 \pm 1.0	151.4 \pm 0.9	141.4 \pm 1.1	143.7 \pm 1.1	156.0 \pm 0.9	153.8 \pm 0.9		
Maze2D	Umaze	131.8 \pm 1.3	117.7 \pm 1.6	136.6 \pm 1.3	133.1 \pm 1.2	124.9 \pm 1.6	139.8 \pm 1.2	146.3 \pm 1.1	141.1 \pm 1.2		
Average		142.4	149.2	163.6	160.9	151.1	162.0	168.4	164.5		

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403