

# Structural Information-based Hierarchical Diffusion for Offline Reinforcement Learning

Xianghua Zeng<sup>1</sup>, Hao Peng<sup>1</sup>, Yicheng Pan<sup>1</sup>, Angsheng Li<sup>1,2</sup>, Guanlin Wu<sup>3</sup>

<sup>1</sup>State Key Laboratory of Software Development Environment, Beihang University, Beijing, China

<sup>2</sup>Zhongguancun Laboratory, Beijing, China

<sup>3</sup>National University of Defense Technology, Changsha, China

**Email:** [zengxianghua@buaa.edu.cn](mailto:zengxianghua@buaa.edu.cn)

**Code:** <https://github.com/SELGroup/SIHD>



# Hierarchical Diffusion-Based Offline RL

## Offline Reinforcement Learning:

- policy learning without any environmental interaction
- extrapolation error caused by out-of-distribution states and actions
- mode-covering challenge from diverse hierarchical trajectories

## Diffusion Models for Offline RL:

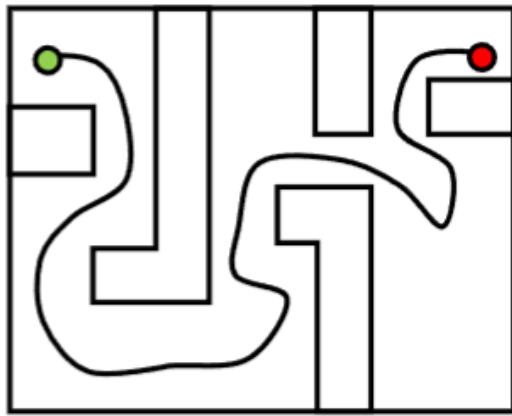
- building diffusion-based policies with strong generative capacity
- reframing decision-making as offline trajectory generation
- exponential increase in estimate variance and high computational cost

## Hierarchical Diffusion for Offline RL:

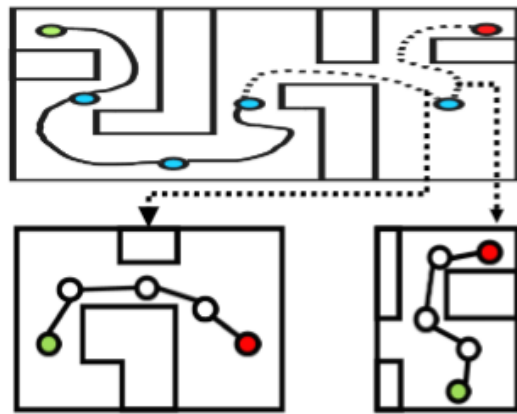
- task decomposition into manageable subproblems
- introducing manually predefined skills into diffusion
- trajectory segmentation for subgoal-directed diffusion

# Hierarchical Diffusion-Based Offline RL

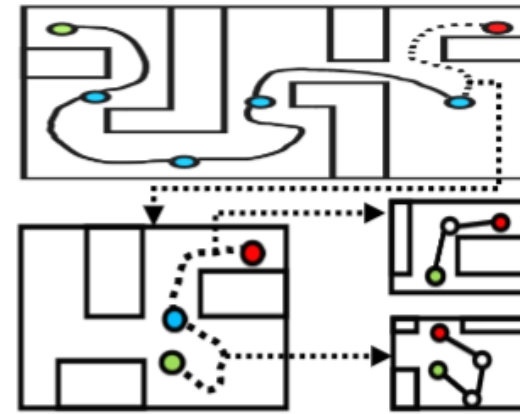
## Illustrative Example For Hierarchical Trajectory Generation:



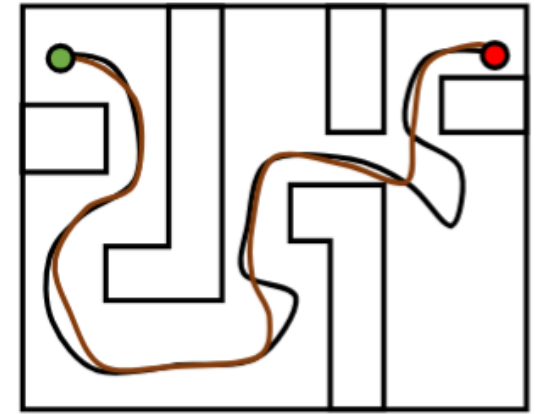
a) offline trajectory



b) 2-layer rigid diffusion



c) multi-scale hierarchy



d) regularized exploration

### Summary:

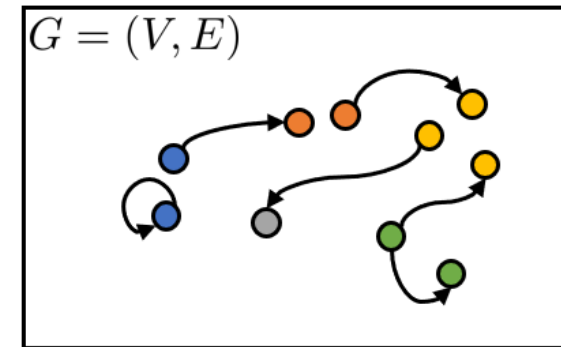
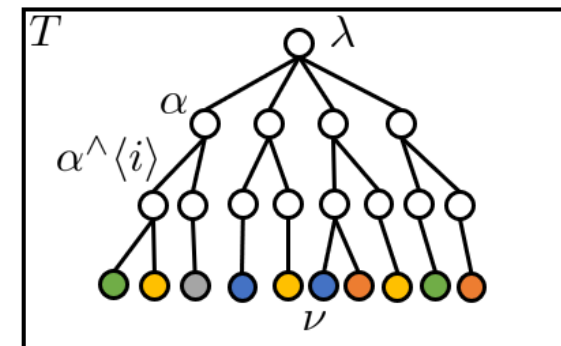
- single predefined temporal scale to segment offline trajectories
- fixed two-layer diffusion hierarchy composed of subgoal and action layers
- hindering adaptability to varying temporal patterns and task-specific complexities

# Structural Information Principles

The structural entropy measures the uncertainty of a graph under a strategy of hierarchical partitioning, which is called an “encoding tree”.

The encoding tree of  $G$  is defined as a rooted tree  $T$  with the following properties:

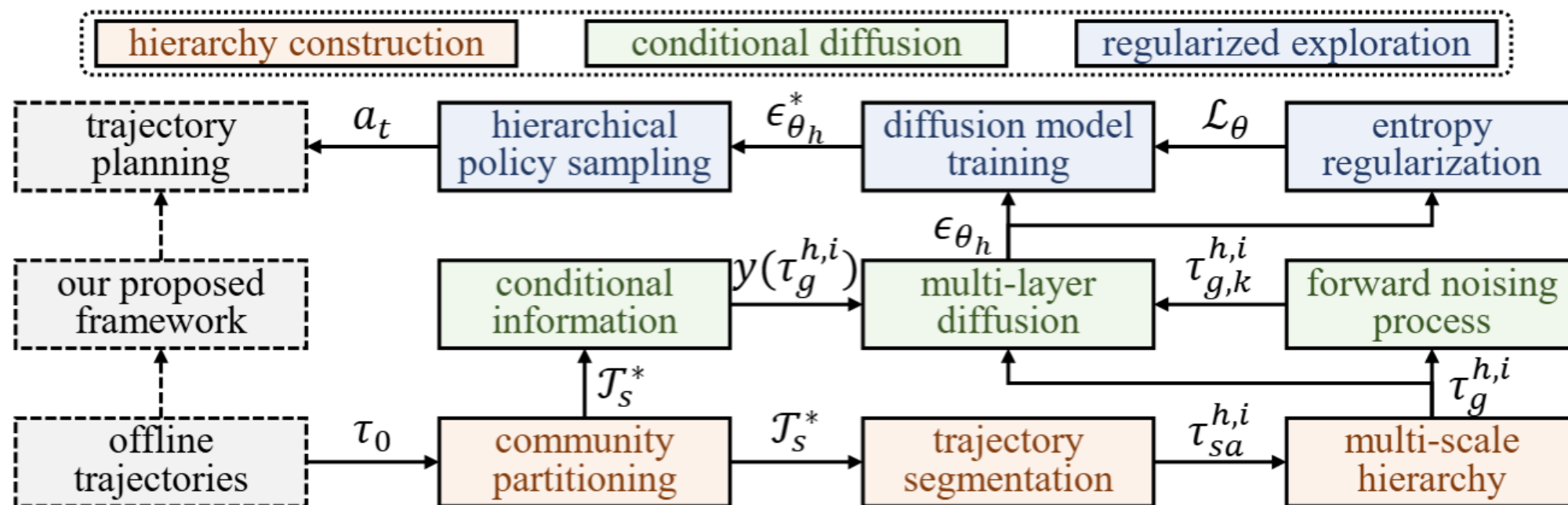
- 1) For each node  $\alpha \in T$ , there is a vertex subset in  $G$  corresponding with  $\alpha$ ,  $T_\alpha \subseteq V$ .
- 2) For the root node  $\lambda$ , we set  $T_\lambda = V$ .
- 3) For each node  $\alpha \in T$ , its children nodes are marked as  $\alpha^{\wedge\langle i \rangle}$  ordered from left to right as  $i$  increases, and  $\alpha^{\wedge\langle i \rangle^-} = \alpha$ .
- 4) For each node  $\alpha \in T$ , we suppose that  $L$  is the number of its children nodes; then all vertex subsets  $T_{\alpha^{\wedge\langle i \rangle}}$  are disjointed, and  $T_\alpha = \bigcup_{i=1}^L T_{\alpha^{\wedge\langle i \rangle}}$ .
- 5) For each leaf node  $\nu$ ,  $T_\nu$  is a singleton containing a single vertex in  $V$ .

Graph  $G$ Encoding Tree  $T$

# Our Proposed Adversarial Attack Framework

## Overall SIHD Framework:

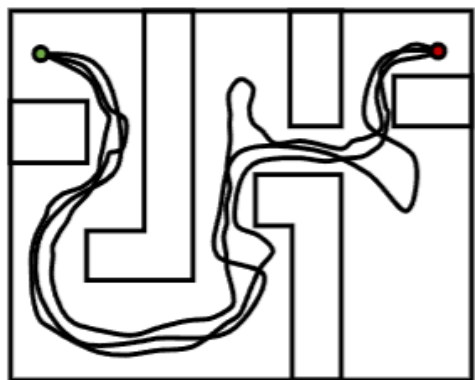
- structural information guided construction of adaptive multi-scale diffusion hierarchy
- conditional forward diffusion and reverse inference at multiple time scales
- encouraging regularized exploration for underrepresented offline states



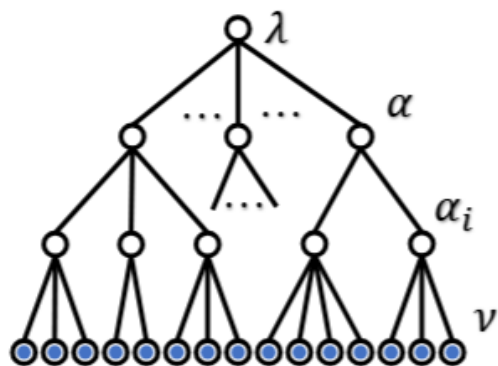
# Our Proposed Hierarchical Diffusion Framework

## Stage 1: Multi-Scale Diffusion Hierarchy

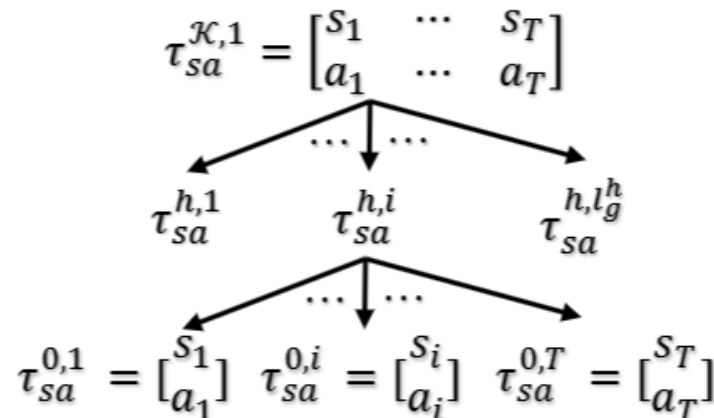
- leverage feature similarity to identify structural relationships among offline states
- derive a tree-structured partitioning of offline state communities
- construct the trajectory-adaptive multi-scale hierarchical diffusion instead of the rigid two-layer diffusion hierarchy operating at a single scale



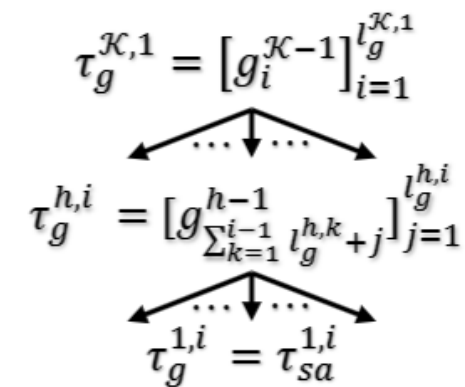
a) offline trajectories



b) community partitioning



c) hierarchical segmentation



d) multi-layer subgoals

# Our Proposed Hierarchical Diffusion Framework

## Stage 2: Conditional Diffusion Model

- generate a subgoal sequence conditioned on task rewards at the top of diffusion hierarchy
- generate subgoal sequences at intermediate levels of the hierarchy or state-action sequences at the base level, with each generation process conditioned on the parent subgoal
- compute the structural information gain as classifier-free guidance to enhance conditional hierarchical generation

## Stage 3: Structural Entropy Regularizer

- introduce a structural entropy-based exploration regularizer
- encourage the hierarchical policy to explore underrepresented regions of the state space
- limit extrapolation errors resulting from deviations from the offline behavioral policy



# Comparative Experiments on the D4RL Benchmark

- Datasets:** Gym-MuJoCo, Maze2D and AntMaze benchmark tasks
- Baselines:** model-free methods (CQL, IQL, and Decision Transformer), model-based methods (MoReL and Trajectory Transformer), diffusion-based methods (Diffuser, HDMI, HD)

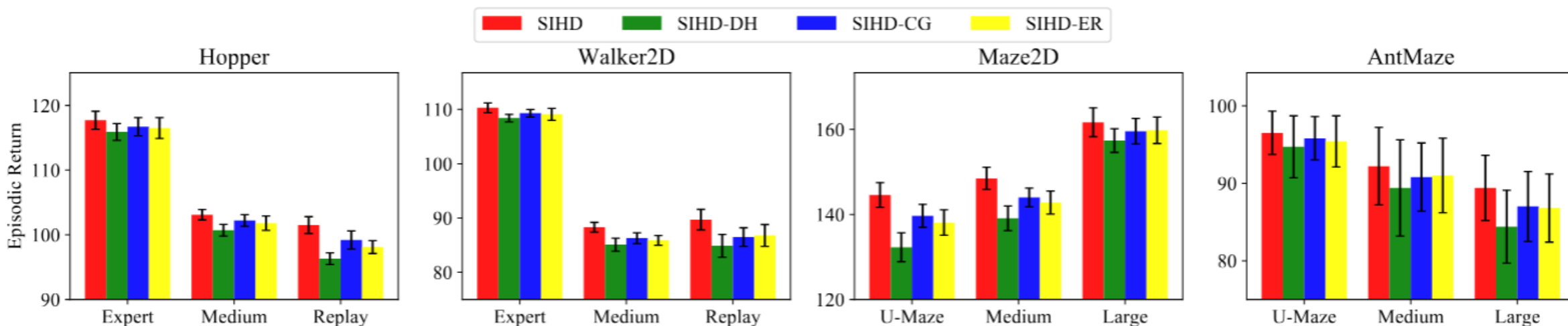
Gym-MuJoCo Tasks	HalfCheetah			Hopper			Walker2D		
	Expert	Medium	Replay	Expert	Medium	Replay	Expert	Medium	Replay
CQL	91.6	44.0	45.5	105.4	58.5	95.0	108.8	72.5	77.2
IQL	86.7	47.4	44.2	91.5	66.3	94.7	109.6	78.3	73.9
DT	86.8	42.6	36.6	107.6	67.6	82.7	108.1	74.0	66.6
MoReL	53.3	42.1	40.2	108.7	95.4	93.6	95.6	77.8	49.8
TT	95.0	46.9	41.9	110.0	61.1	91.5	101.9	79.0	82.6
Diffuser	88.9 $\pm$ 0.3	42.8 $\pm$ 0.3	37.7 $\pm$ 0.5	103.3 $\pm$ 1.3	74.3 $\pm$ 1.4	93.6 $\pm$ 0.4	106.9 $\pm$ 0.2	79.6 $\pm$ 0.6	70.6 $\pm$ 1.6
HDMI	92.1 $\pm$ 1.4	48.0 $\pm$ 0.9	44.9 $\pm$ 2.0	113.5 $\pm$ 0.9	76.4 $\pm$ 2.6	99.6 $\pm$ 1.5	107.9 $\pm$ 1.2	79.9 $\pm$ 1.8	80.7 $\pm$ 2.1
HD	92.5 $\pm$ 0.3	46.7 $\pm$ 0.2	38.1 $\pm$ 0.7	115.3 $\pm$ 1.1	99.3 $\pm$ 0.3	94.7 $\pm$ 0.7	107.1 $\pm$ 0.1	84.0 $\pm$ 0.6	84.1 $\pm$ 2.2
SIHD	<b>94.4 <math>\pm</math> 0.5</b>	<b>48.7 <math>\pm</math> 1.1</b>	<b>47.0 <math>\pm</math> 0.4</b>	<b>117.7 <math>\pm</math> 1.4</b>	<b>103.1 <math>\pm</math> 0.8</b>	<b>101.5 <math>\pm</math> 1.3</b>	<b>110.3 <math>\pm</math> 0.9</b>	<b>88.3 <math>\pm</math> 0.9</b>	<b>89.7 <math>\pm</math> 1.9</b>
Abs.(%) Avg.	1.9(2.1)	0.7(1.5)	1.5(3.3)	2.4(2.1)	3.8(3.8)	1.9(1.9)	0.7(0.6)	4.3(5.1)	5.6(6.7)

Gym-MuJoCo Tasks	Single-task Maze2D			Multi-task Maze2D			AntMaze		
	U-Maze	Medium	Large	U-Maze	Medium	Large	U-Maze	Medium	Large
IQL	47.4	34.9	58.6	24.8	12.1	13.9	62.2	70.0	47.5
MPPI	33.2	10.2	5.1	41.2	15.4	8.0	-	-	-
Diffuser	113.9 $\pm$ 3.1	121.5 $\pm$ 2.7	123.0 $\pm$ 6.4	128.9 $\pm$ 1.8	127.2 $\pm$ 3.4	132.1 $\pm$ 5.8	76.0 $\pm$ 7.6	31.9 $\pm$ 5.1	-
IRIS	-	-	-	-	-	-	89.4 $\pm$ 2.4	64.8 $\pm$ 2.6	43.7 $\pm$ 1.3
HiGoC	-	-	-	-	-	-	91.2 $\pm$ 1.9	79.3 $\pm$ 2.5	67.3 $\pm$ 3.1
HDMI	120.1 $\pm$ 2.5	121.8 $\pm$ 1.6	128.6 $\pm$ 2.9	131.3 $\pm$ 1.8	131.6 $\pm$ 1.9	135.4 $\pm$ 2.5	-	-	-
HD	128.4 $\pm$ 3.6	135.6 $\pm$ 3.0	155.8 $\pm$ 2.5	144.1 $\pm$ 1.2	140.2 $\pm$ 1.6	165.5 $\pm$ 0.6	94.0 $\pm$ 4.9	88.7 $\pm$ 8.1	83.6 $\pm$ 5.8
SIHD	<b>144.6 <math>\pm</math> 2.9</b>	<b>148.5 <math>\pm</math> 2.6</b>	<b>161.7 <math>\pm</math> 3.4</b>	<b>157.0 <math>\pm</math> 0.6</b>	<b>156.8 <math>\pm</math> 1.7</b>	<b>169.4 <math>\pm</math> 2.7</b>	<b>96.5 <math>\pm</math> 2.8</b>	<b>92.2 <math>\pm</math> 5.0</b>	<b>89.4 <math>\pm</math> 4.2</b>
Abs.(%) Avg.	16.2(12.6)	12.9(9.5)	5.9(3.8)	12.9(9.0)	16.6(11.8)	3.9(2.4)	2.5(2.7)	3.5(3.9)	5.8(6.9)



## Ablation Study on Primary SIHD Modules

- SIHD-DH: replacing the multi-scale hierarchy with a rigid two-layer single-scale diffuser
- SIHD-CG: removing the classifier-based guidance in conditional diffusion
- SIHD-ER: removing the structural entropy regularizer in policy learning
- highlighting the importance of each component and the critical role of the multi-scale diffusion hierarchy in SIHD
- becoming even more pronounced in long-horizon decision-making tasks



## Conclusion and Future Works

- This work proposes SIHD, a novel hierarchical diffusion framework that leverages structural information embedded in historical trajectories to construct an adaptive multi-scale diffusion hierarchy and promote exploration of underrepresented states in offline datasets, thereby enabling effective policy learning.
- Extensive evaluations on the challenging D4RL benchmark demonstrate the superior decision-making performance and generalization capabilities of SIHD across diverse offline RL tasks.
- In future work, we aim to refine the hierarchical diffusion framework by further exploring how to more effectively represent and integrate subgoal constraints as conditional information.
- We also plan to extend the SIHD framework to other offline RL environments and broader diffusion-based generative modeling domains.

# Structural Information-based Hierarchical Diffusion for Offline Reinforcement Learning

Xianghua Zeng<sup>1</sup>, Hao Peng<sup>1</sup>, Yicheng Pan<sup>1</sup>, Angsheng Li<sup>1,2</sup>, Guanlin Wu<sup>3</sup>

<sup>1</sup>State Key Laboratory of Software Development Environment, Beihang University, Beijing, China

<sup>2</sup>Zhongguancun Laboratory, Beijing, China

<sup>3</sup>National University of Defense Technology, Changsha, China

**Email:** [zengxianghua@buaa.edu.cn](mailto:zengxianghua@buaa.edu.cn)

**Code:** <https://github.com/SELGroup/SIHD>