



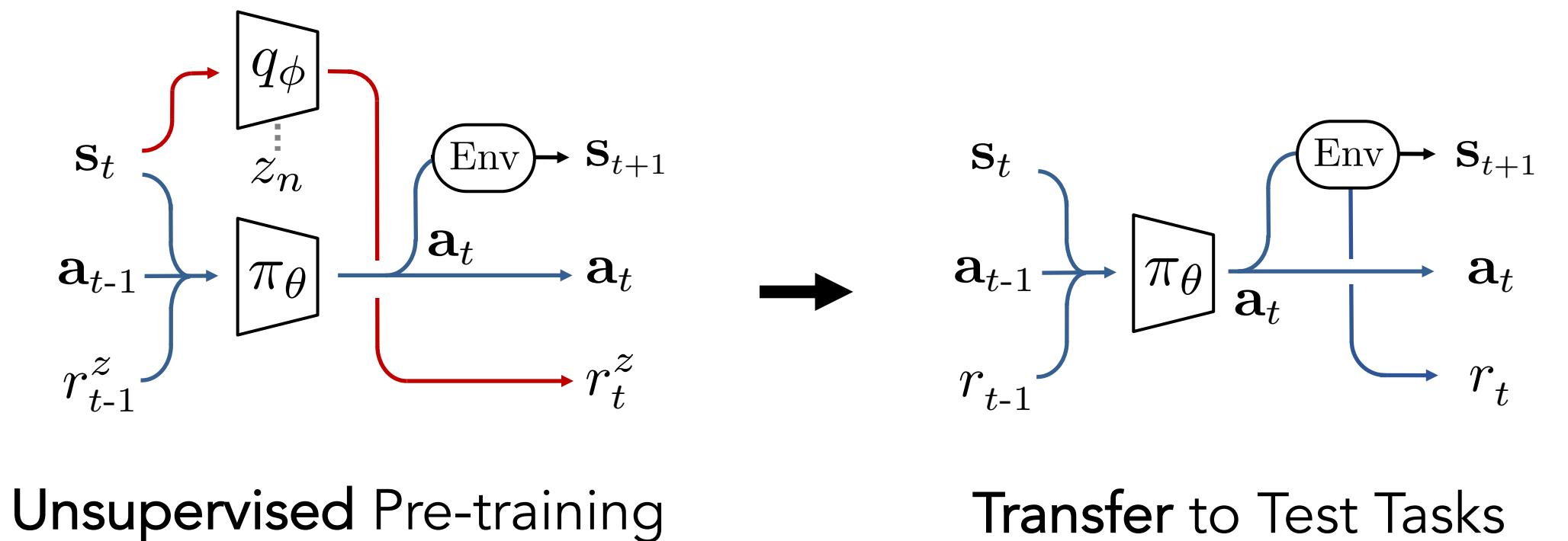
Unsupervised Curricula for Visual Meta-Reinforcement Learning

Allan Jabri, Kyle Hsu, Ben Eysenbach, Abhishek Gupta, Sergey Levine, Chelsea Finn

Motivation

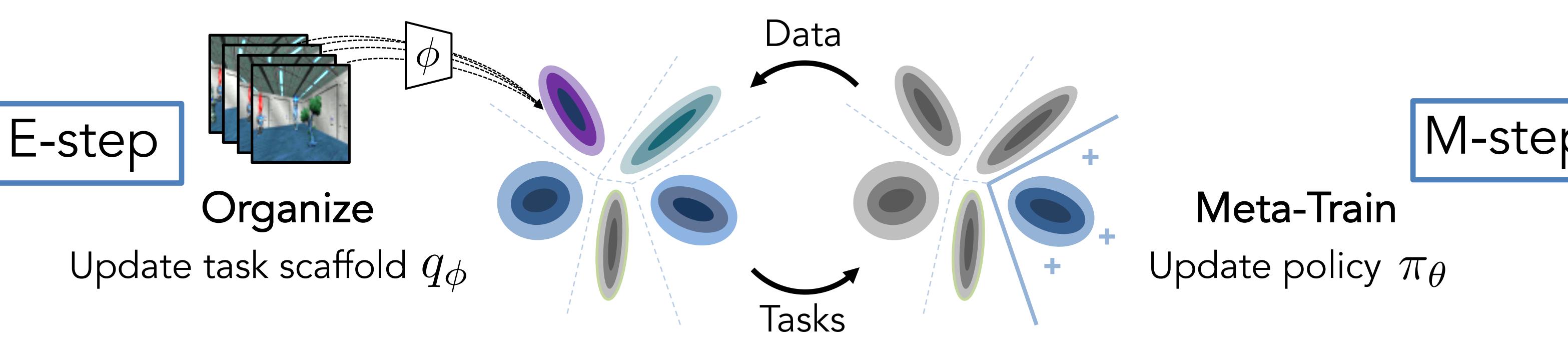
Meta-RL provides a formulation for transfer between different tasks, but relies on specification of a task distribution

→ Can useful meta-RL strategies be acquired by pretraining with task distribution formed through unsupervised interaction?



Method

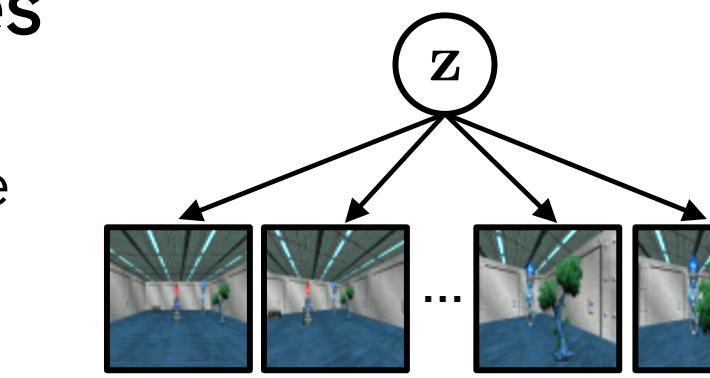
Main Idea: Jointly discover and meta-learn tasks by alternating between task acquisition and meta-reinforcement learning of updated task distribution.



(E) Trajectory-level generative model with discriminative features

$$\max_{\phi} \mathbb{E}_{z \sim q_\phi(z), \tau \sim \mathcal{D}} [\log q_\phi(\tau|z)]$$

w/ conditional independence
 $q_\phi(s) = \sum_z q_\phi(s|z)p(z)$



(M) Policy meta-learns task distribution

$$\max_{\theta} \mathbb{E}_{z \sim q_\phi(z), s \sim \pi_\theta(s|z)} [\log q_\phi(s|z) - \log q_\phi(s)]$$

i.e. learn tasks indexed by z
 $r_z(s) = \log q_\phi(s|z) - \log q_\phi(s)$

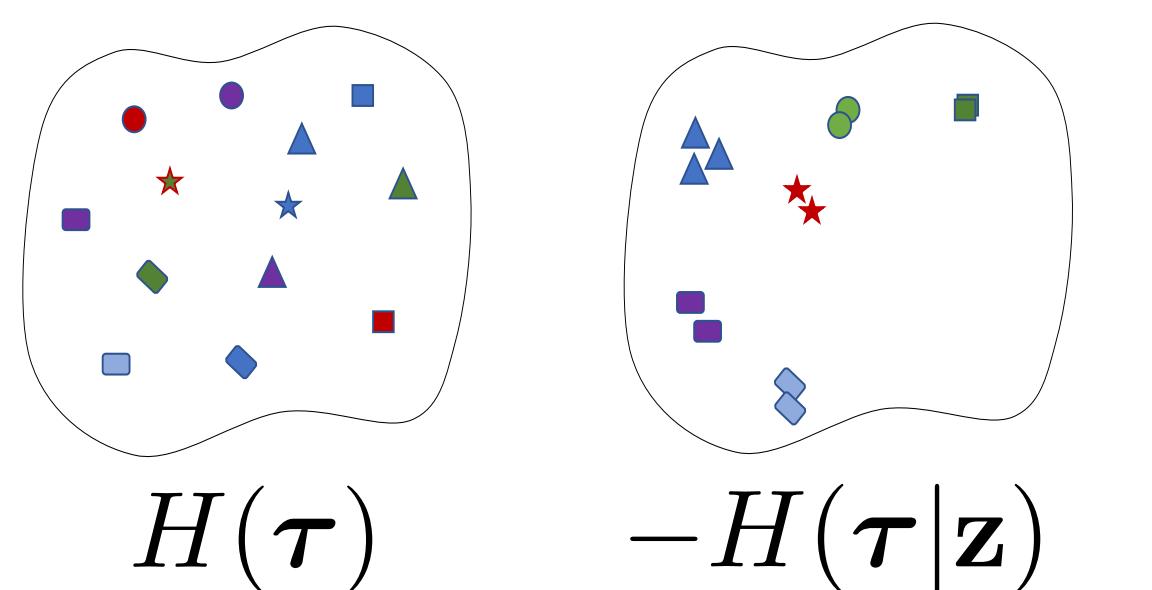
Key Challenges

1. Task acquisition with high-dimensional observations, i.e. pixels.
2. Adaptation of task distribution and meta-learner, i.e. curriculum.

Formulation

Intuition: Prepare meta-RL agent with task inference and execution strategies that transfer

→ Task dist. should balance Diversity and Structure



$$\max_{\theta, \phi} I(\tau; z)$$

π_θ Meta-RL policy

q_ϕ Learned task scaffold

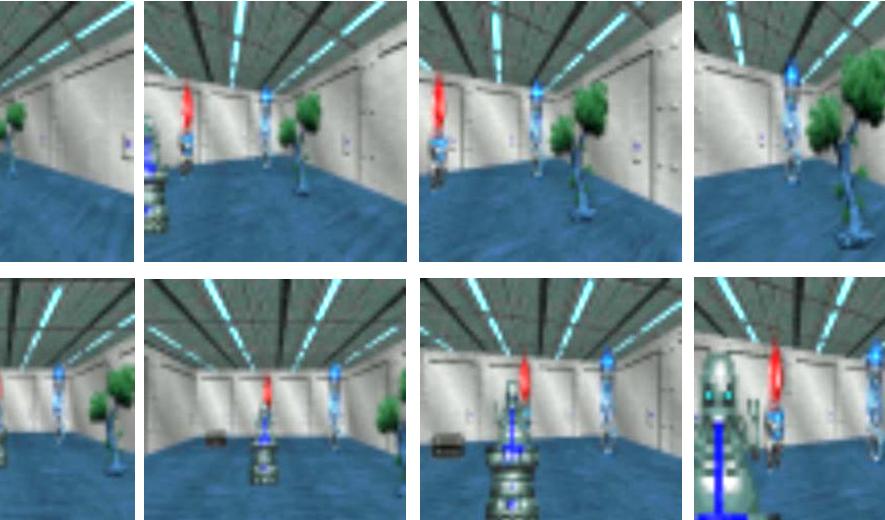
τ Post-update trajectories

z Task latent variable

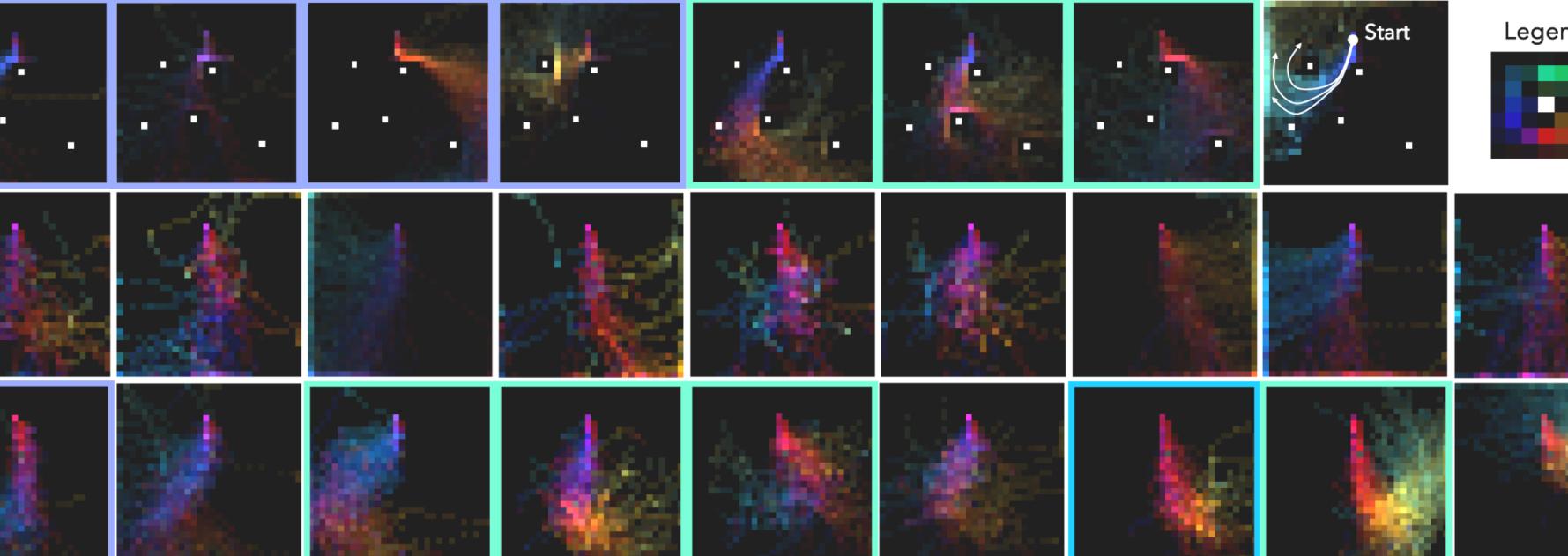
Information Maximization
via Variational EM

Environments and Discovered Tasks

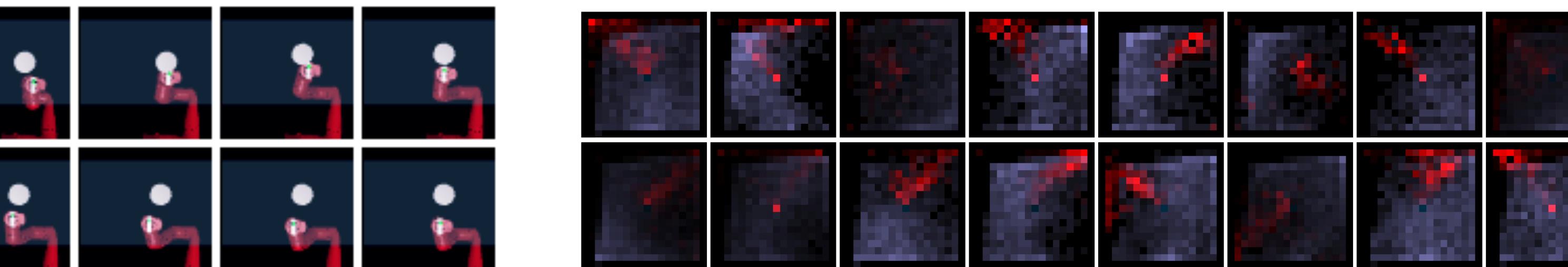
Visual Navigation in ViZDoom



Visualization of Discovered Task Distributions



Visuomotor Control w/ MuJoCo Sawyer

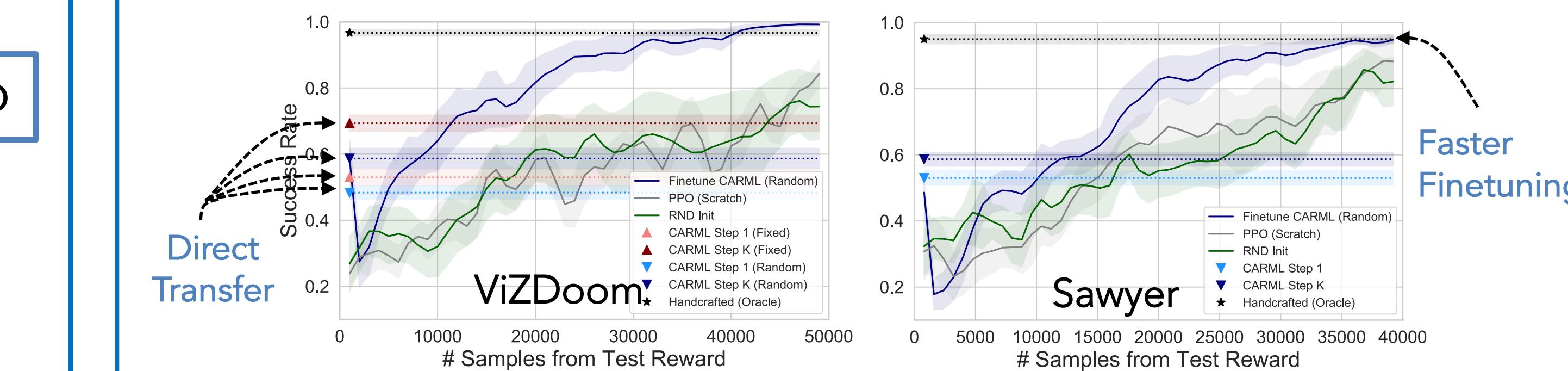


Experiments

Do acquired meta-RL strategies transfer to test task distributions?

Direct Transfer: Apply policy to test tasks without finetuning

Finetuning: Update parameters of policy to specific test tasks

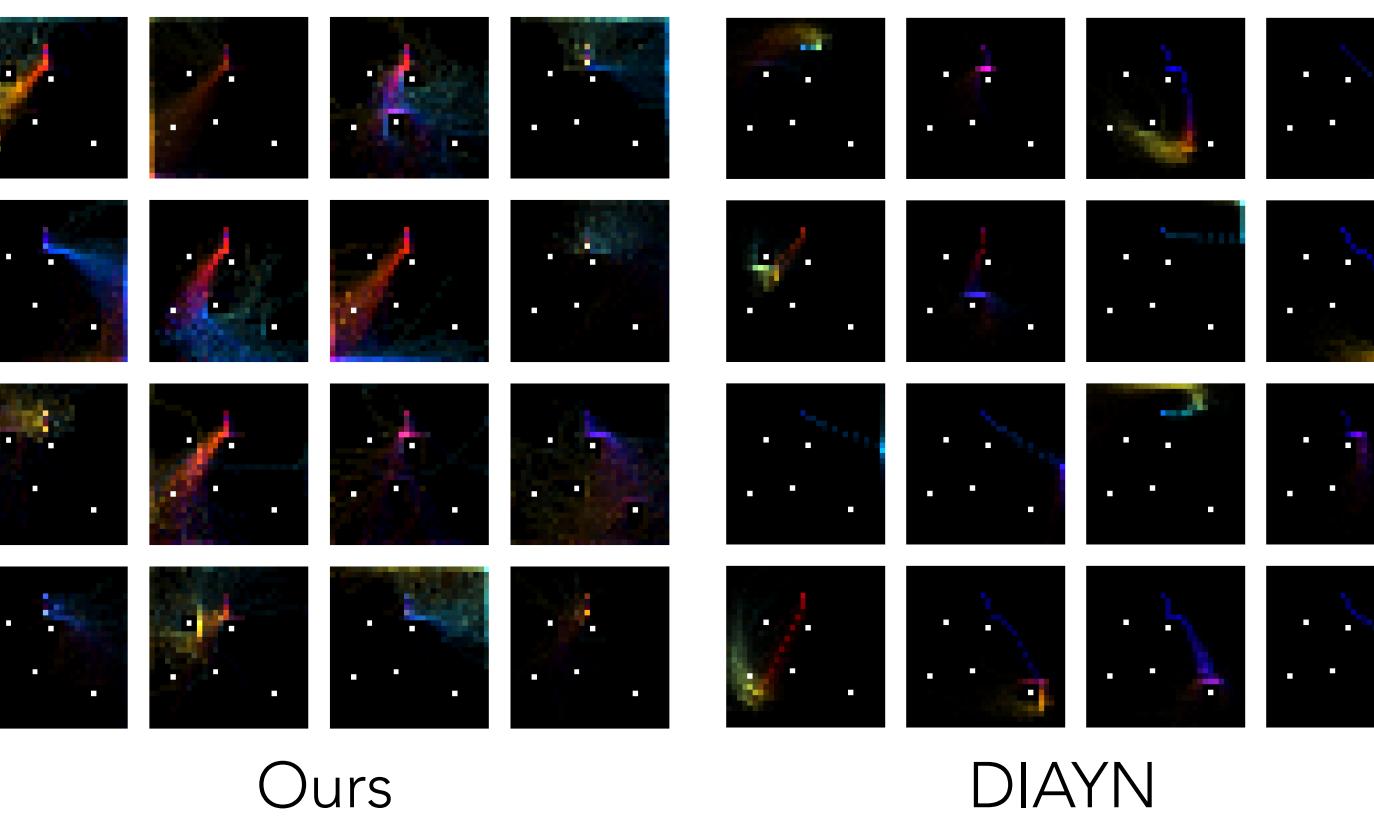


Comparing Variants: Effect of Task Acquisition and Joint Optimization

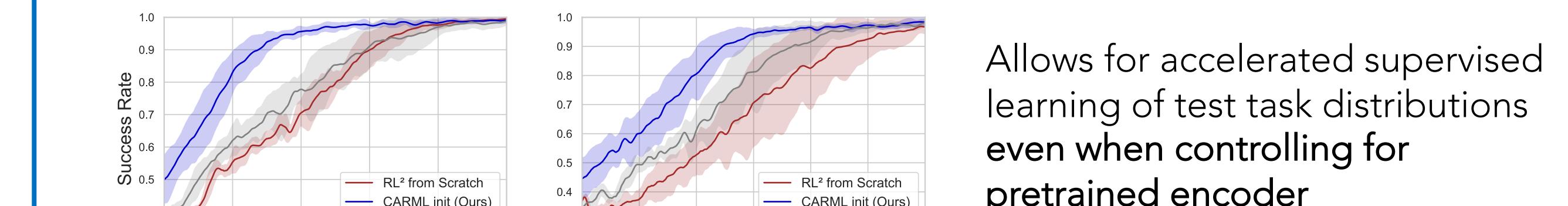
Online disc.: Discriminative q_ϕ (DIAYN)

Pipelined CARML: 1 E-step, 1 M-step

Avoiding mode collapse in task distribution



Meta-Pretraining: Accelerated Supervised Meta-RL



Allows for accelerated supervised learning of test task distributions even when controlling for pretrained encoder

Takeaways

- Proposed task acquisition strikes balance between generative and discriminative approaches for modeling pixel trajectories, avoiding task mode-collapse
- Task distribution supports meta-RL of strategies that transfer to related but unseen hand-crafted tasks
- Adapting tasks with the policy eases meta-learning with curriculum
- Need more task-specific bias for direct transfer, more complex tasks