

Unsupervised Curricula for Visual Meta-Reinforcement Learning

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

NeurIPS 2019

A Need to Generalize

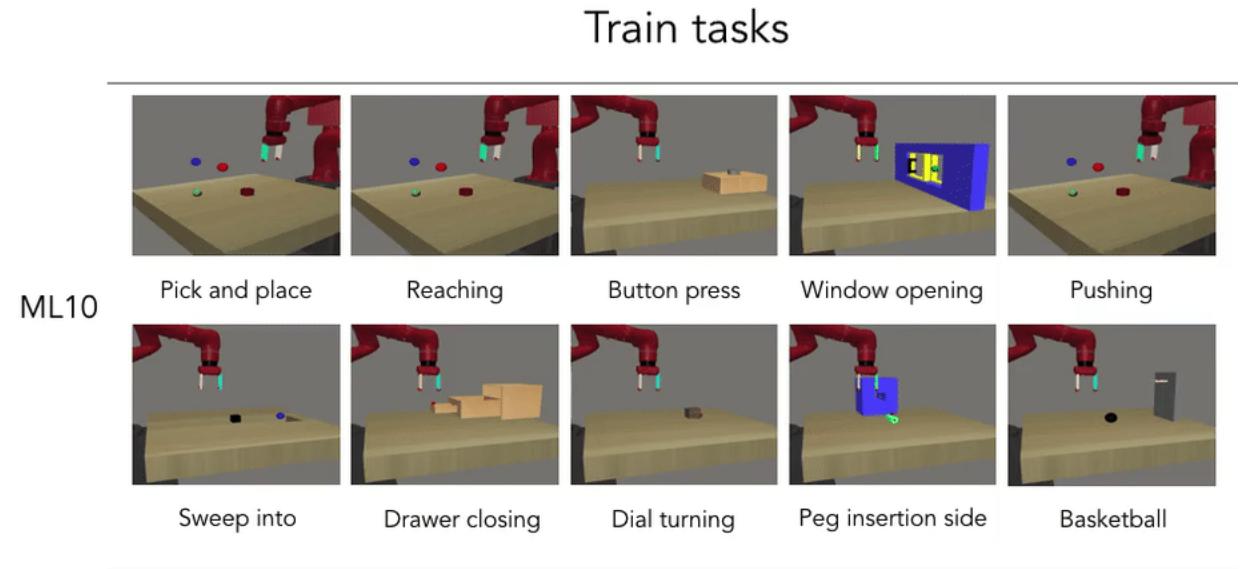
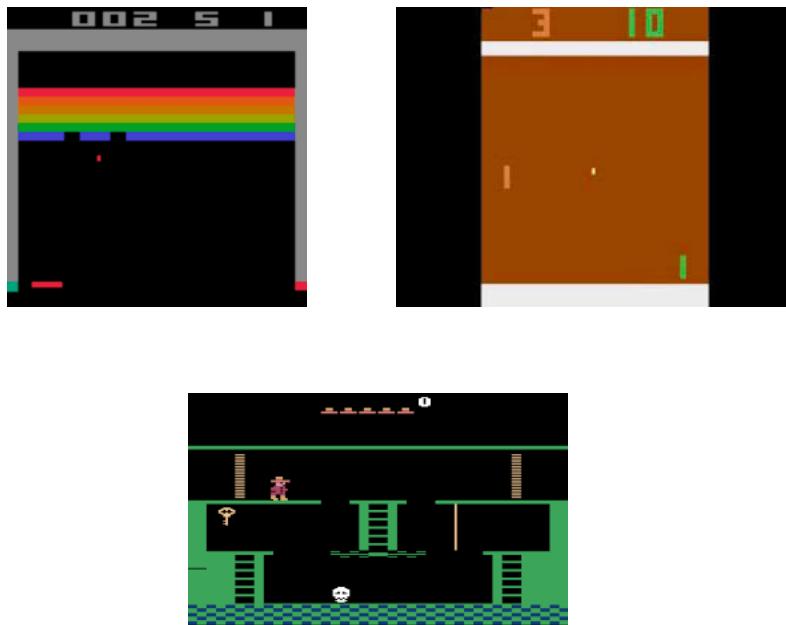


C. elegans, foraging



New Caledonian Crow
applying tools

From Specialist to Generalist



Source: Meta-World
meta-world.github.io

Multi-task Reinforcement Learning

Contextual Policies

$$\pi(a|o, z)$$

Task description is given

e.g. a goal



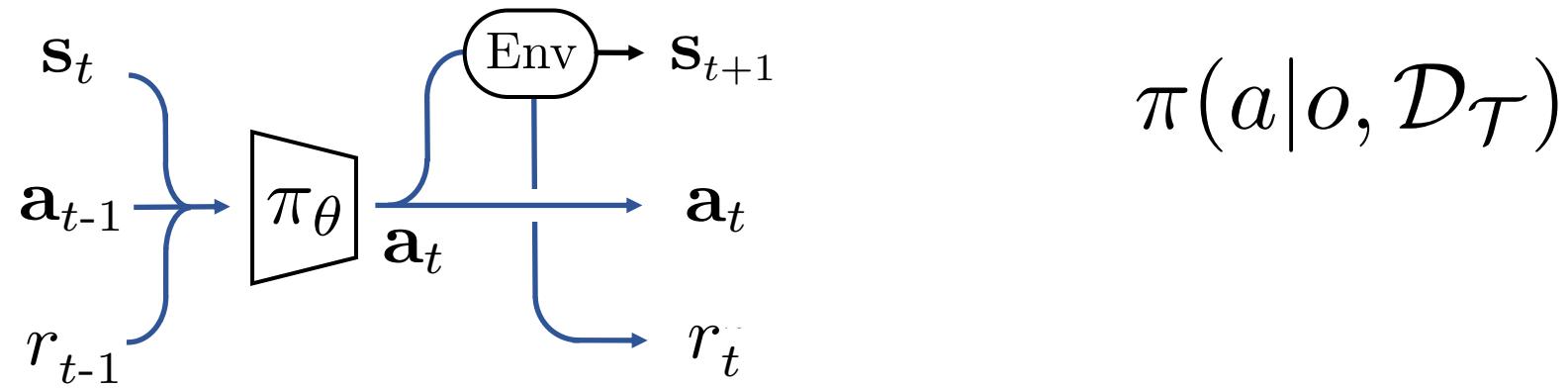
Meta-learning for RL

$$\pi(a|o, \mathcal{D}_{\mathcal{T}})$$

Task inferred from data
collected by policy

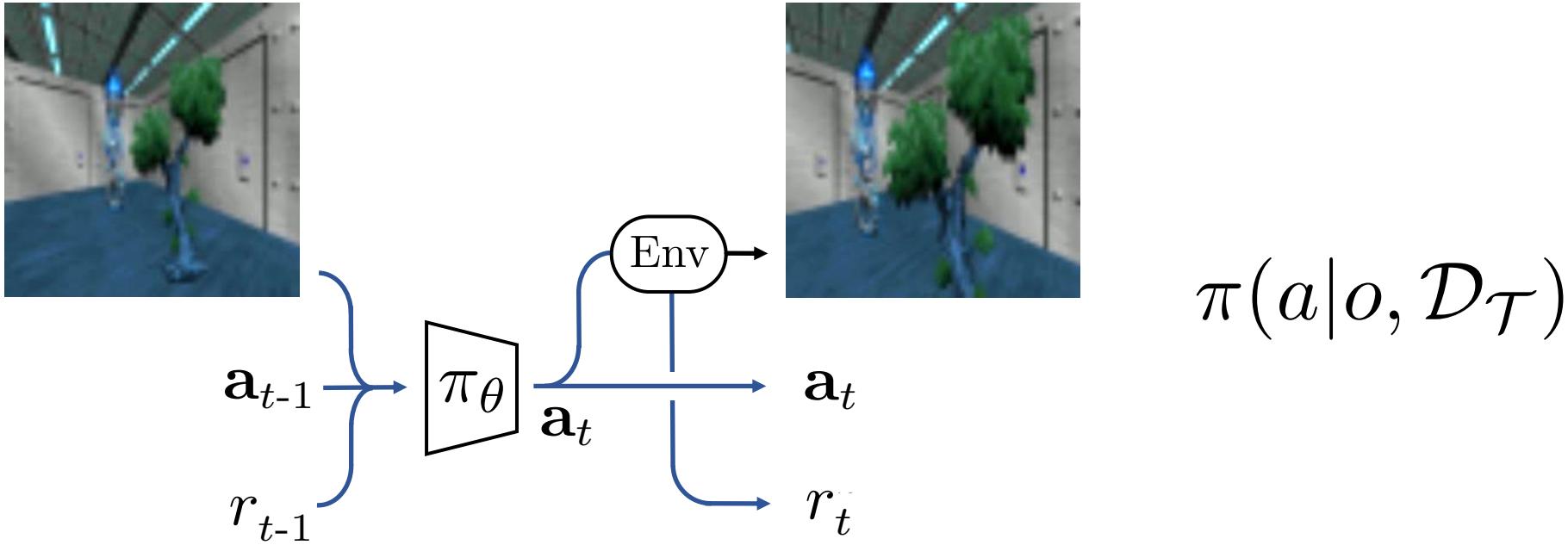
more general

Meta-Reinforcement-Learning



Recurrent policy learns to **infer task** by collecting the right data

Visual Meta-Reinforcement-Learning

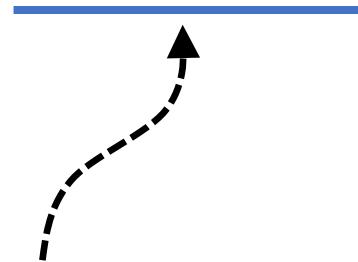


Search for and associate **stimulus** and **reward**.

The Task Distribution

$$\arg \max_{\theta} \sum_{i=1}^n \mathbb{E}_{\pi_{\theta}(\mathcal{D}_{\mathcal{M}_i})}[R(\tau)]$$

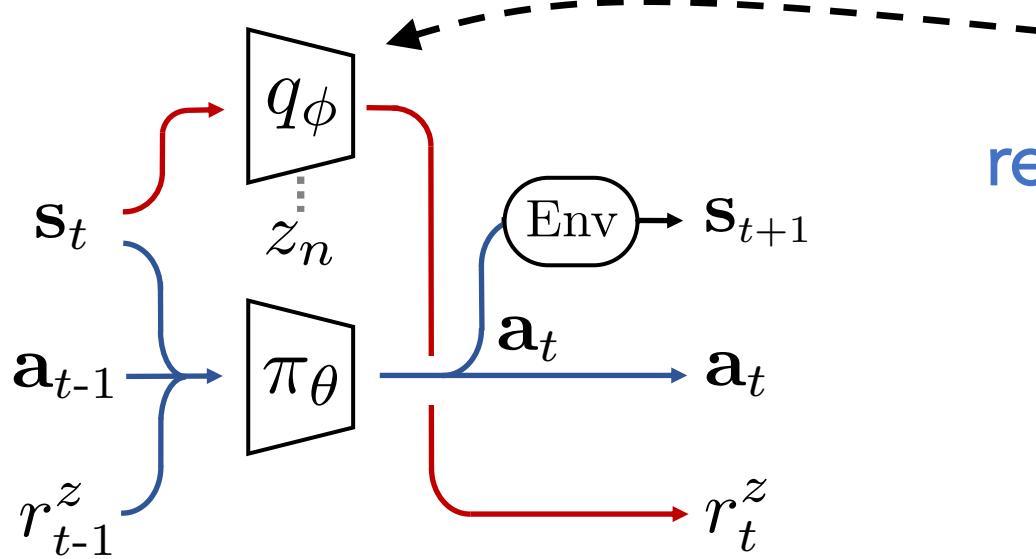
where $\mathcal{M}_i \sim p(\mathcal{M})$



Meta-training tasks give rise to
task inference and execution strategies

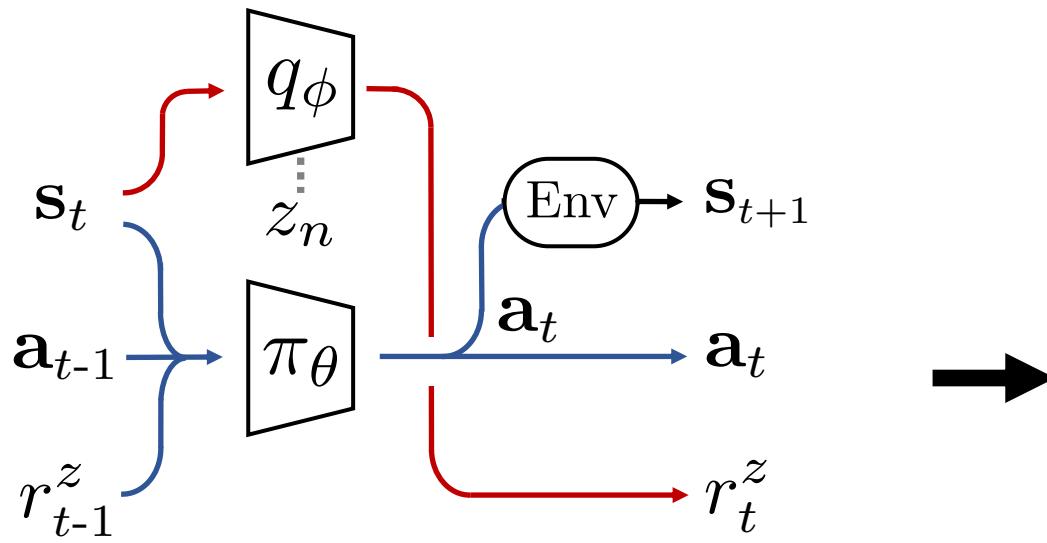
Can we learn **useful** meta-RL strategies
with tasks formed **without supervision**?

“Meta-Pre-training”

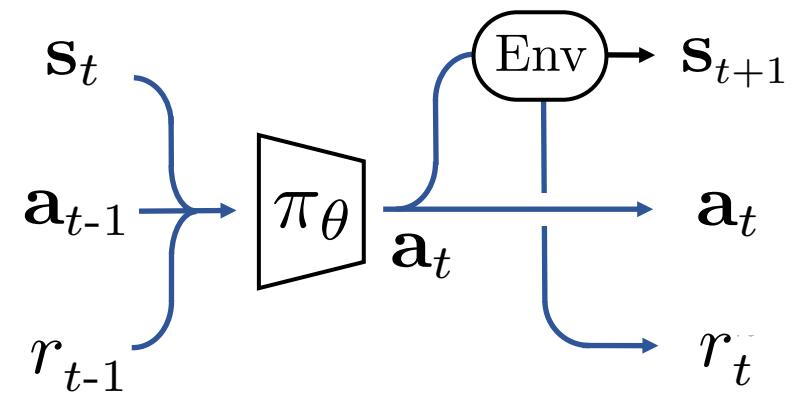


learned model providing
reward functions for meta-learning

“Meta-Pre-training”



Unsupervised Pre-training



Transfer to Test Tasks

Task Acquisition

Unsupervised discovery of tasks

Tasks
→

Meta-learning

Learn to learn to solve tasks

Task Acquisition

Unsupervised discovery of tasks

Tasks

.....

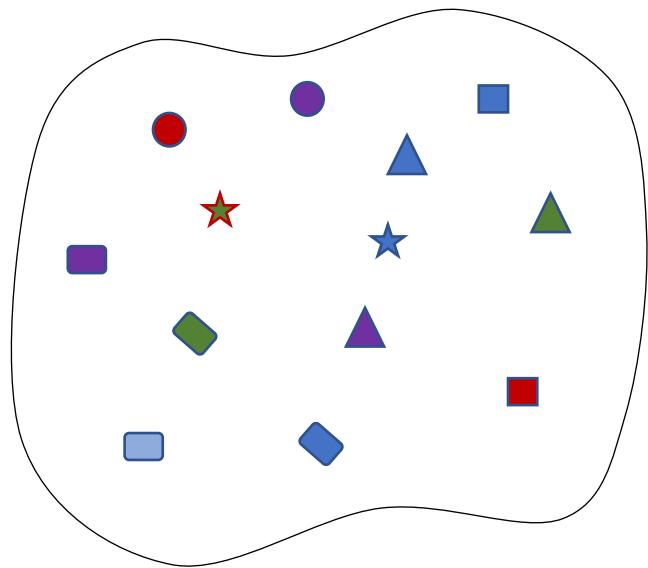
Data?

Meta-learning

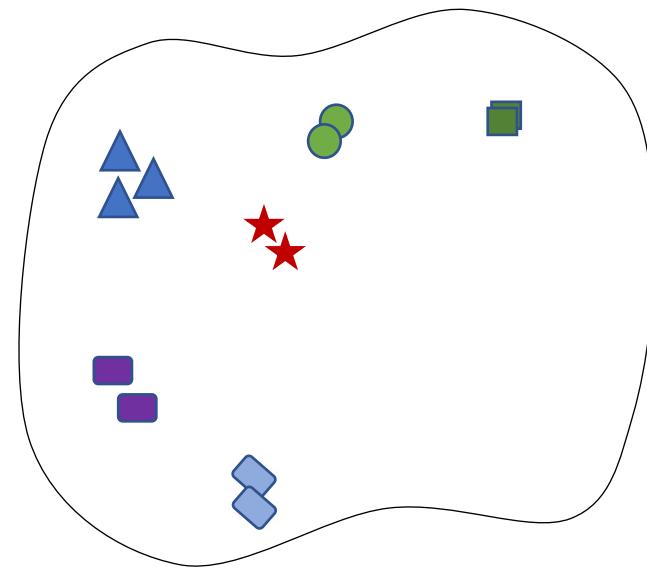
Learn to learn to solve tasks

Should co-adapt

Criteria for Task Distribution

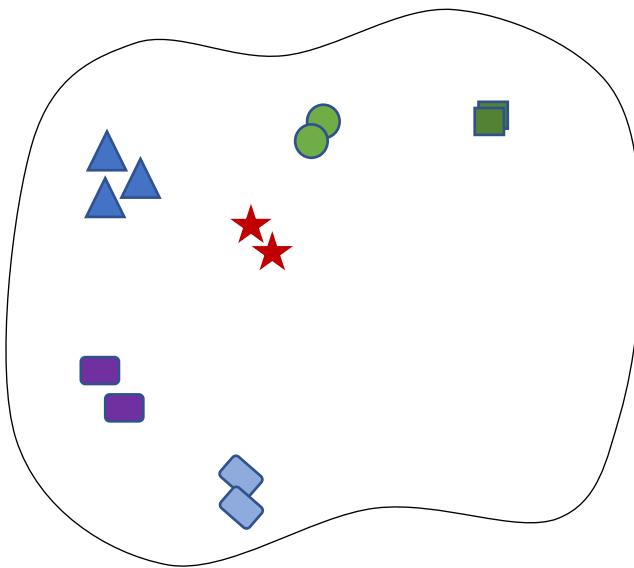
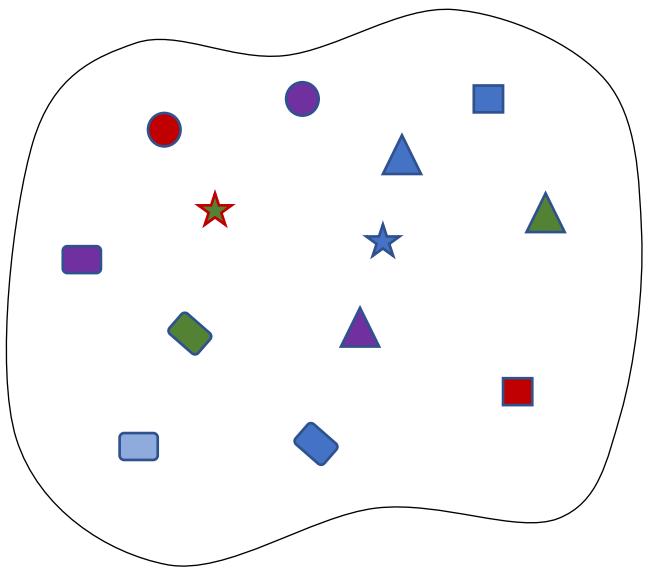


Diversity



Structure

Criteria for Task Distribution



$$\text{Diversity } H(\boldsymbol{\tau}) - H(\boldsymbol{\tau}|\mathbf{z}) \text{ Structure} \\ = I(\boldsymbol{\tau}; \mathbf{z})$$

Formulation

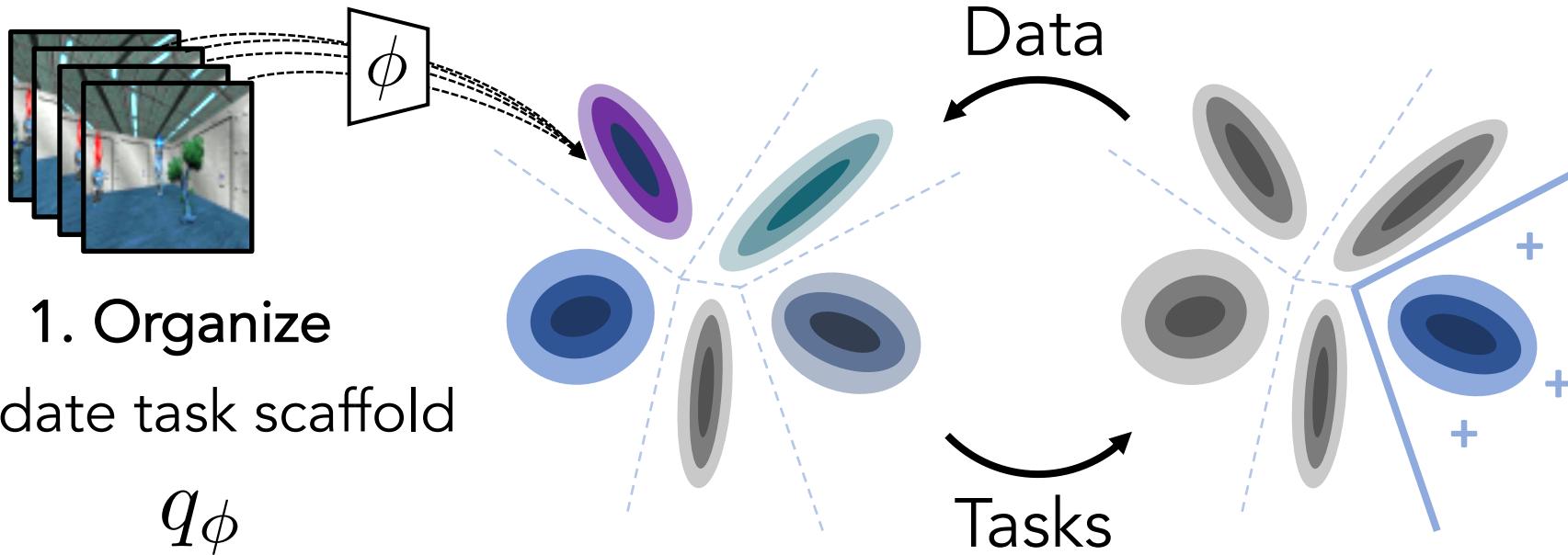
$$\max_{\theta, \phi} I(\tau; \mathbf{z})$$

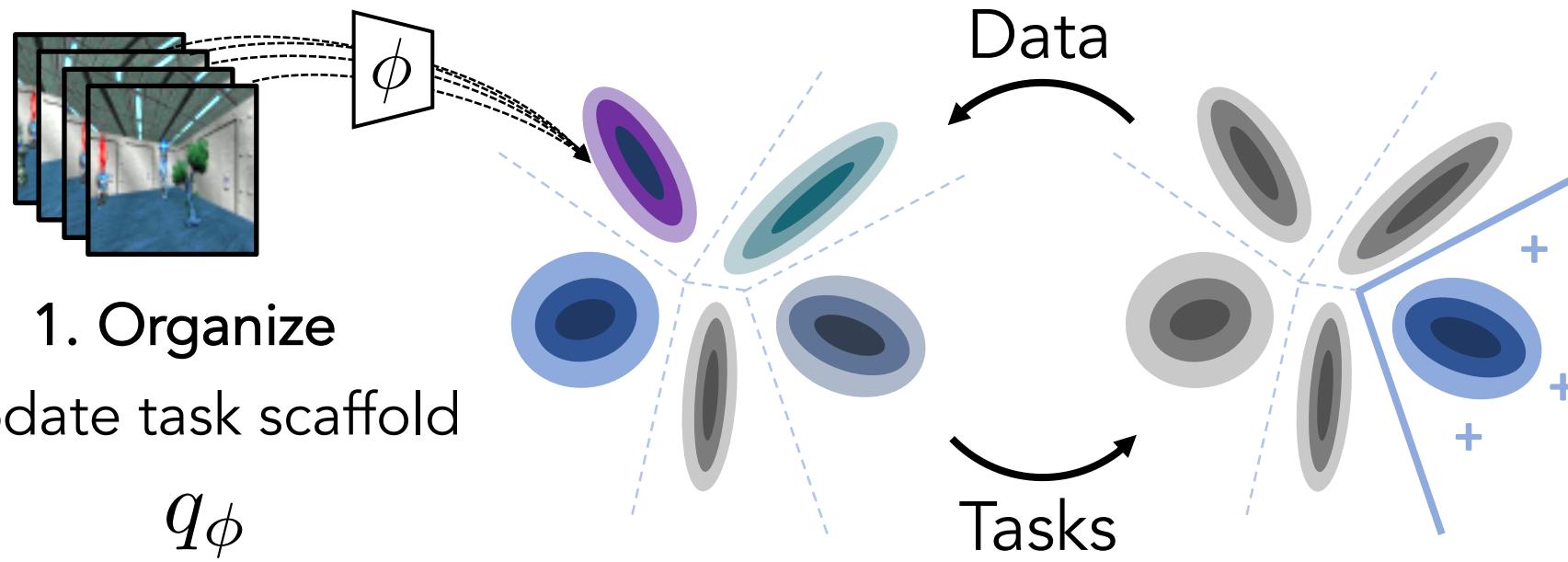
Policy π_θ

τ Post-update trajectories

Task scaffold q_ϕ

\mathbf{z} Task latent variable

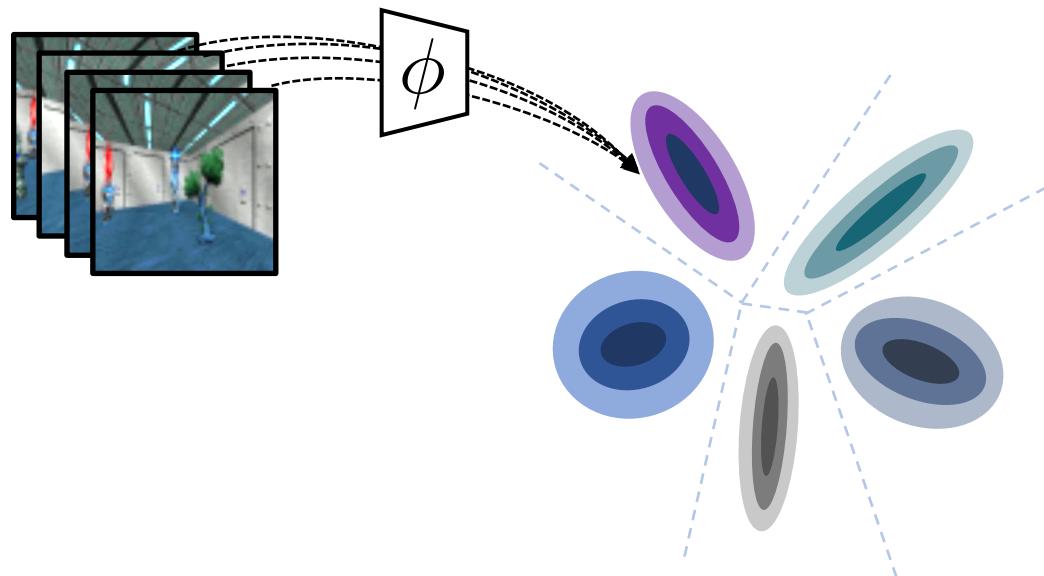




E-step

M-step

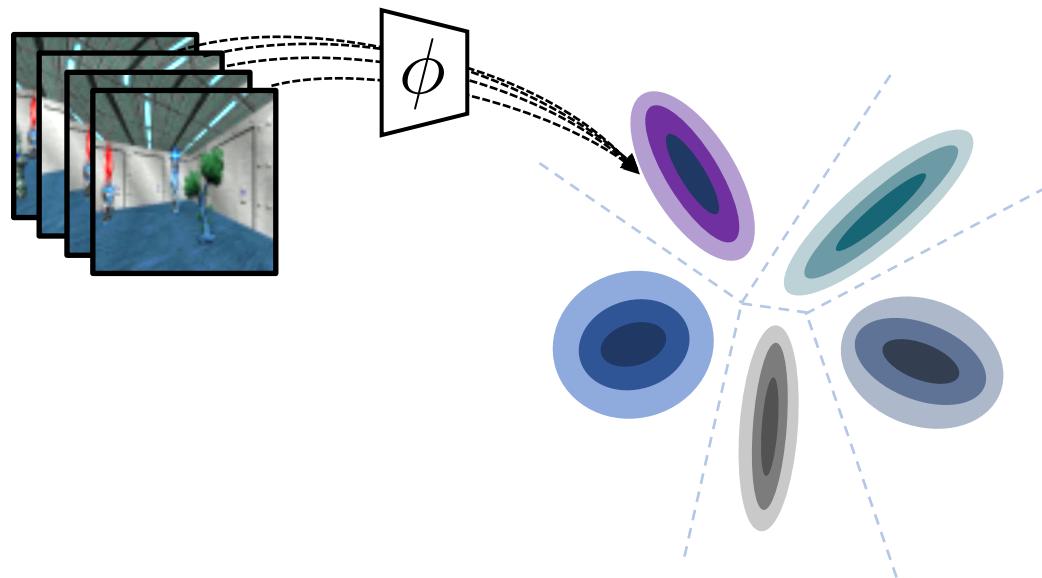
Task Acquisition (E-step)



Trajectory-level **discriminative clustering**

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

Task Acquisition (E-step)

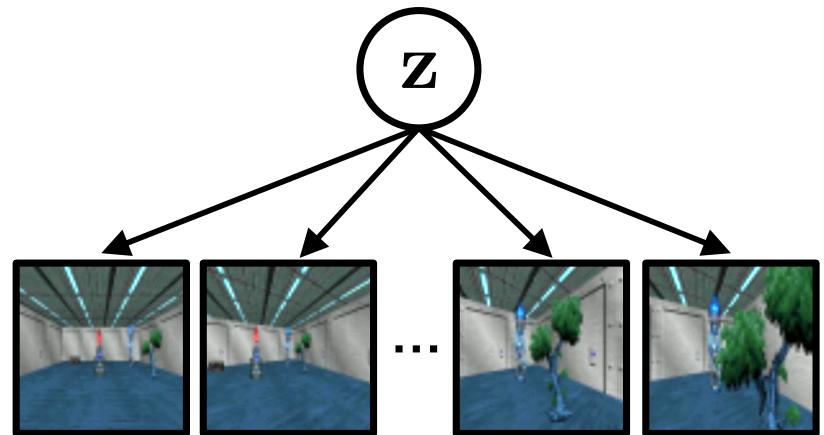


Trajectory-level **discriminative clustering**

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

Conditional independence assumption

$$q_{\phi}(\mathbf{s}) = \sum_{\mathbf{z}} q_{\phi}(\mathbf{s} | \mathbf{z}) p(\mathbf{z})$$

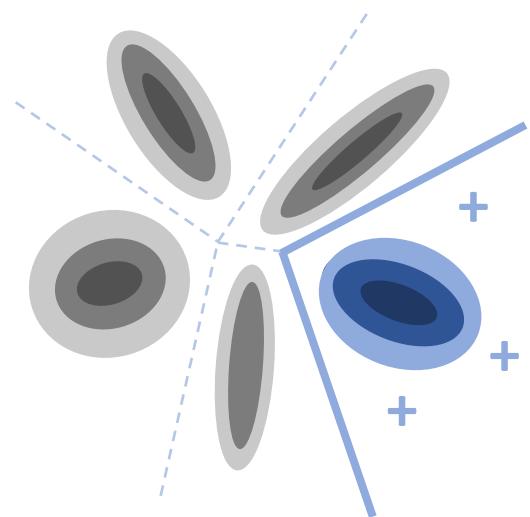


Meta-train (M-step)

Policy **learns updated task distribution**

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

i.e. $r_{\mathbf{z}}(\mathbf{s}) = \log q_{\phi}(\mathbf{s}|\mathbf{z}) - \log q_{\phi}(\mathbf{s})$



Meta-train (M-step)

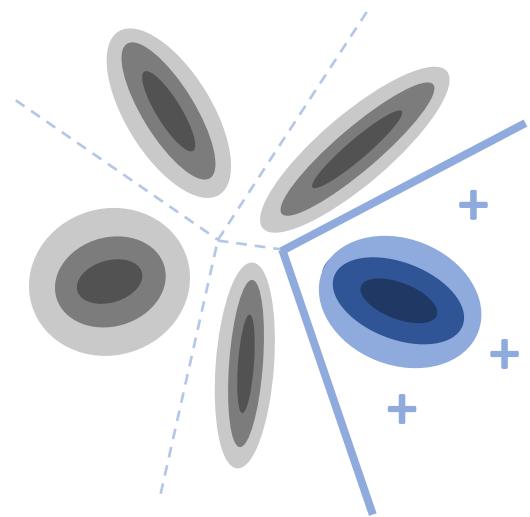
Policy **learns updated task distribution**

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

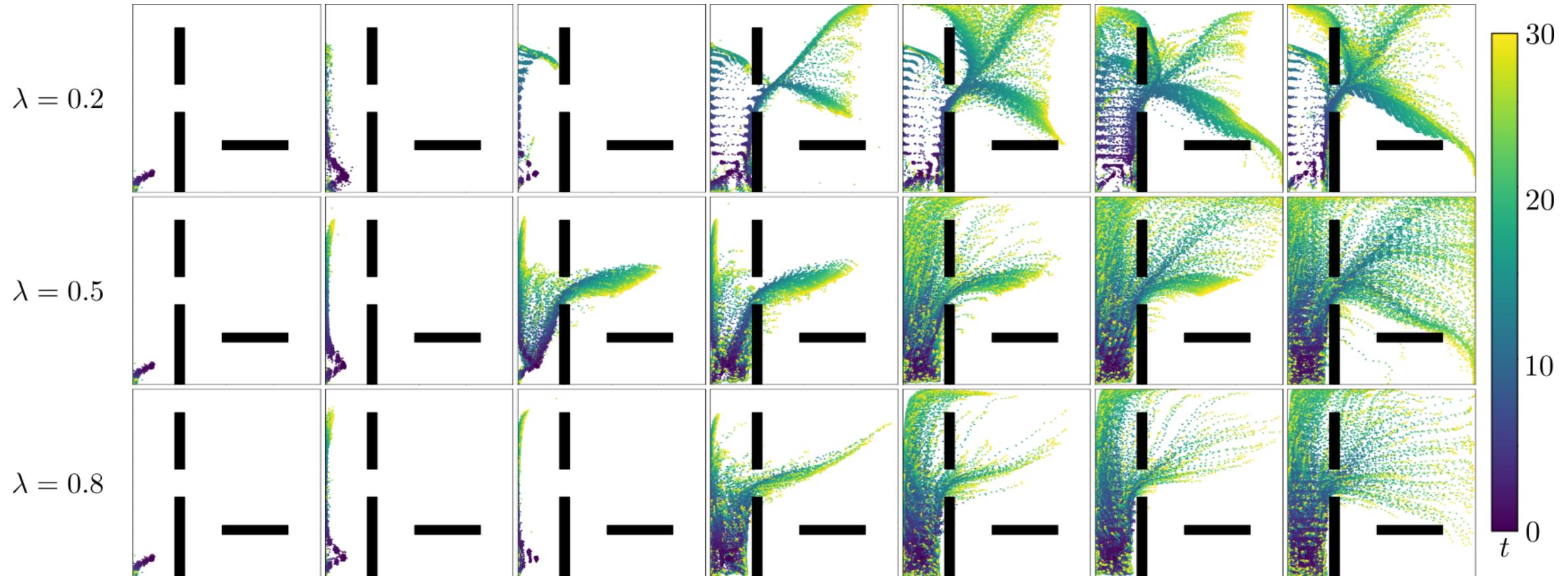
i.e. $r_{\mathbf{z}}(\mathbf{s}) = \log q_{\phi}(\mathbf{s}|\mathbf{z}) - \log q_{\phi}(\mathbf{s})$

Density-based Exploration

$$\begin{aligned} r_{\mathbf{z}}(\mathbf{s}) &= \lambda \log q_{\phi}(\mathbf{s}|\mathbf{z}) - \log q_{\phi}(\mathbf{s}) \\ &= \underline{\lambda - 1} \log q_{\phi}(\mathbf{s}|\mathbf{z}) + \log q_{\phi}(\mathbf{z}|\mathbf{s}) + C \end{aligned}$$

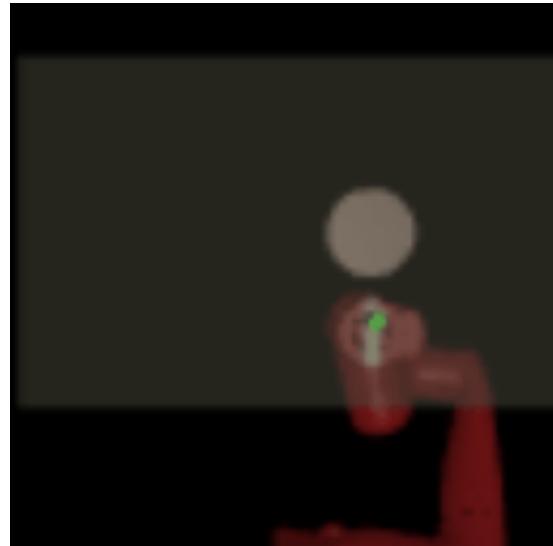
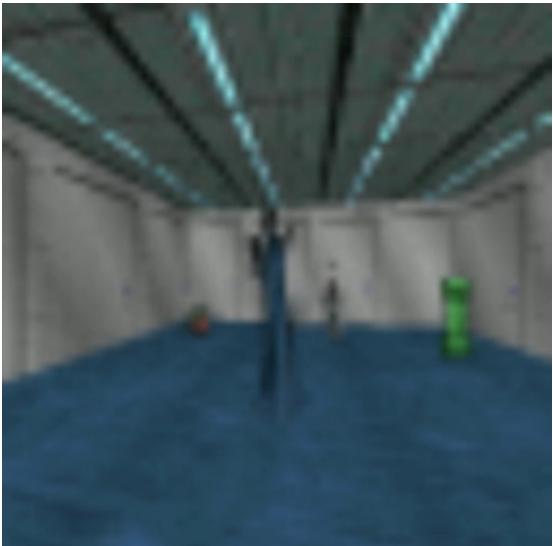


Effect of λ



Structure v. Diversity

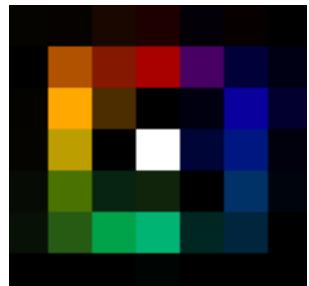
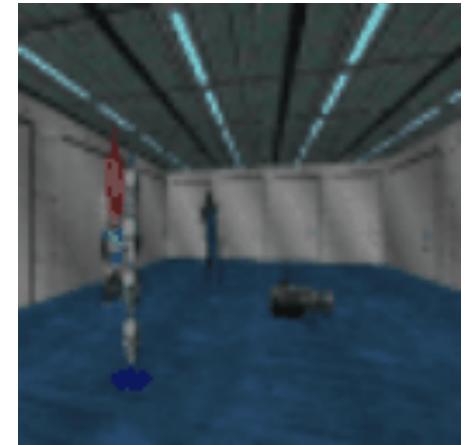
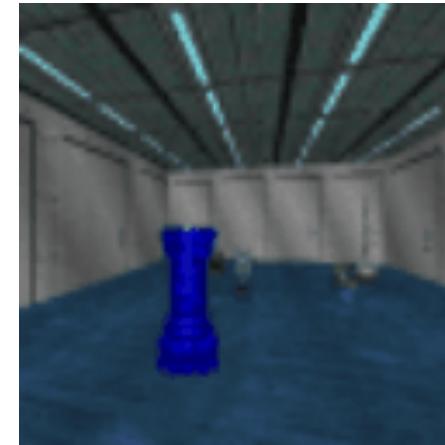
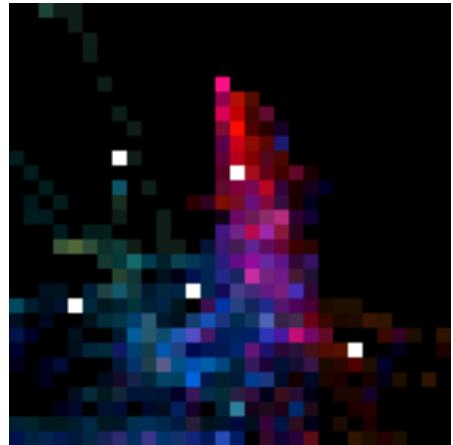
Experimental Setting



Visual Navigation
in VizDoom

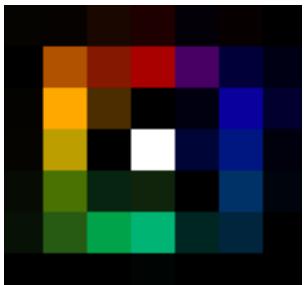
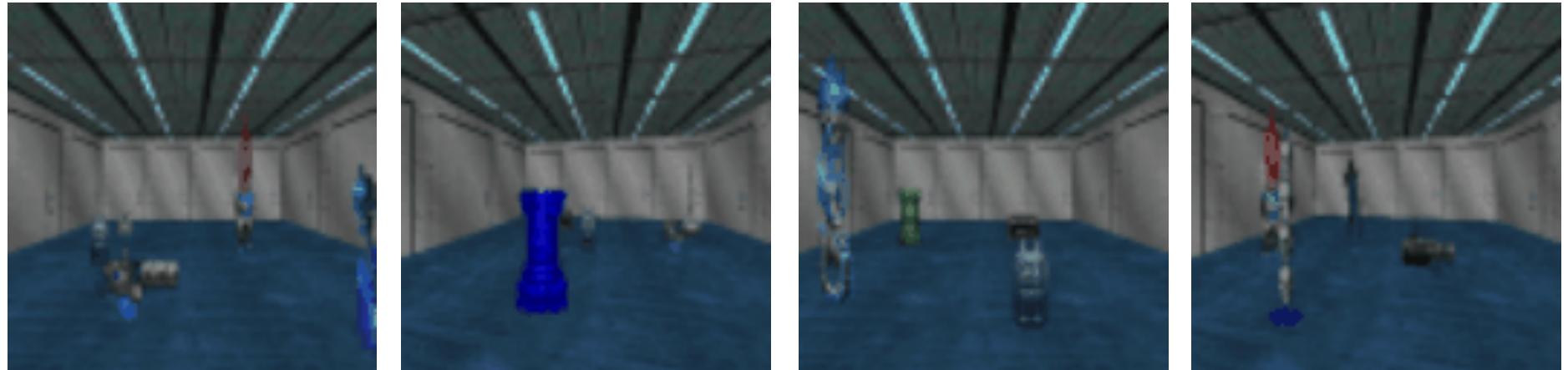
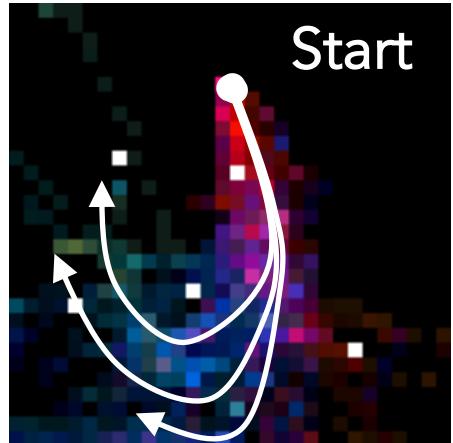
Object Pushing
with Sawyer in MuJoCo

What kind of tasks are discovered?



Direction encoded as color

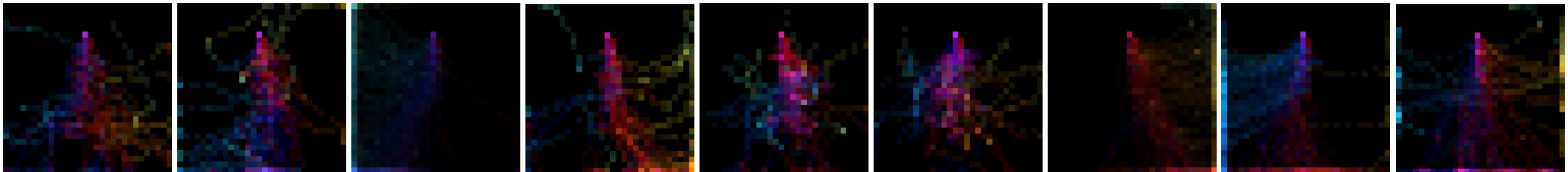
What kind of tasks are discovered?



Direction encoded as color

What kind of tasks are discovered?

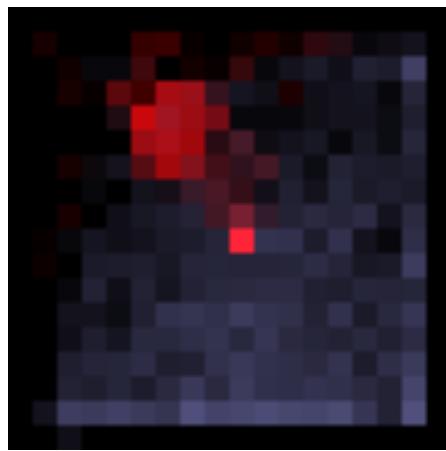
Step 1



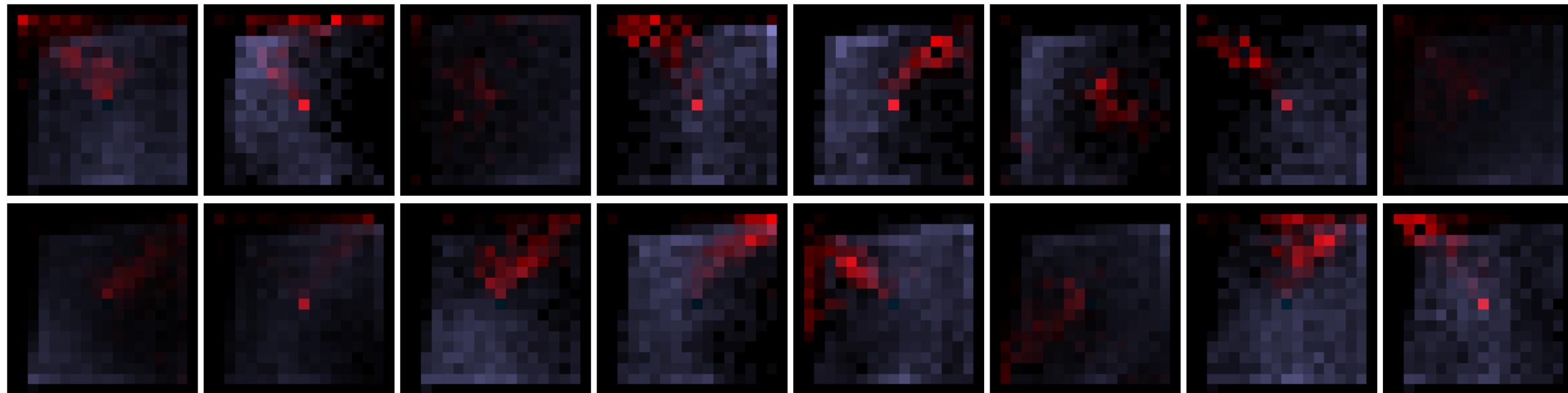
Step 5



What kind of tasks are discovered?

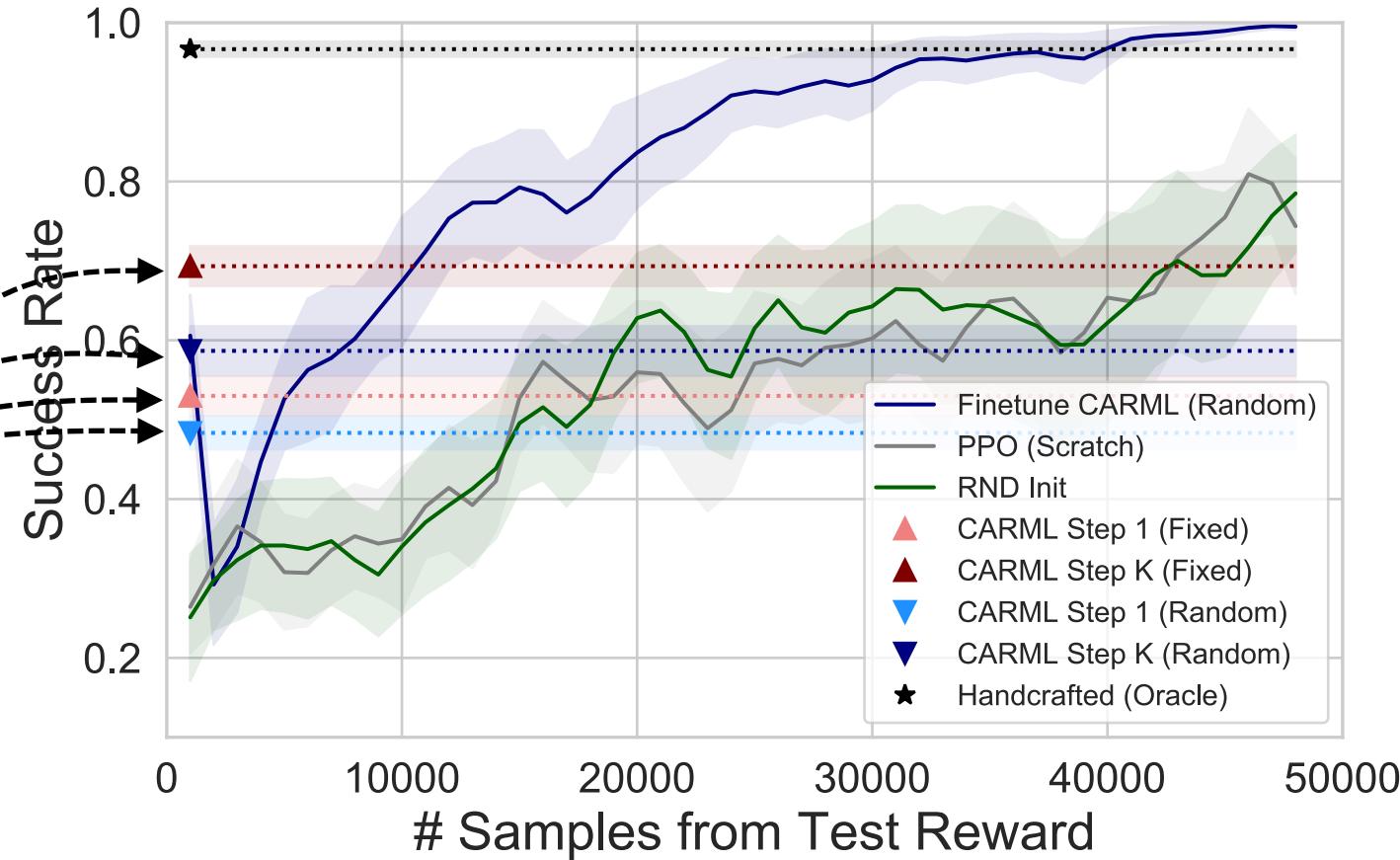


What kind of tasks are discovered?

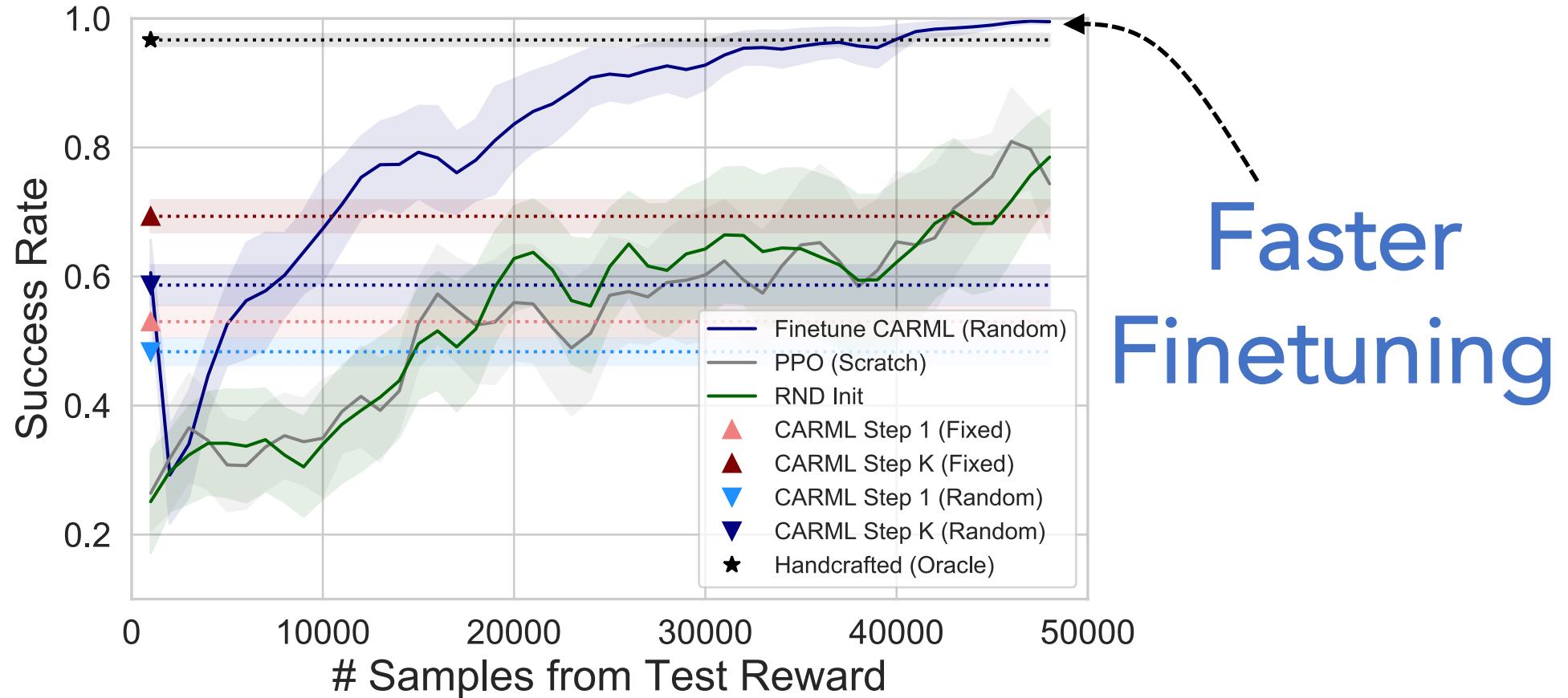


Transfer to Test Tasks – VizDoom

Direct
Transfer



Transfer to Test Tasks – VizDoom



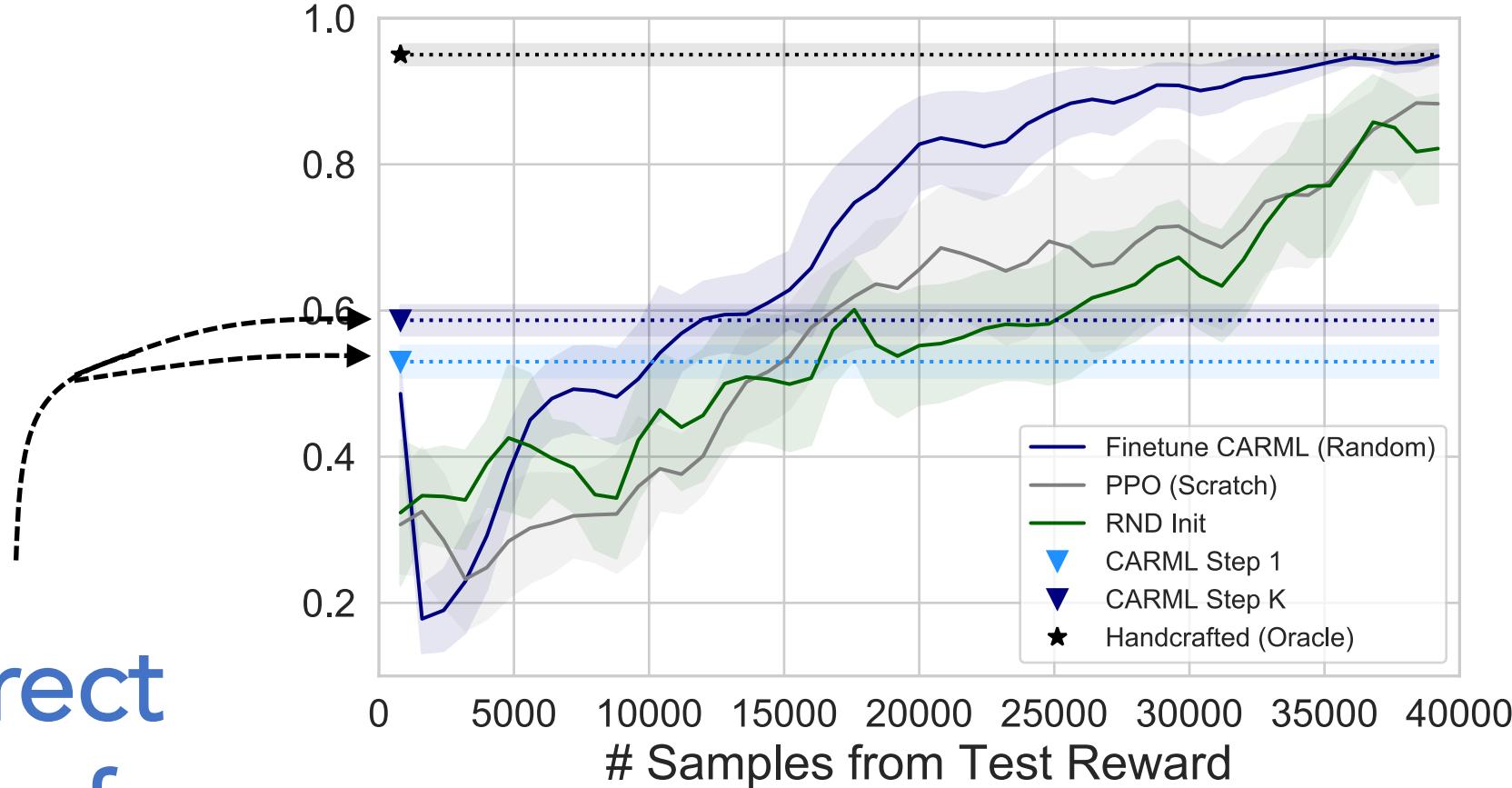
Transfer to Test Tasks – VizDoom

Reward
Episode
boundary

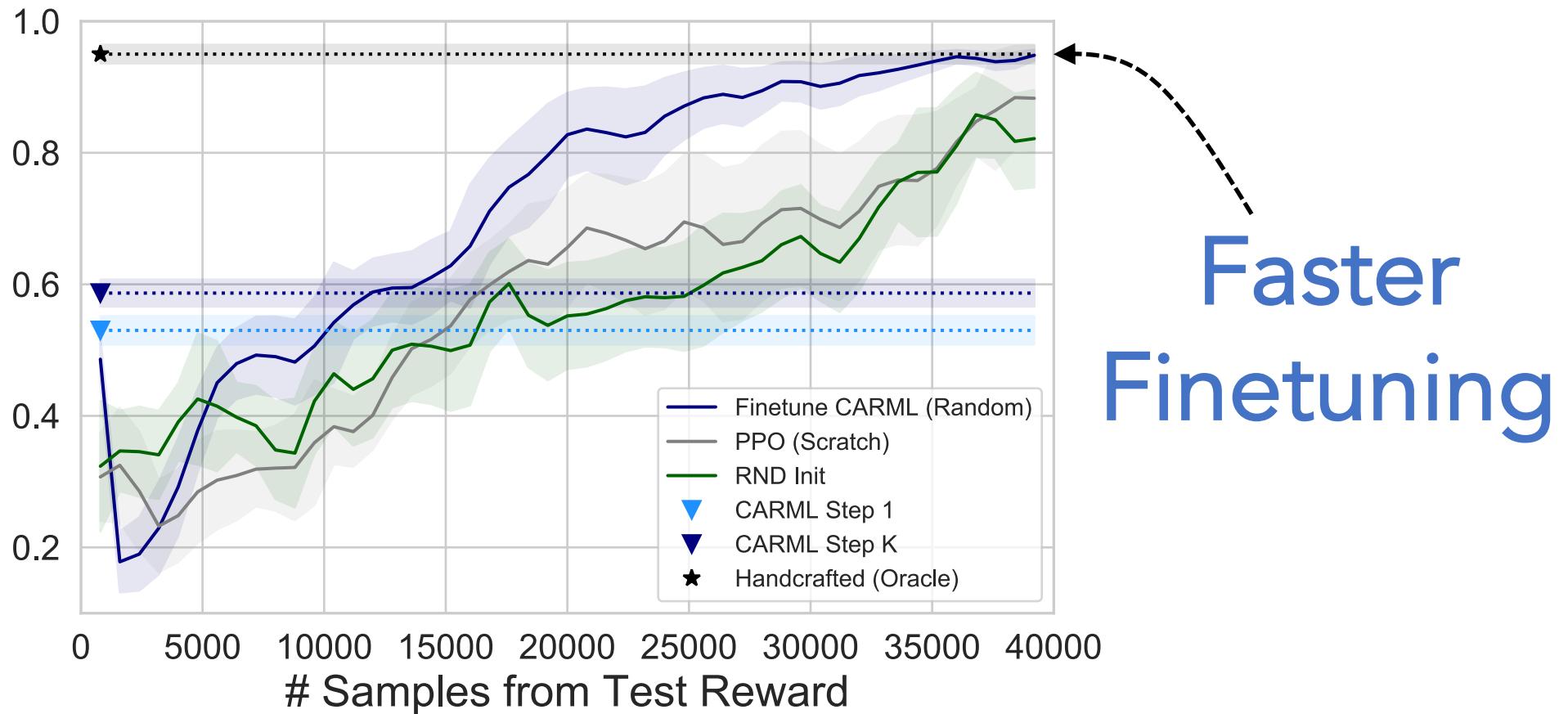


Transfer to Test Tasks – Sawyer

Direct
Transfer



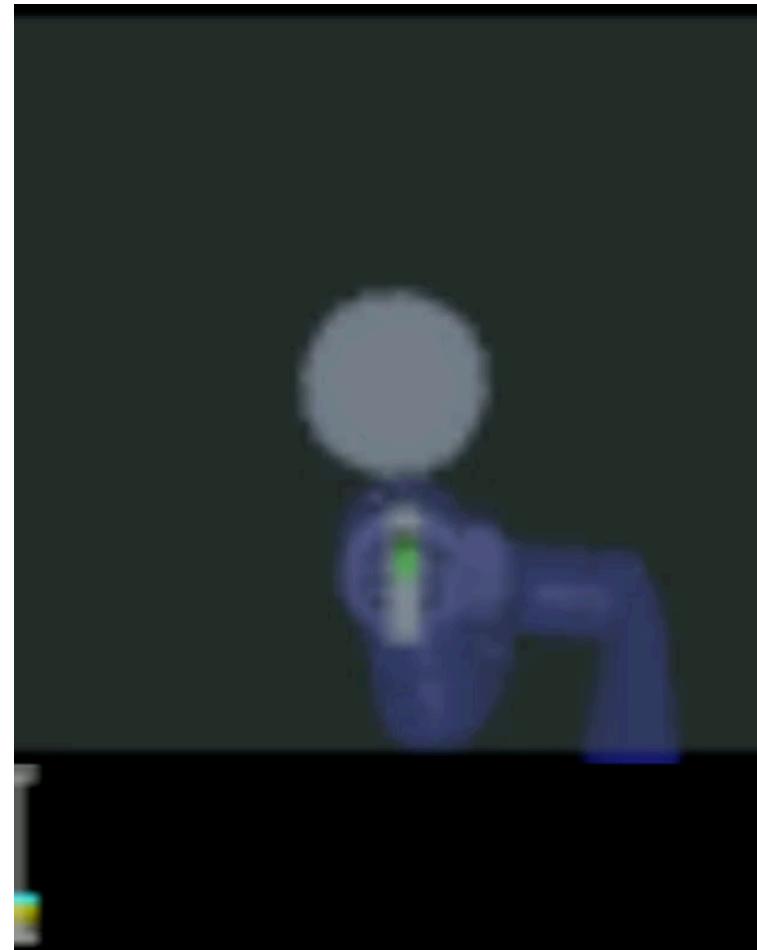
Transfer to Test Tasks – Sawyer



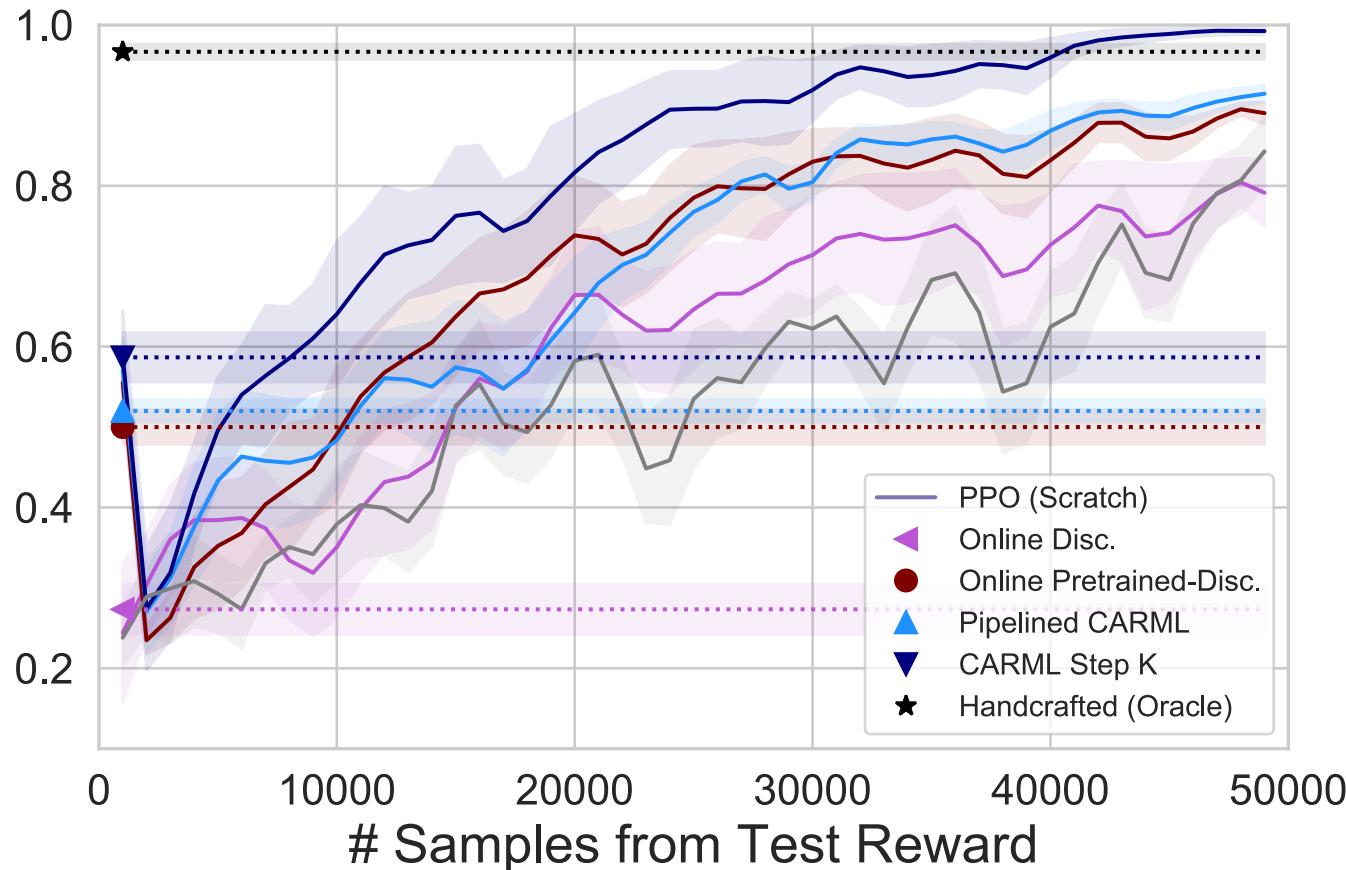
Faster
Finetuning

Transfer to Test Tasks – Sawyer

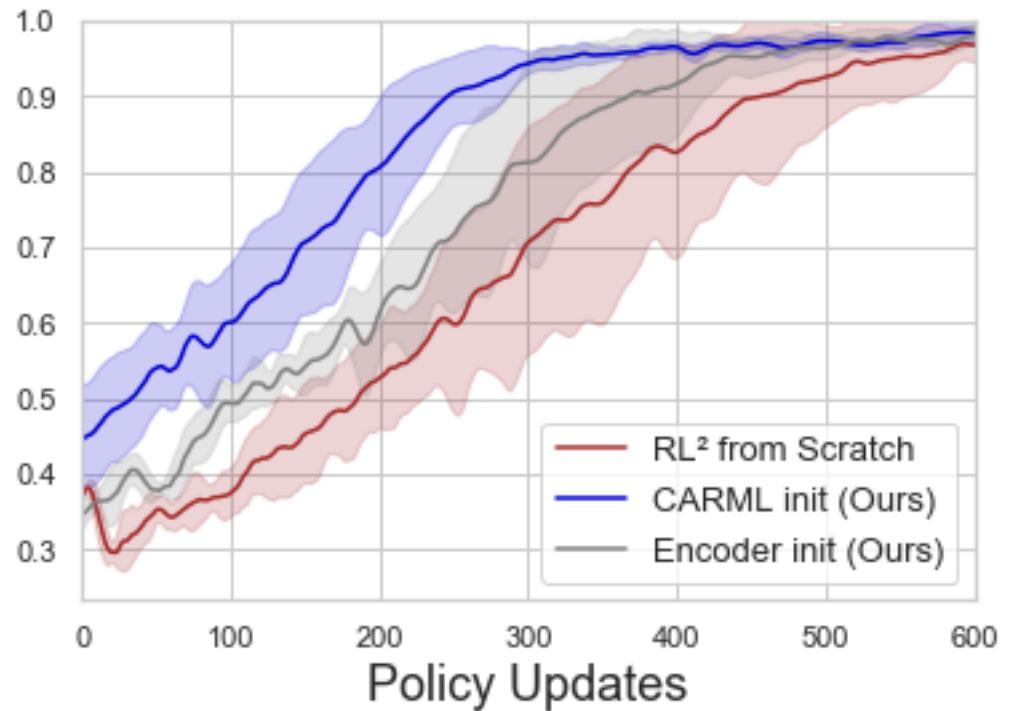
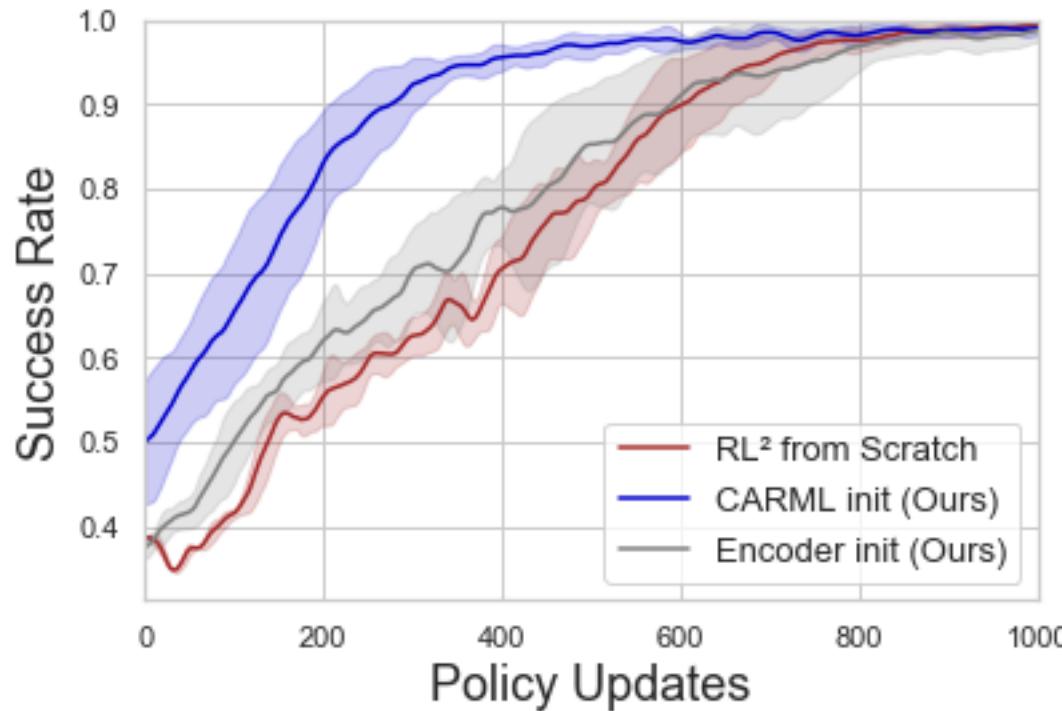
Reward
Episode
boundary



Variants



Faster Supervised Meta-RL



Thank You



Kyle Hsu



Ben Eysenbach



Abhishek Gupta



Sergey Levine



Chelsea Finn

Poster #35, East Exhibition Hall B + C

<https://sites.google.com/view/carm/>