DISCOVER: Automated Curricula for Sparse-Reward RL

Leander Diaz-Bone ^{*, 1}, Marco Bagatella ^{*, 1, 2}, Jonas Hübotter ^{*, 1}, Andreas Krause ¹



¹ETH Zürich, ²Max Planck Institute for Intelligent Systems

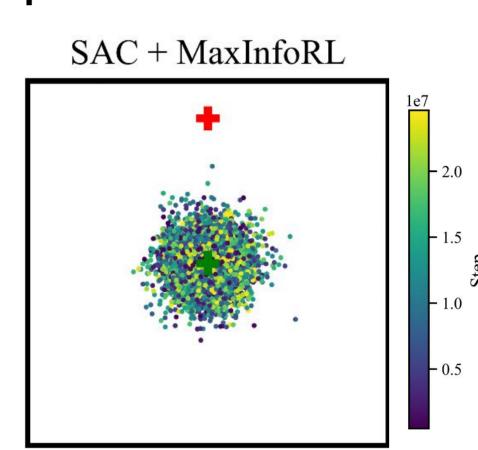


Background

- Exploration in long-horizon, sparse-reward RL is challenging.

 SAC+MaxInfoRL
- Standard RL (+ curiosity)
 does not explore effectively.

How can an agent explore towards an unseen target?



Method

 We propose DISCOVER, which selects exploratory goals trading off achievability, novelty & relevance to the target:

$$g_t = rg \max_{g \in \mathcal{G}_{\mathsf{ach}}} \ lpha_t igg[\underbrace{V(s_0, g)}_{\mathsf{Achievability}} + eta_t \underbrace{\sigma(s_0, g)}_{\mathsf{Novelty}} igg] + \left(1 - lpha_t\right) igg[\underbrace{V(g, g^\star) + eta_t \, \sigma(g, g^\star)}_{\mathsf{Relevance}} igg]$$

- All components are estimated using an ensemble of goal-conditioned value functions, learning how goals relate.
- We automatically adapt α to match 50% goal achievement rate.
- We connect DISCOVER to principled exploration in bandits, bounding the target task achieval time:

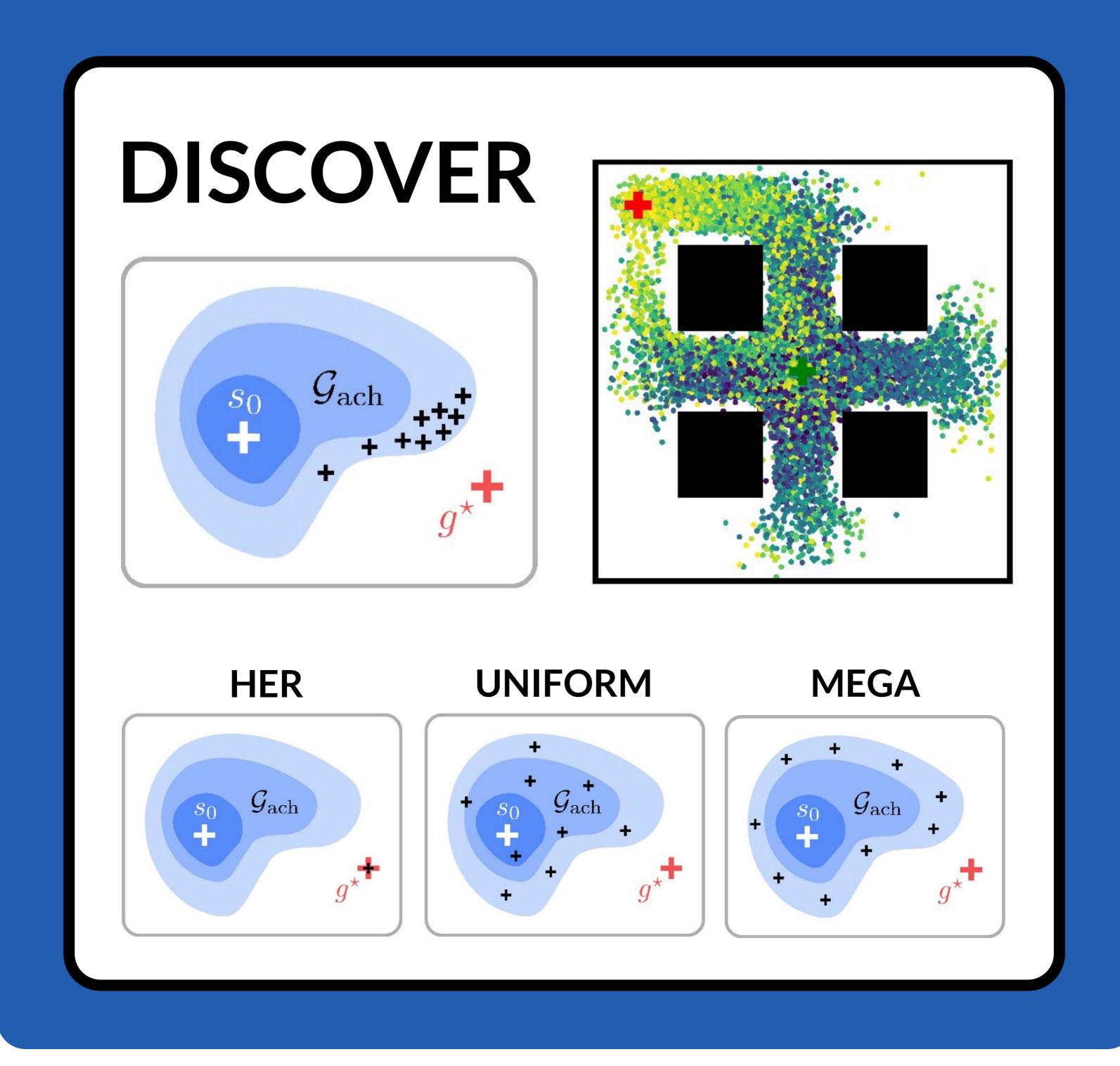
$$N \le \widetilde{O}\left(\frac{Dd^2}{\kappa^3}\right)$$

N: Episodes until g^* becomes reachable $D = V^*(s_0, g^*)$: Optimal distance between s_0 and g^* κ : Expansion rate

d: Feature dimensionality

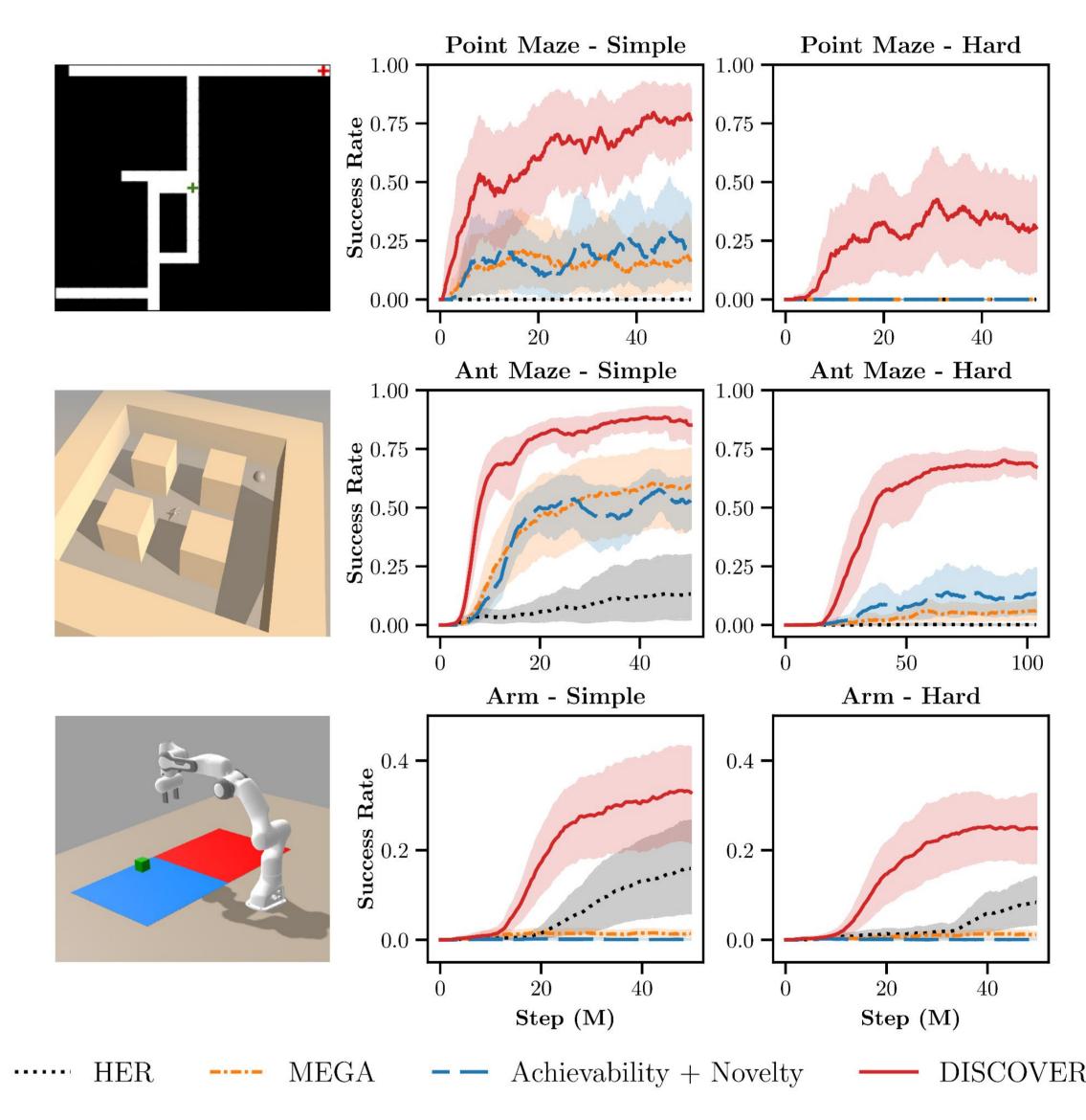
Long-horizon, sparse-reward RL tasks are challenging.

Only agents that direct exploration can solve them.



Results

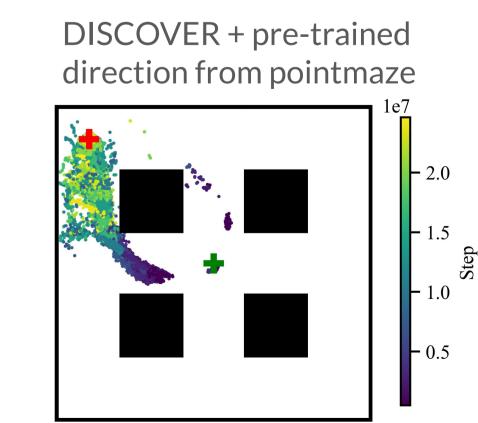
 DISCOVER outperforms baselines on complex sparse-reward control tasks.



• DISCOVER effectively directs exploration.

Dimension23456HER ∞ ∞ ∞ ∞ ∞ MEGA4.8 ∞ ∞ ∞ ∞ Ach. + Nov.5.2 ∞ ∞ ∞ DISCOVER2.93.17.45.418.7

Number of Steps (M) required to reach 10%



Components of DISCOVER after 8M steps (antmaze-hard)

