

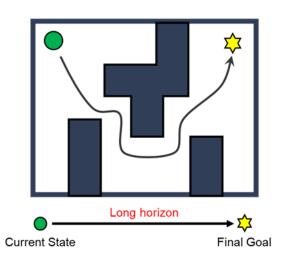
# Graph-Assisted Stitching for Offline Hierarchical Reinforcement Learning

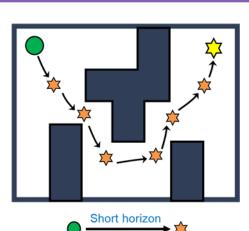
(paper, code, video)

**Project website** 

Seungho Baek, Taegeon Park, Jongchan Park, Seungjun Oh, Yusung Kim Sungkyunkwan University

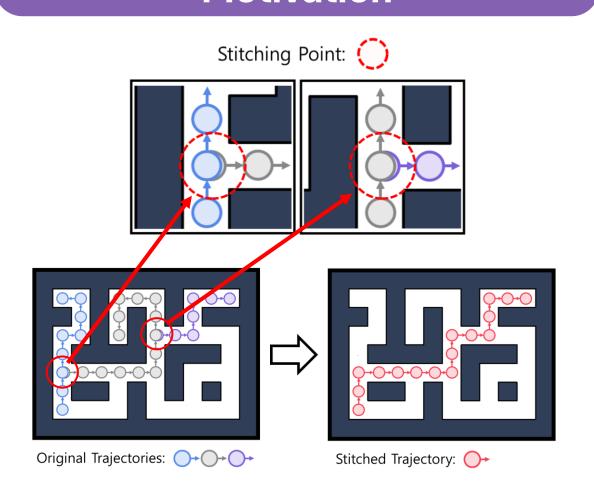
# Introduction





- Offline HRL introduces a two-level decision-making framework, where a high-level policy generates subgoals and a low-level policy executes primitive actions to reach them.
- This hierarchical structure decomposes complex long-horizon problems into manageable short-horizon subproblems.

## **Motivation**

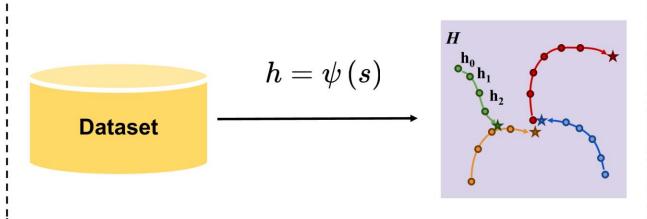


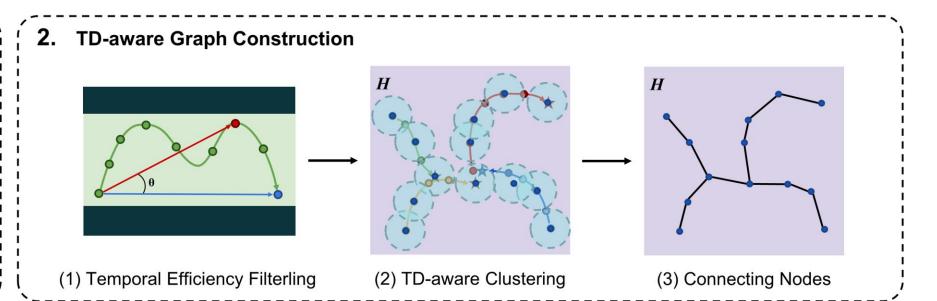
### Why is Trajectory Stitching Crucial in Offline HRL?

- Stitching composes new trajectories by connecting partial segments from different goal-oriented trajectories.
- This enables agents to utilize transitions that are **temporally** disjoint and not observed within a single trajectory.
- Such stitching facilitates generalization across goals, especially in sparse-reward and long-horizon tasks.
- However, existing offline HRL methods typically lack mechanisms for cross-goal stitching, limiting their ability to compose effective subgoal sequences from suboptimal trajectories.

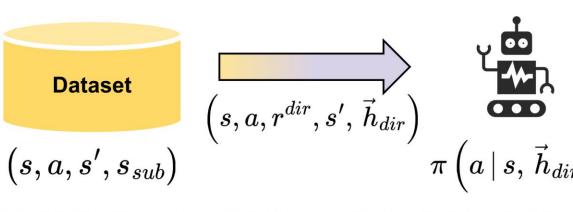
# **Overview of GAS**

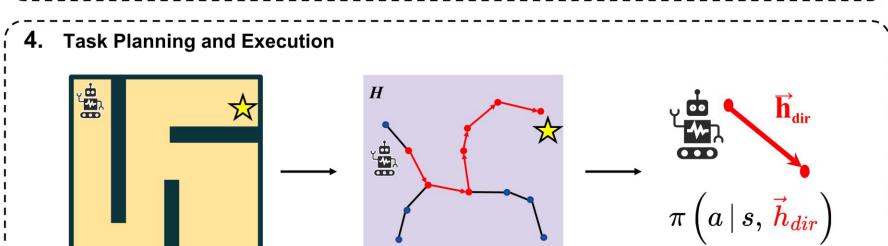
### 1. Pre-Training Temporal Distance Representation





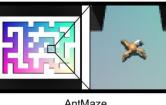






# **Experiments**

#### **Datasets**









### Benchmark results: state-based environments

| <b>Dataset Type</b> | Dataset   | GCBC  | GCIQL   | QRL  | CRL  | <b>HGCBC</b>                                  | HHILP  | HIQL   | GAS (ours)   |
|---------------------|---|---|---|--|--|---|--|--|--|
| Locomotion          | antmaze-medium-navigate<br>antmaze-large-navigate<br>antmaze-giant-navigate | $33.1 \pm 5.6 \\ 23.4 \pm 3.2 \\ 0.0 \pm 0.0$               | $74.6 \pm 4.8 \\ 32.6 \pm 4.7 \\ 0.1 \pm 0.4$                             | $81.9 \pm 8.2$<br>$74.9 \pm 4.4$<br>$14.3 \pm 3.6$ | $95.3 \pm 1.0$<br>$85.5 \pm 5.3$<br>$15.0 \pm 5.7$ | $58.1 \pm 5.5$ $44.3 \pm 4.1$ $7.2 \pm 1.7$   | $96.3 \pm 0.4$<br>$86.8 \pm 3.6$<br>$53.1 \pm 2.6$ | $95.3 \pm 1.3$<br>$89.9 \pm 2.2$<br>$67.3 \pm 5.5$                         | $96.3 \pm 1.3$<br>$93.2 \pm 0.5$<br>$77.6 \pm 2.9$ |
| Stitching           | antmaze-medium-stitch<br>antmaze-large-stitch<br>antmaze-giant-stitch       | $43.2 \pm 7.7 \\ 2.3 \pm 3.6 \\ 0.0 \pm 0.0$                | $\begin{array}{c} 26.6 \pm 6.8 \\ 9.6 \pm 3.1 \\ 0.0 \pm 0.0 \end{array}$ | $67.0 \pm 10.6 \\ 20.2 \pm 1.7 \\ 0.4 \pm 0.3$     | $57.0 \pm 7.9 \\ 14.4 \pm 5.9 \\ 0.0 \pm 0.0$      | $65.9 \pm 5.7 \\ 10.7 \pm 5.8 \\ 0.0 \pm 0.0$ | $96.0 \pm 1.2$<br>$34.1 \pm 3.0$<br>$0.0 \pm 0.0$  | $\begin{array}{c} 92.0 \pm 2.8 \\ 71.7 \pm 4.8 \\ 1.0 \pm 1.2 \end{array}$ | $98.1 \pm 1.2$ $96.3 \pm 0.9$ $88.3 \pm 3.6$       |
| Exploratory         | antmaze-medium-explore<br>antmaze-large-explore                             | $\begin{array}{c} 2.7  \pm 2.8 \\ 0.0  \pm 0.0 \end{array}$ | $\begin{array}{c} 11.7 \pm 1.3 \\ 0.6 \pm 0.5 \end{array}$                | $1.4 \pm 1.2$<br>$0.3 \pm 1.0$                     | $1.0 \pm 1.6$<br>$0.0 \pm 0.0$                     | $15.0 \pm 8.2 \\ 0.0 \pm 0.0$                 | $39.9 \pm 7.4$<br>$2.4 \pm 1.9$                    | $32.2 \pm 3.0 \\ 2.9 \pm 4.3$  | $98.1 \pm 0.4$<br>$94.2 \pm 3.0$                   |
| Manipulation        | scene-play  | $5.4 \pm 0.9$<br>$69.5 \pm 14.1$                            | 50.4 ± 1.4  | $5.1 \pm 1.7$ $61.9 \pm 85$                        | $19.2 \pm 3.0$ $32.7 \pm 11.7$                     | $4.6 \pm 1.3$ $71.1 \pm 62$                   | $43.4 \pm 5.2$ $66.7 \pm 9.0$                      | $40.0 \pm 9.6$<br>$73.1 \pm 2.4$   | $73.6 \pm 8.0$ $87.3 \pm 8.8$                      |

### **Benchmark results:** pixel-based environments

| <b>Dataset Type</b> | Dataset  | GCIQL  | QRL   | CRL  | HHILP  | HIQL  | GAS (ours)   |
|---------------------|--|--|---|--|--|---|--|
| Locomotion          | visual-antmaze-medium-navigate<br>visual-antmaze-large-navigate<br>visual-antmaze-giant-navigate | $19.1 \pm 1.6$ $4.6 \pm 1.9$ $1.5 \pm 0.8$                         | $0.0 \pm 0.0$<br>$0.0 \pm 0.0$<br>$0.2 \pm 0.8$ | $93.7 \pm 1.2$<br>$79.5 \pm 7.5$<br>$43.4 \pm 5.9$ | $94.1 \pm 1.2$<br>$85.6 \pm 2.5$<br>$42.4 \pm 1.9$ | $95.5 \pm 0.8$<br>$80.0 \pm 2.1$<br>$34.1 \pm 14.0$ | $96.4 \pm 0.5$<br>$87.0 \pm 1.2$<br>$59.0 \pm 2.1$ |
| Stitching           | visual-antmaze-medium-stitch<br>visual-antmaze-large-stitch<br>visual-antmaze-giant-stitch       | $\begin{array}{c} 4.2\pm1.6 \\ 0.2\pm0.3 \\ 0.0\pm0.0 \end{array}$ | $0.0 \pm 0.0$<br>$0.1 \pm 0.5$<br>$0.2 \pm 0.6$ | $68.0 \pm 8.3 \\ 14.7 \pm 7.1 \\ 0.0 \pm 0.0$      | $92.4 \pm 1.2$<br>$33.8 \pm 1.2$<br>$3.6 \pm 1.3$  | $90.4 \pm 4.1$<br>$38.5 \pm 5.7$<br>$0.9 \pm 1.1$   | $90.0 \pm 3.0$<br>75.2 ± 4.4<br>55.8 ± 3.5         |
| Exploratory         | visual-antmaze-medium-explore<br>visual-antmaze-large-explore                                    | $0.0 \pm 0.0 \\ 0.0 \pm 0.0$                                       | $0.1 \pm 0.3$<br>$0.0 \pm 0.0$                  | $0.0 \pm 0.0$<br>$0.0 \pm 0.0$                     | $0.0 \pm 0.0 \\ 0.0 \pm 0.0$                       | $0.9 \pm 1.4$<br>$0.0 \pm 0.0$                      | $65.9 \pm 6.8$<br>$15.1 \pm 6.8$                   |
| Manipulation        | visual-scene-play  | $10.6 \pm 2.7$   | $13.5 \pm 2.8$                                  | $8.4\pm 0.9$                                       | $35.6 \pm 4.9$                                     | $47.9 \pm 3.9$                                      | <b>54.4</b> ± 6.2                                  |

# **Ablation:** Temporal Efficiency Filtering

| Dataset                | # States in Dataset | TE-Filtered States (%) |      | # Nodes in Graph |      | Normalized Return |                        |                   |
|------------------------|---------------------|------------------------|------|------------------|------|-------------------|------------------------|-------------------|
| Dataset                |                     | All States             | Ours | All States       | Ours | All States        | Ours                   | $\Delta \uparrow$ |
| antmaze-giant-navigate | 1M                  | 100                    | 6    | 2092             | 978  | $63.4 \pm 3.7$    | <b>77.6</b> ± 2.9      | +14.2             |
| antmaze-giant-stitch   | 1 <b>M</b>          | 100                    | 8    | 3490             | 1966 | $75.3 \pm 5.7$    | <b>88.3</b> $\pm$ 3.6  | +13.0             |
| antmaze-large-explore  | 5M                  | 100                    | 2    | 6213             | 2499 | $75.4\pm 4.3$     | $94.2 \pm 3.0$         | +18.8             |
| scene-play             | 1 <b>M</b>          | 100                    | 6    | 2809             | 725  | $63.5\pm 5.5$     | $\textbf{73.6}\pm 8.0$ | +10.1             |

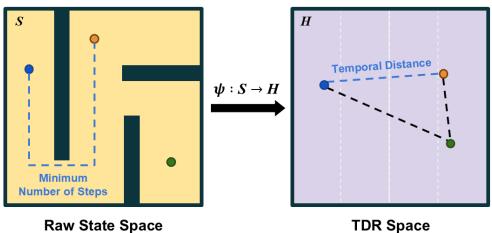
• This table compares the performance of GAS with TE filtering to a variant that constructs the graph from all states without filtering.

#### **Performance Highlights**

- Q. Does GAS excel at long-horizon reasoning?
- A. Yes, GAS shows strong performance on antmaze-giant-navigate and scene-play, which require substantial long-horizon reasoning capabilities in navigation and manipulation domains, respectively.
- Q. Does GAS demonstrate effective stitching ability?
- A. Yes, GAS outperforms baselines on antmaze-{medium, large, giant}-stitch, where the datasets consist of short goal-reaching trajectories.
- Q. Can GAS effectively learn from suboptimal datasets?
- A. Yes, GAS achieves the best performance on antmaze-{medium, large}-explore, where the datasets consist of extremely low-quality data.
- Q. Can GAS effectively handle image-based tasks?
- A. Yes, GAS demonstrates strong performance not only in state-based environments but also in *pixel-based environments*.

# Keyldeas

### **Temporal Distance Representation (TDR)**

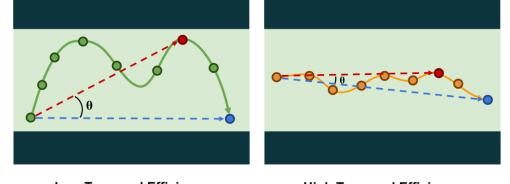


**Raw State Space** 

• TDR<sub>[1]</sub> $\psi$  embeds states into a latent space H, where the **Euclidean distance** between any two points corresponds to the minimum number of steps (i.e., optimal temporal distance) required to transition from one state to another in the raw state space S.

[1] Park et al., "Foundation Policies with Hilbert Representations", ICML 2024.

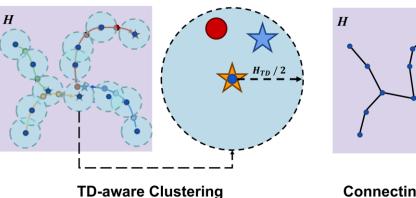
### Temporal Efficiency (TE)

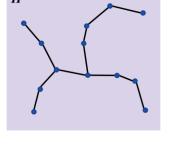




- TE measures the directional alignment between the actual and optimal transitions over a fixed temporal distance. (Sreached: state observed  $H_{TD}$  steps after  $s_{cur}$ ) (Soptimal: state at  $H_{TD}$  temporal distance from  $s_{cur}$ )
- Before graph construction, filtering low-TE states reduces construction overhead and improves graph quality.

### **TD-aware Graph Construction**





**Connecting Nodes** 

- GAS clusters states in the TDR space at regular temporal distance intervals  $H_{TD}$ , grouping semantically similar states from different trajectories.
- Each cluster center becomes a graph node, and edges are added between nodes if their temporal distance is below  $H_{-}TD$ , enabling cross-goal stitching across disconnected trajectories.