



# Temporal Distance-aware Transition Augmentation for Offline Model-based Reinforcement Learning

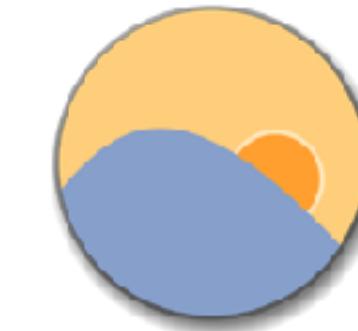
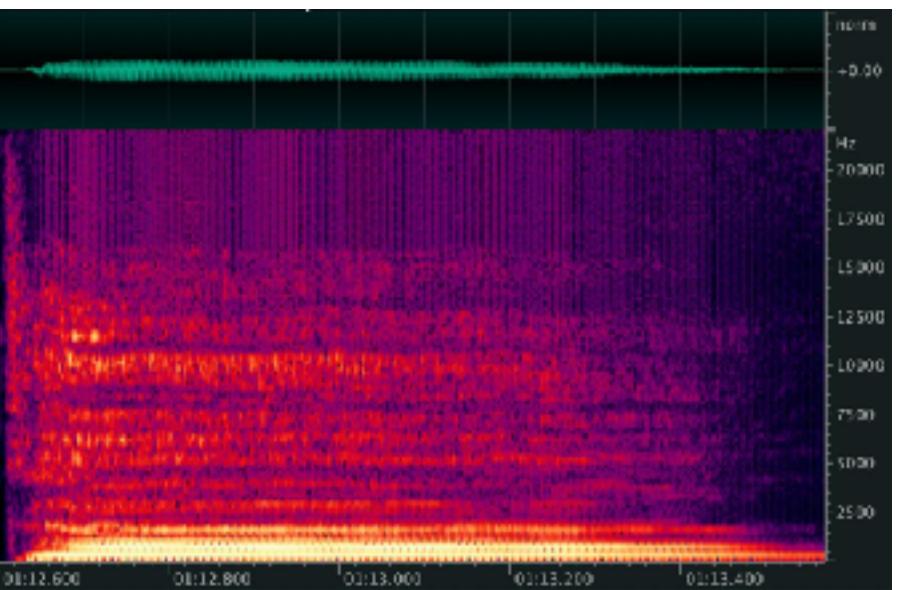
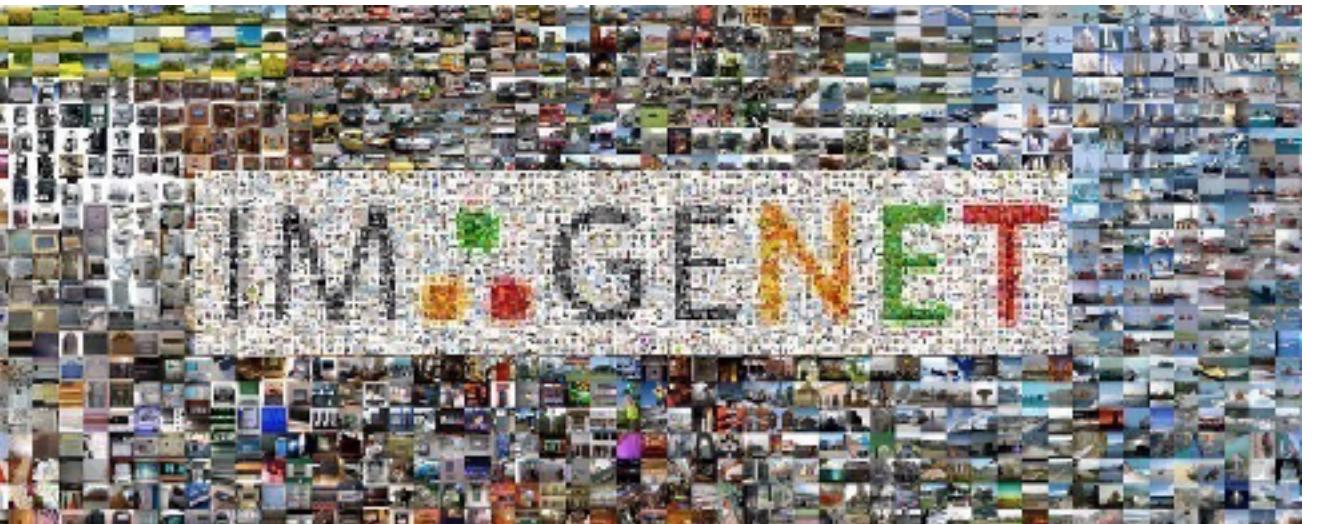
Dongsu Lee and Minhae Kwon

Brain and Machine Intelligence Lab. (BMIL)  
Soongsil University  
Seoul, Republic of Korea

This work was done while Dongsu Lee was visiting at Carnegie Mellon University

# Data-driven machine learning

discipline  
translation tool learning computational  
languages early likewise modern possible  
tools browse assist real-life combinations speech  
using corpora whether machines learner  
allows words simple co-occur world search  
corpus use well started corpora  
interface correctly



f.lux

# In reinforcement learning


$$\pi(a|s)$$

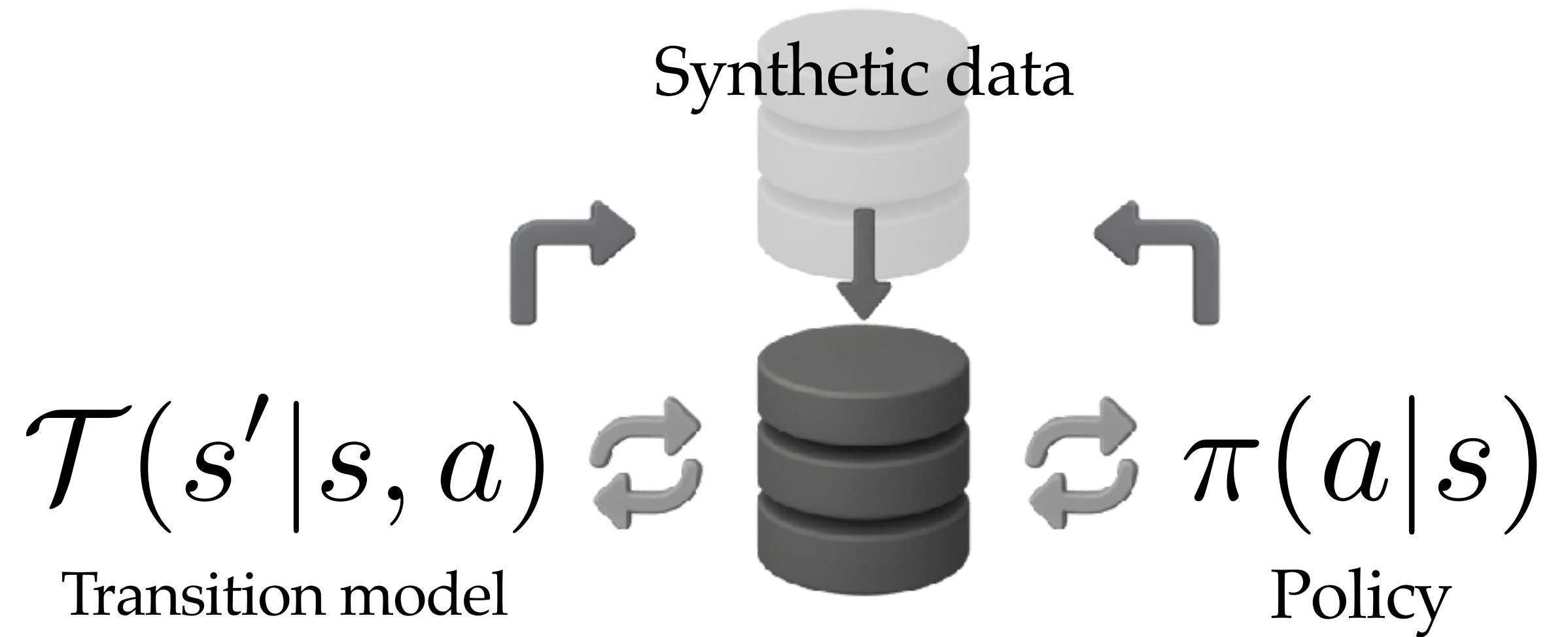

Trajectories


$$\pi(a|s)$$

Online RL

Offline RL

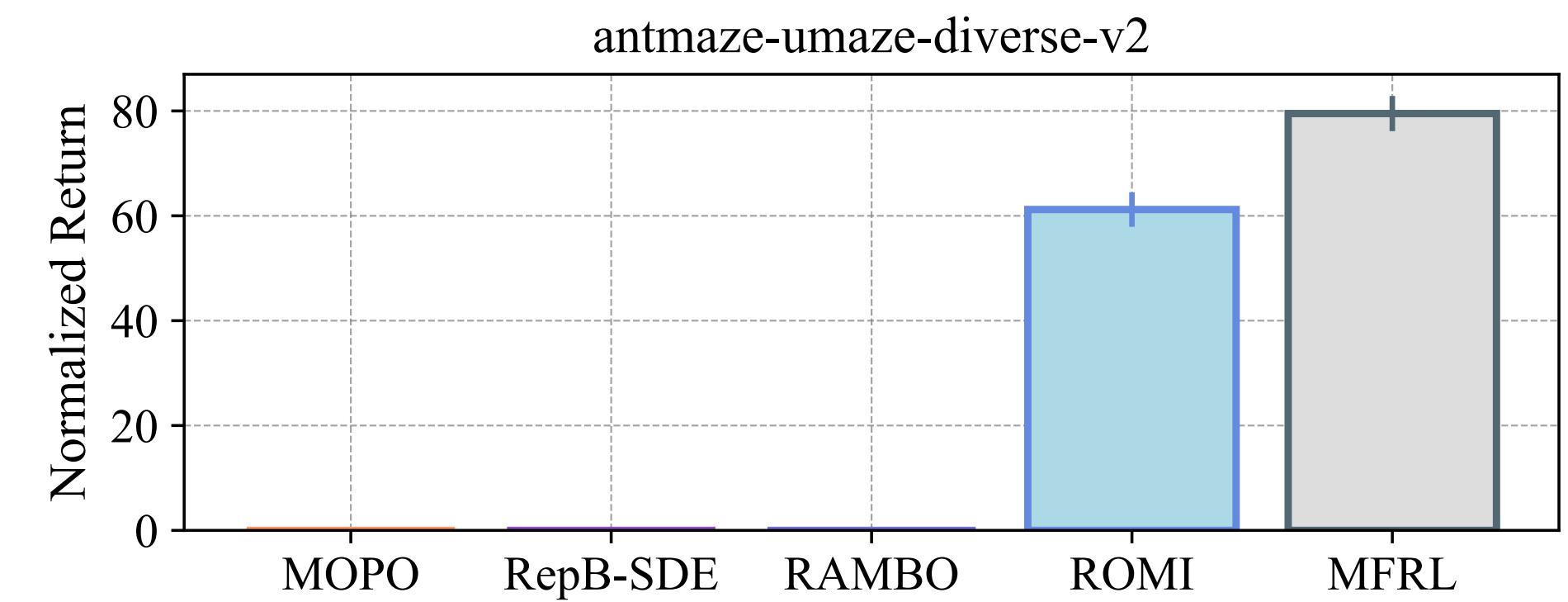
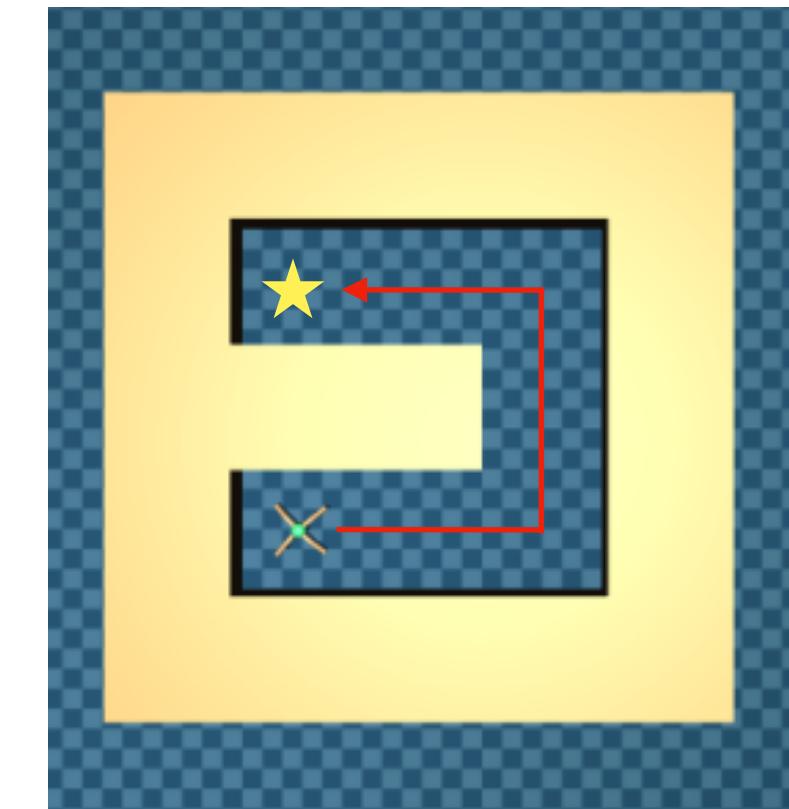
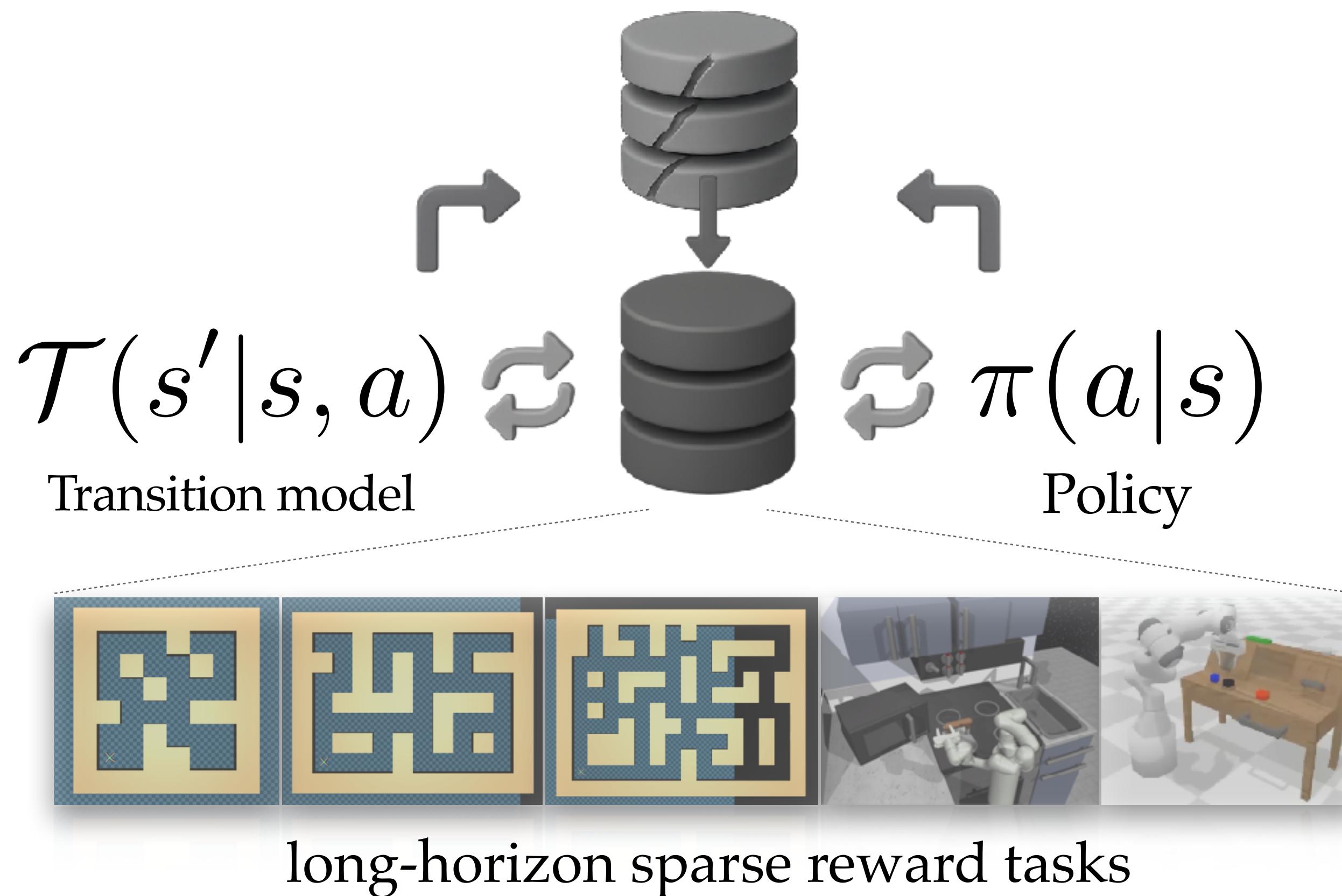
# Offline model-based reinforcement learning



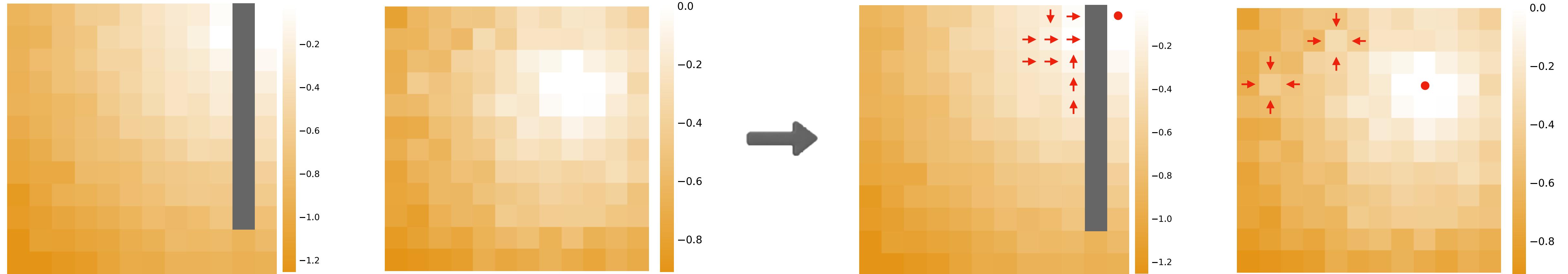
# Offline model-based reinforcement learning

However, in long-horizon sparse reward tasks...

Unfaithful and useless synthesized dataset



# Why do prior offline MBRLs fail in such tasks?

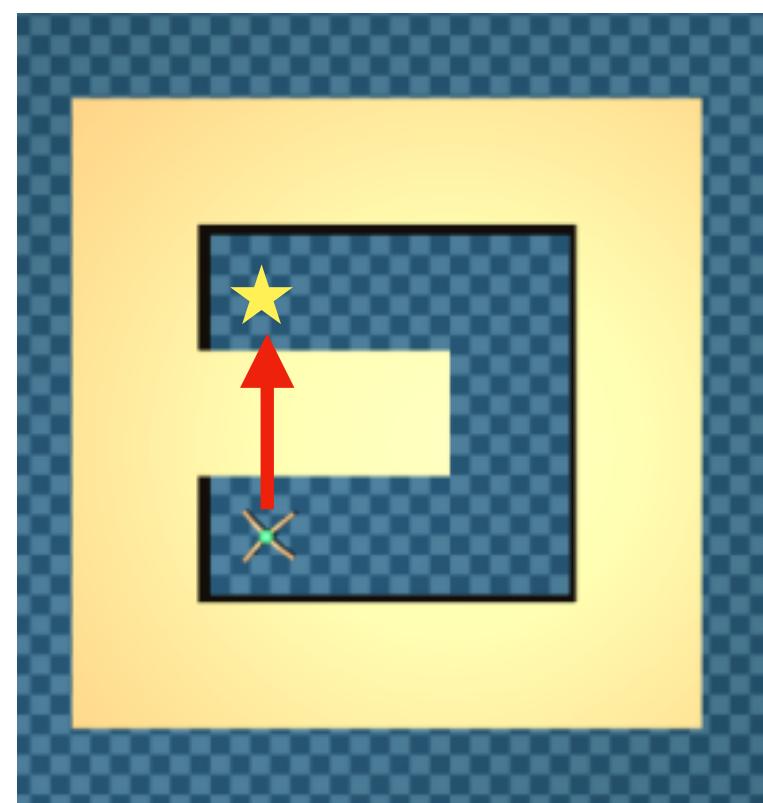


Overgeneralization  
on Euclidean space

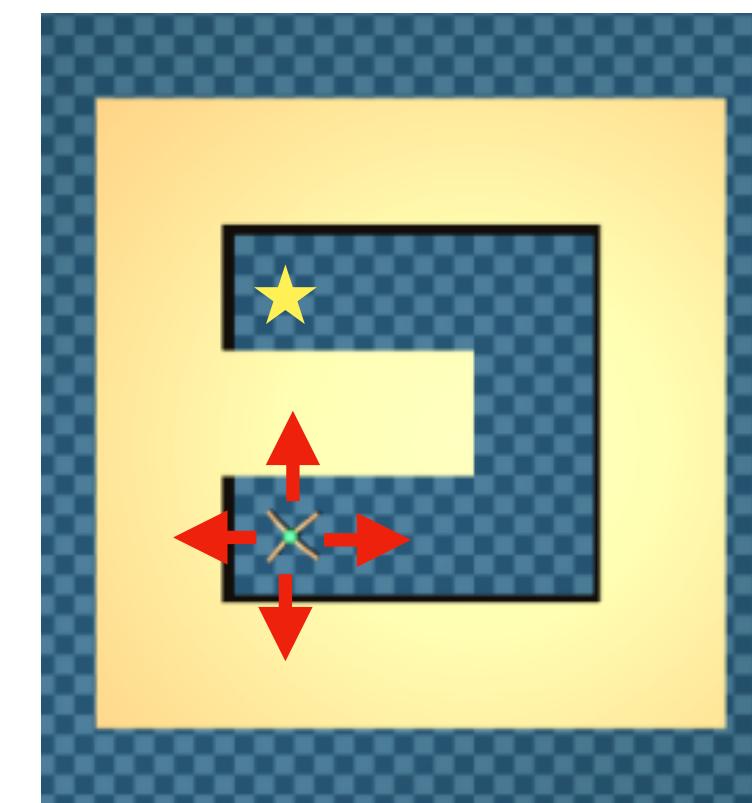
Noisy value function

Unstable policy optimization

Rollout of **erroneous** policy in Euclidean raw state space



Out of reach



Purposeless

Synthesized transition data not only fails  
to eliminate OOD, but actually **makes the**  
**value function and policy worse**

# Our solution: Temporal distance-aware state embedding

1. Building an autoencoder to encode raw state and decode latent state

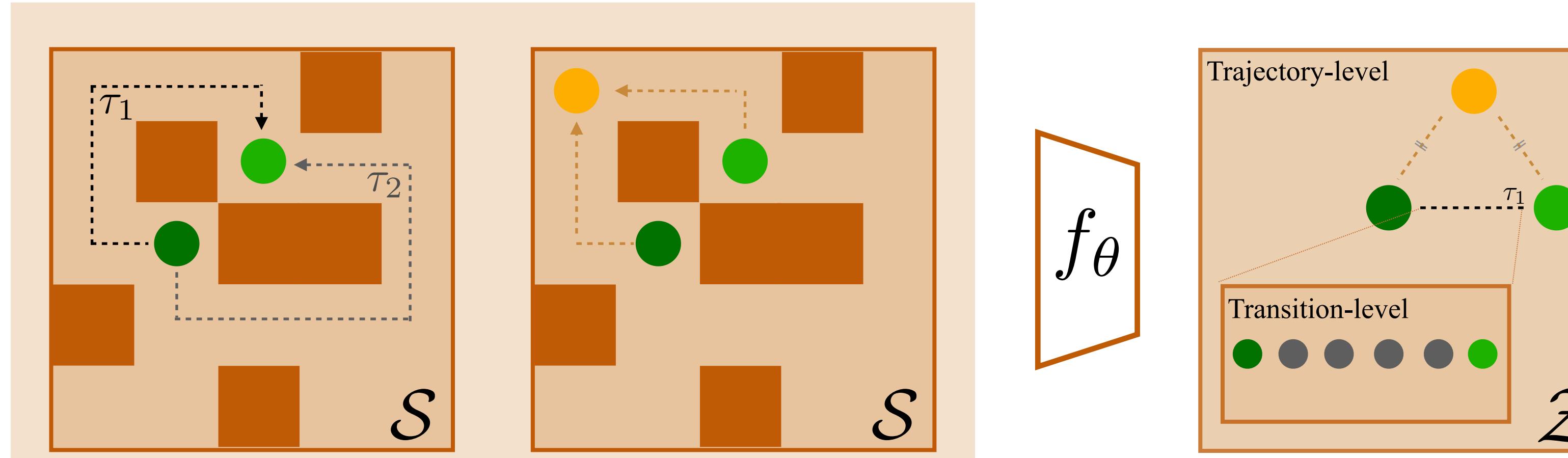
$$\mathcal{L}_{rec} = \arg \min_{\theta} \mathbb{E}_{s \sim \mathcal{D}} \left[ \|s - h \circ f(s; \theta)\| \right]$$

2. Capturing a trajectory-level temporal distance (TD) between any two states

$$\mathcal{L}_{traj} = \mathbb{E}_{\substack{(s, s') \sim \mathcal{D} \\ s_{goal} \sim p_{goal}}} \left[ L_2^\tau \left( \mathcal{B}d - d(f(s; \theta), f(s_{goal}; \theta)) \right) \right]$$

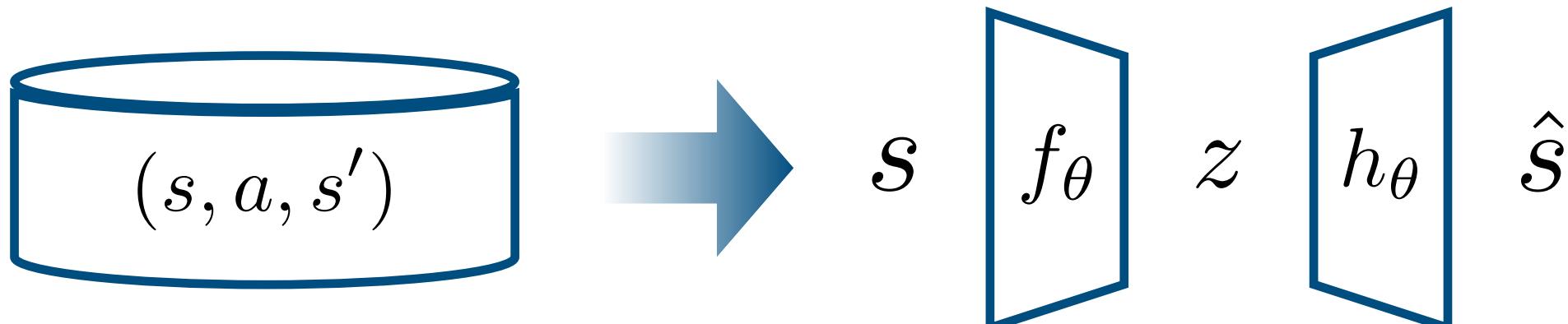
3. Capturing a transition-level temporal distance (TD)

$$\mathcal{L}_{tran} = \mathbb{E}_{(s, s') \sim \mathcal{D}} \left[ L_2^{\tau=1} \left( d(f(s; \theta), f(s'; \theta) - d_0) \right) \right]$$

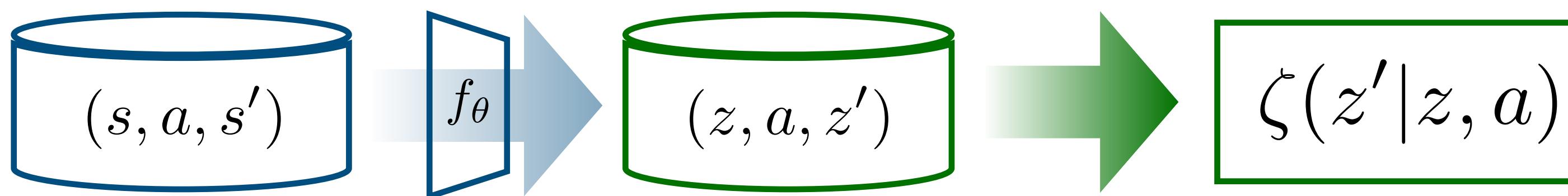


# TempDATA: Temporal distance-aware transition augmentation

**Step 1:** Building temporal distance-aware autoencoder



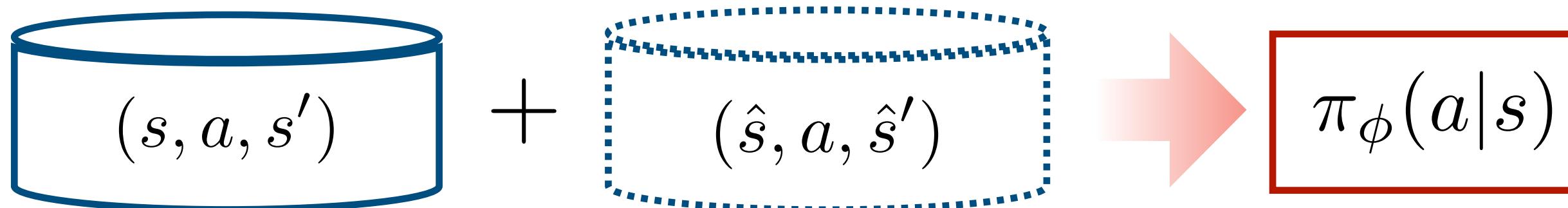
**Step 2:** Building latent dynamic model



**Step 3:** Rolling out synthesized transition



**Step 4:** Extracting policy through off-the-shelf offline RL

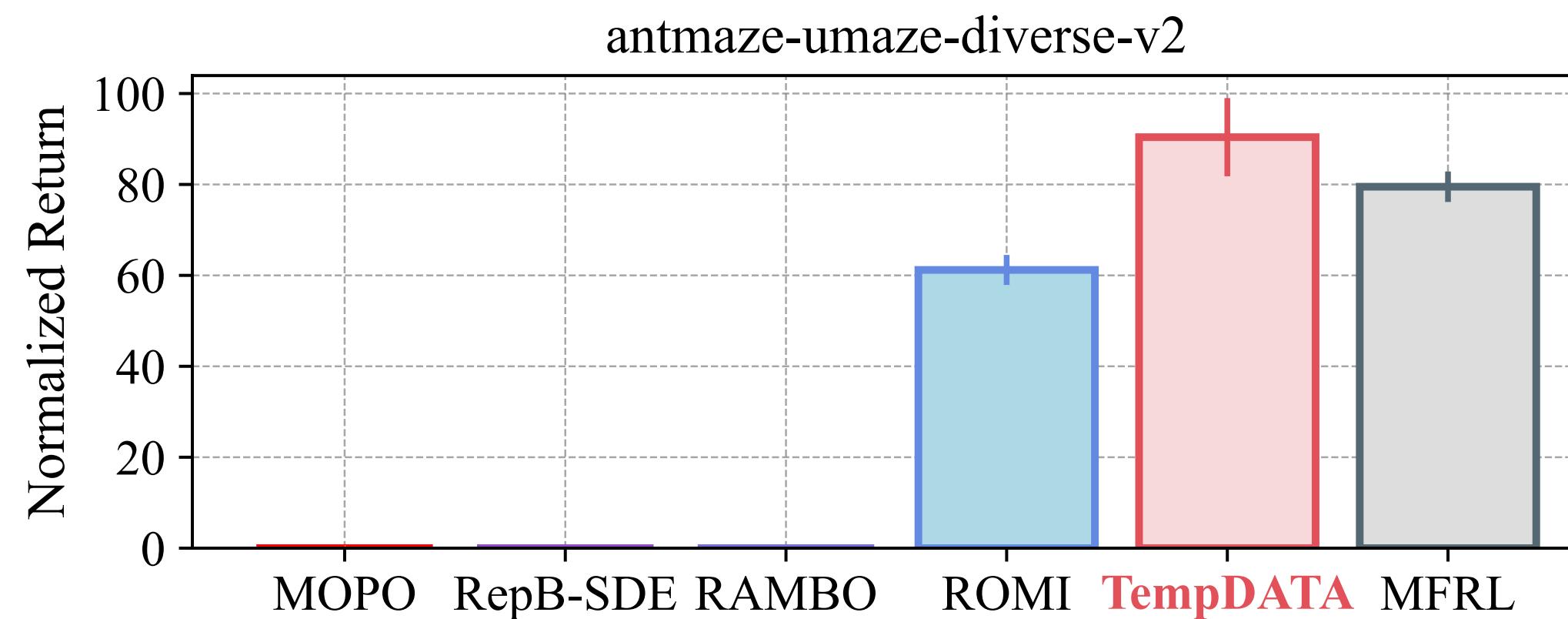
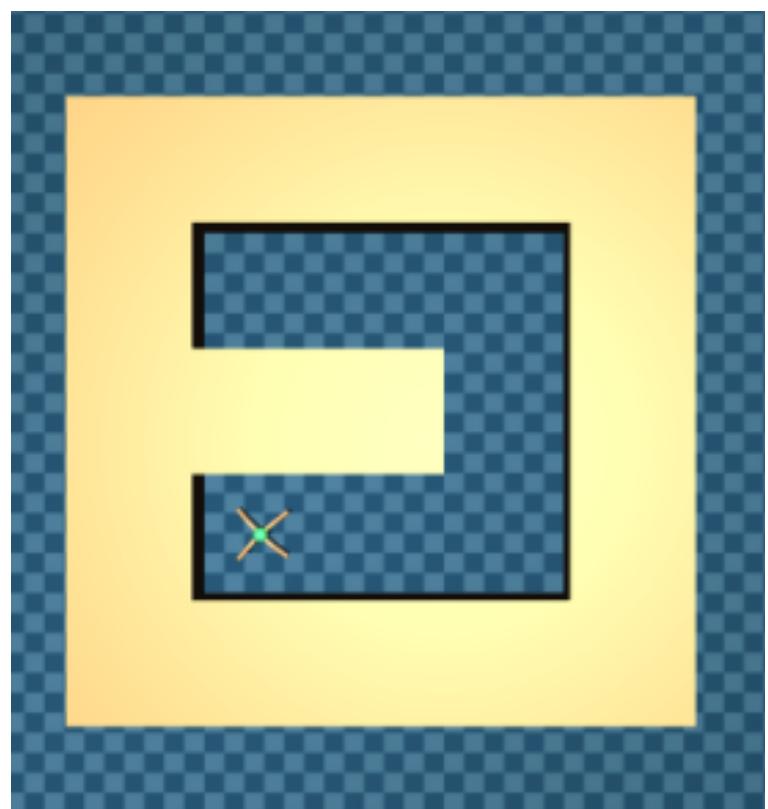


# Closing remarks

**TempDATA** makes offline MBRL more practical in long-horizon sparse reward tasks

- Build a state abstraction capturing trajectory- and transition-level temporal distance
- Roll out a synthesized transition in representation space, instead of raw state space

**Our solution can be adopted in any type (MBRL and MFRL) and technique of RL (Skill RL, GCRL, RL)**



**Project Website**



AntMaze Dataset	TARL-based methods			MBRL-based methods					Proposed
	S4RL <sup>†</sup>	SynthER <sup>†</sup>	GTA <sup>†</sup>	MOPO <sup>†</sup>	RepB-SDE <sup>†</sup>	COMBO <sup>†</sup>	RAMBO <sup>†</sup>	ROMI <sup>†</sup>	
umaze	55.00±21.0	17.1±12.9	66.5±13.8	0.0	0.0	80.3±18.5	25.0±12.0	68.7±2.7	<b>96.3±0.0</b>
umaze-diverse	51.6±23.4	23.9±23.6	57.9±19.0	0.0	0.0	57.3±33.6	0.0	61.2±3.3	<b>90.4±8.6</b>
medium-play	80.9±10.4	41.0±41.2	<b>81.9±8.4</b>	0.0	0.0	0.0	16.4±17.9	35.3±1.3	74.8±8.3
medium-diverse	74.0±19.4	40.1±28.4	<b>78.1±15.8</b>	0.0	0.0	0.0	23.2±14.2	27.3±3.9	69.5±10.8
large-play	42.9±17.4	37.5±13.0	44.4±9.3	0.0	0.0	0.0	0.0	20.2±14.8	<b>56.5±14.1</b>
large-diverse	46.1±16.7	37.5±16.7	47.8±13.4	0.0	0.0	0.0	2.4±3.3	41.2±4.2	<b>44.2±15.3</b>
ultra-play	—	—	—	—	—	0.0±0.0	0.0±0.0	4.9±2.1	<b>53.2±18.2</b>
ultra-diverse	—	—	—	—	—	0.0±0.0	0.0±0.0	8.8±6.8	<b>35.3±10.9</b>
Total score w/o ultra	350.5	236.2	376.5	0.0	0.0	137.6	88.6	253.9	<b>431.7</b>
Total score	—	—	—	—	—	137.6	88.6	272.6	<b>520.2</b>