

# Efficient RL via Disentangled Environment and Agent Representations



Kevin Gmelin\*



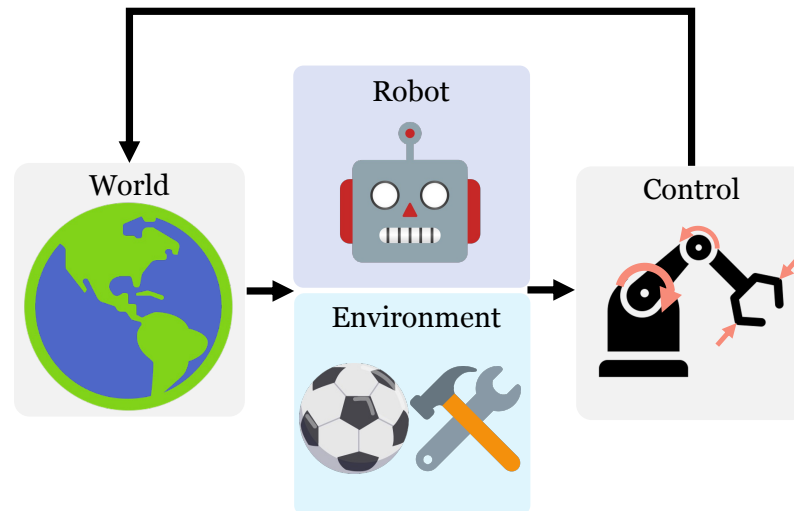
Shikhar Bahl\*



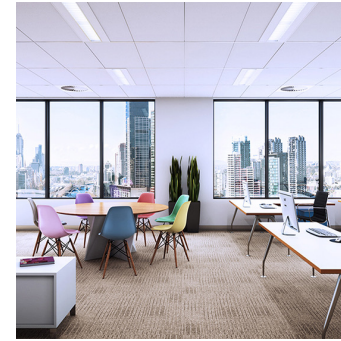
Russell Mendonca



Deepak Pathak



# Goal: General-Purpose Robots

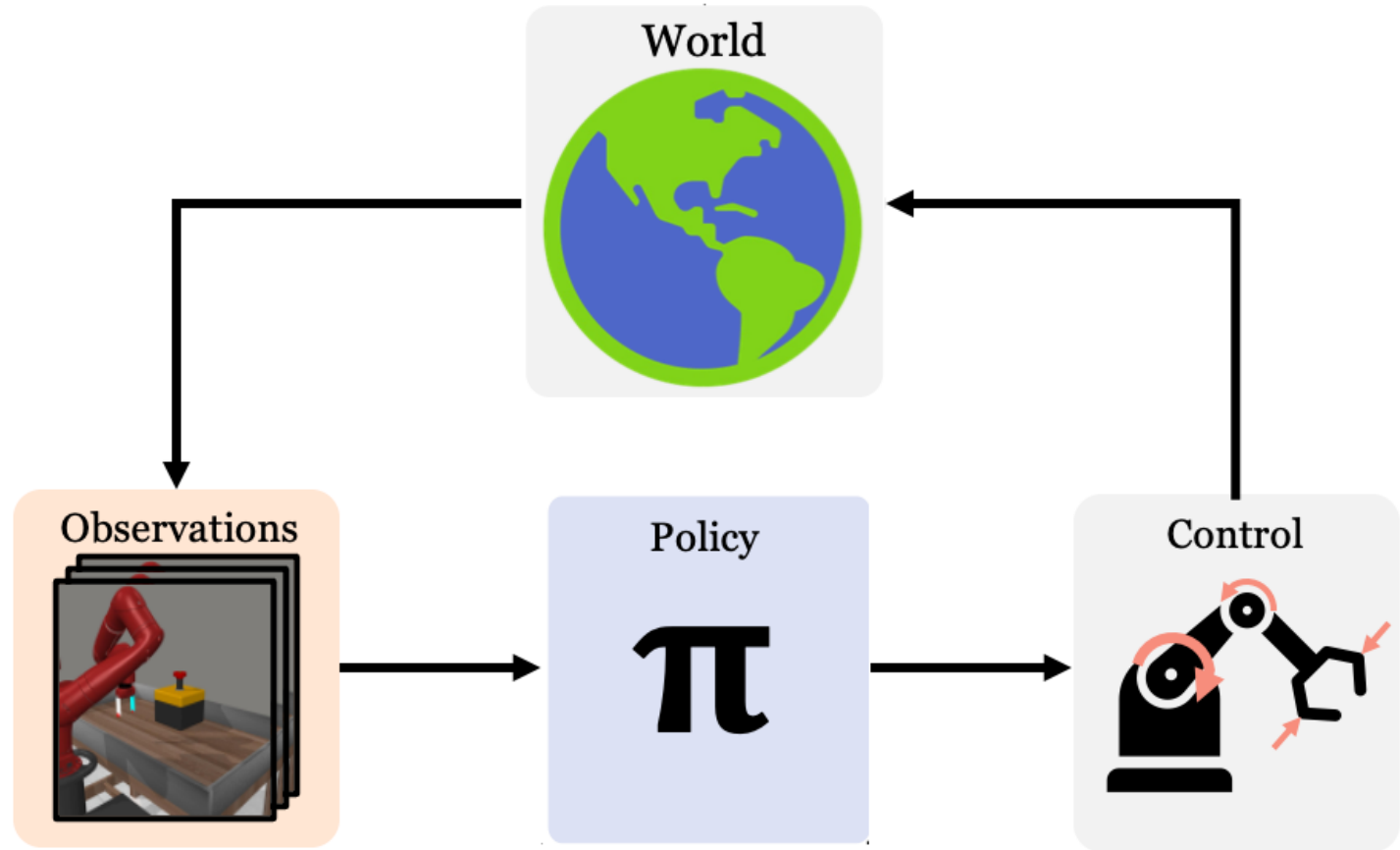


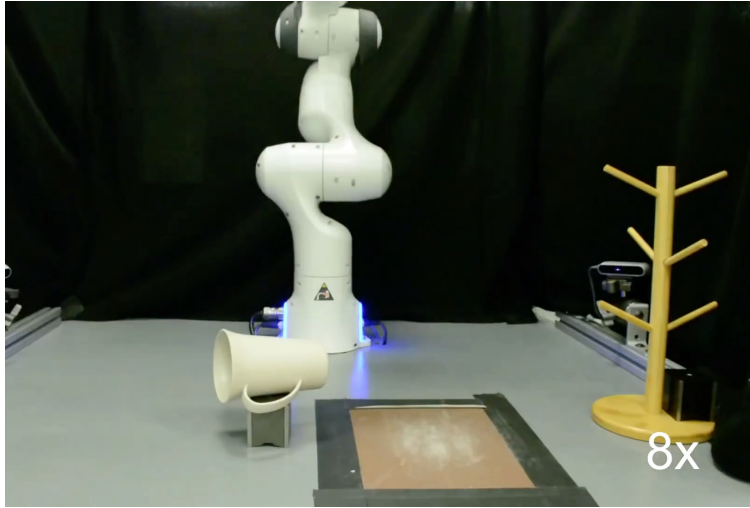
Robots that can do *thousands of tasks in thousands of environments*

# Visual Reinforcement Learning

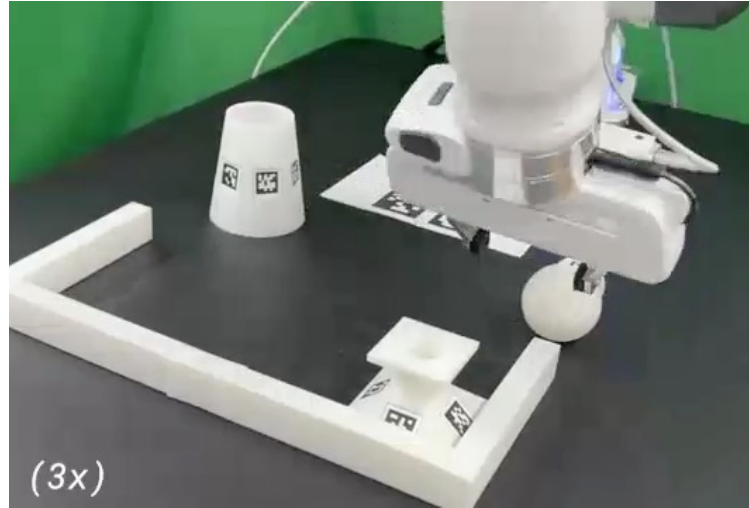
**Goal:** Learn mapping from images to actions

Want **sample-efficient** algorithms

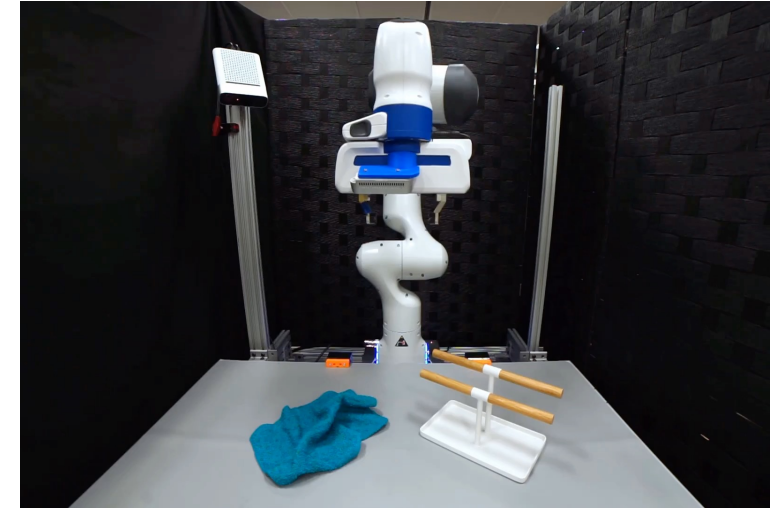




Simeonov et al., 2022



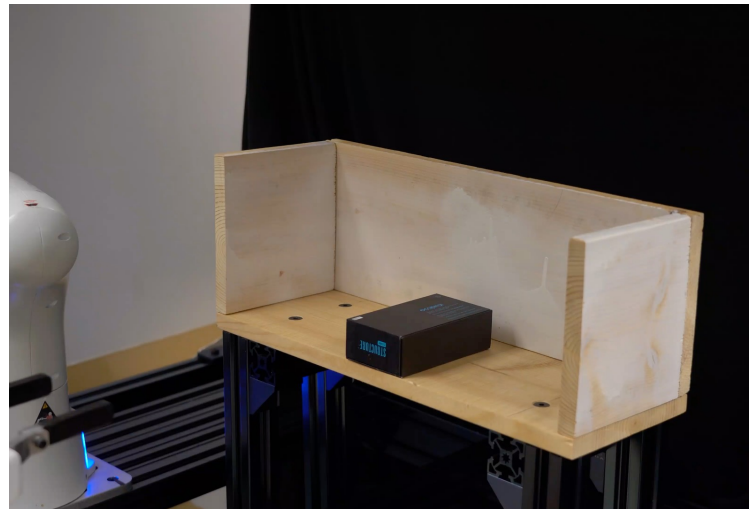
Heo et al., 2023



Huang et al., 2023



Mendonca et al., 2023



Liang et al., 2022



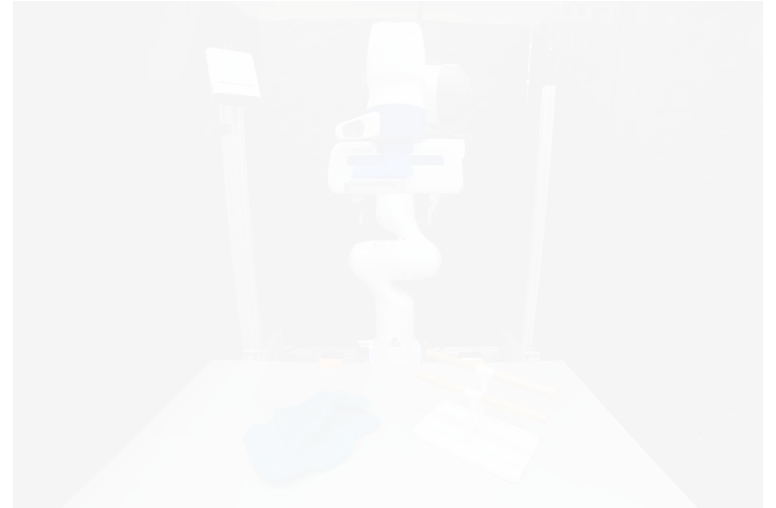
Zhou et al., 2023



Simeonov et al., 2022



Heo et al., 2023



Huang et al., 2023

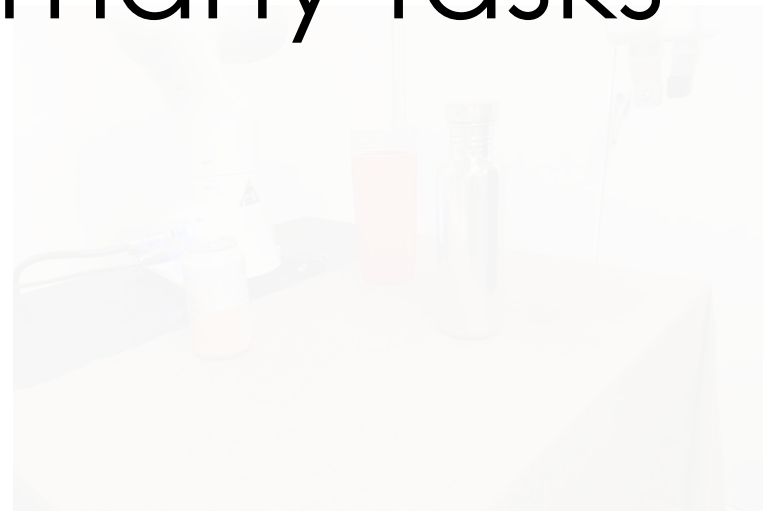
**Agent** is common across many tasks



Mendonca et al., 2023

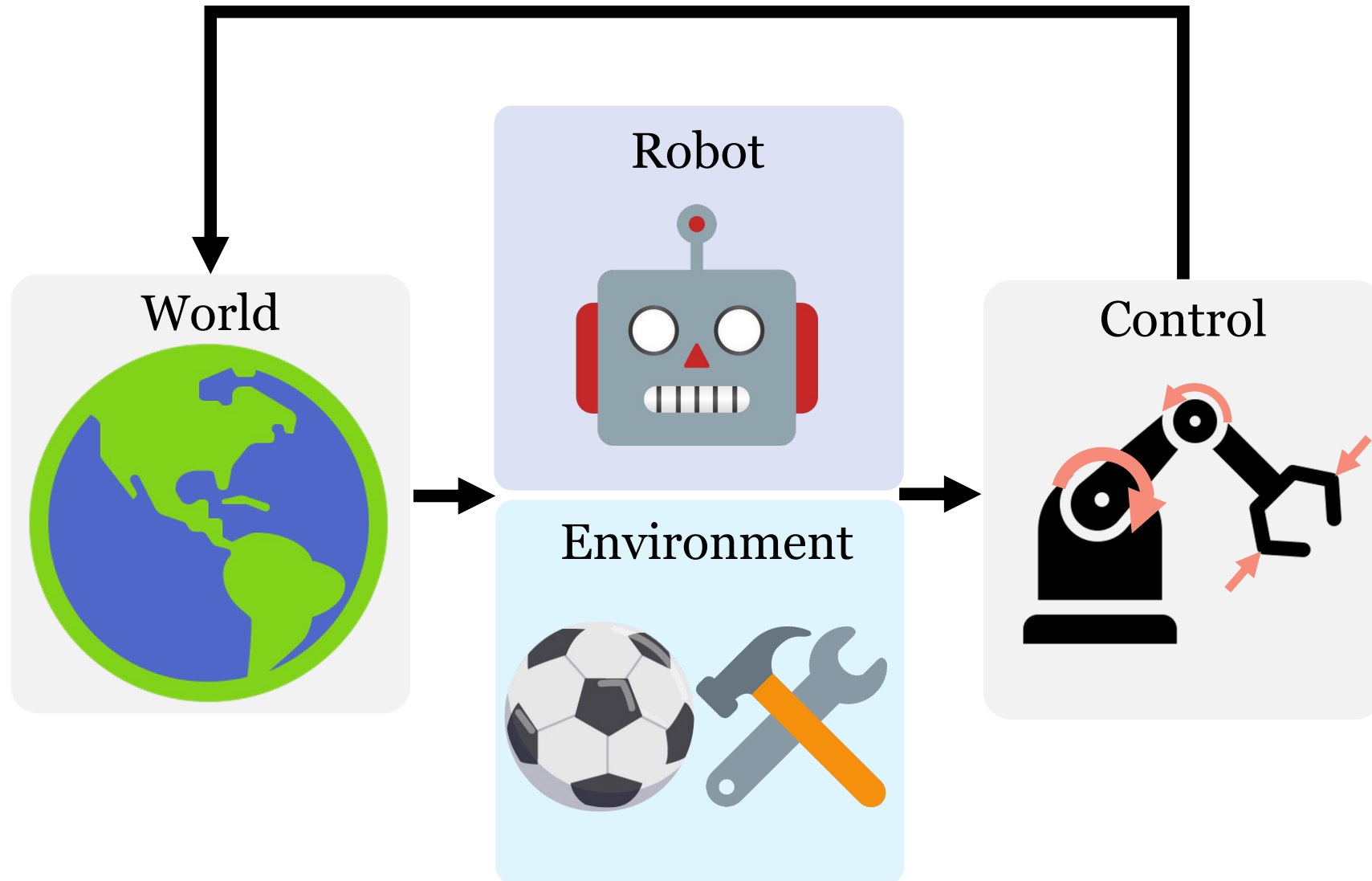


Liang et al., 2022



Zhou et al., 2023

# Agent vs Environment Rep. for Control



# Agent vs Environment Rep. for Control

How can we learn disentangled **agent** & **environment** representations?



# How to obtain supervision?

## Environment?

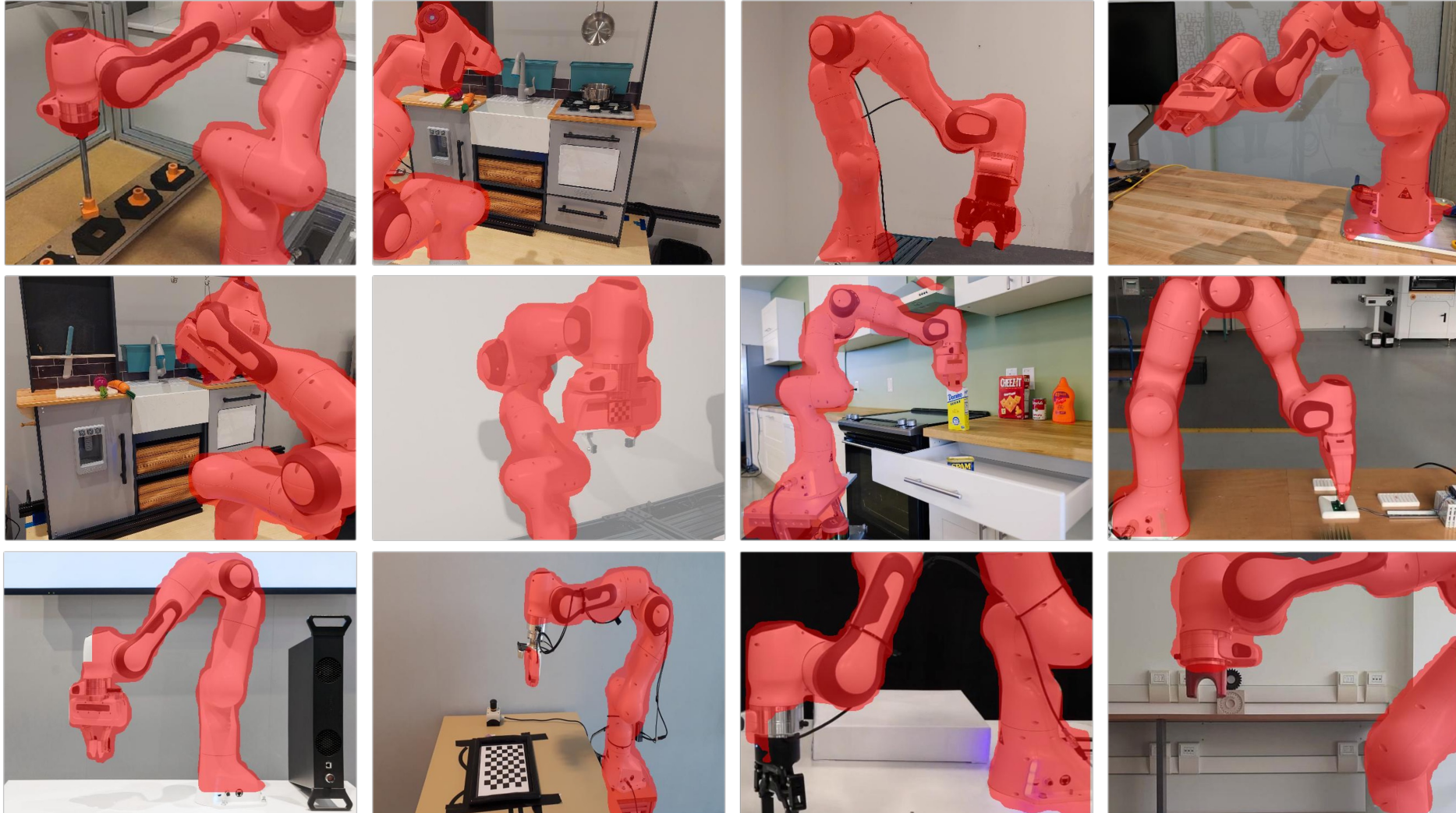
- Can directly use scene observation

## What about for the agent?

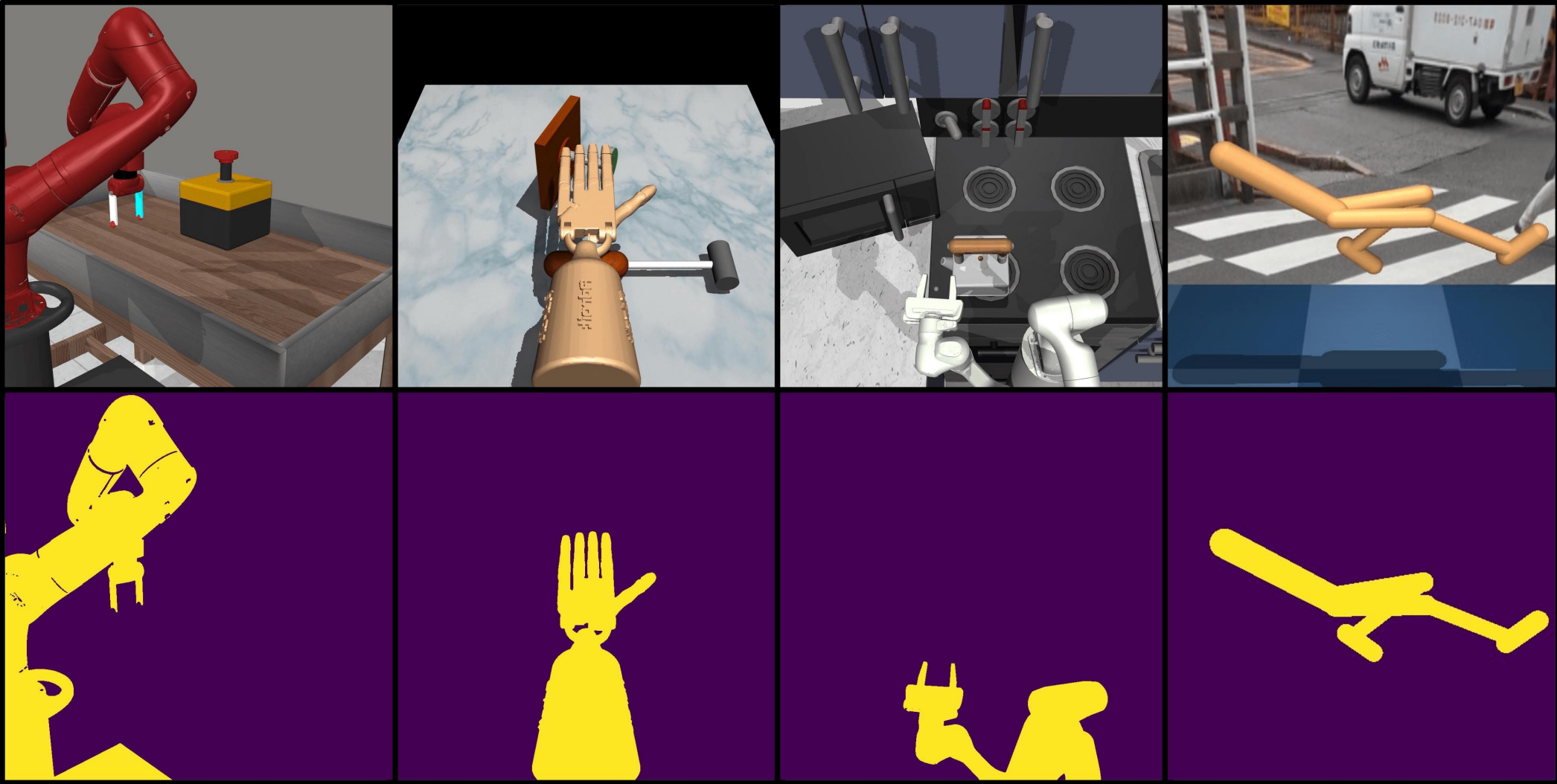
- Common to know what robot looks like
- **Masks** are a good proxy for agent info
- Use **off-the-shelf models**



# Real World Masks



# Simulation



# Shelf-Supervised Agent Masks



Segment Anything  
Kirillov et al., 2023



Mask2Former  
Cheng et al., 2022

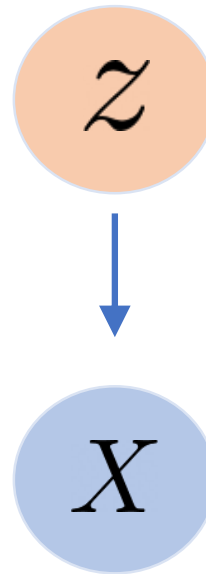
How do we incorporate this into RL?

# Agent/Environment Latent Structure

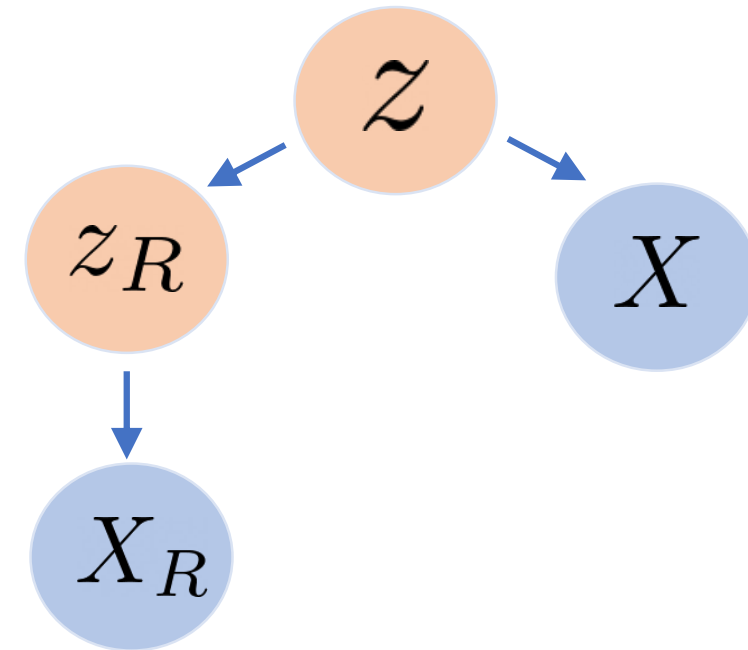
$$p(X, X_R, z, z_R) = p(z) \cdot p(z_R|z) \cdot p(X|z) \cdot p(X_R|z_R)$$

$$\mathcal{L} = \mathbb{E}_{z, z_R \sim q} [\log p(X|z)] + \mathbb{E}_{z, z_R \sim q} [\log p(X_R|z_R)] \\ - D_{KL}(q(z, z_R|X) || (p(z, z_R)))$$

Prior

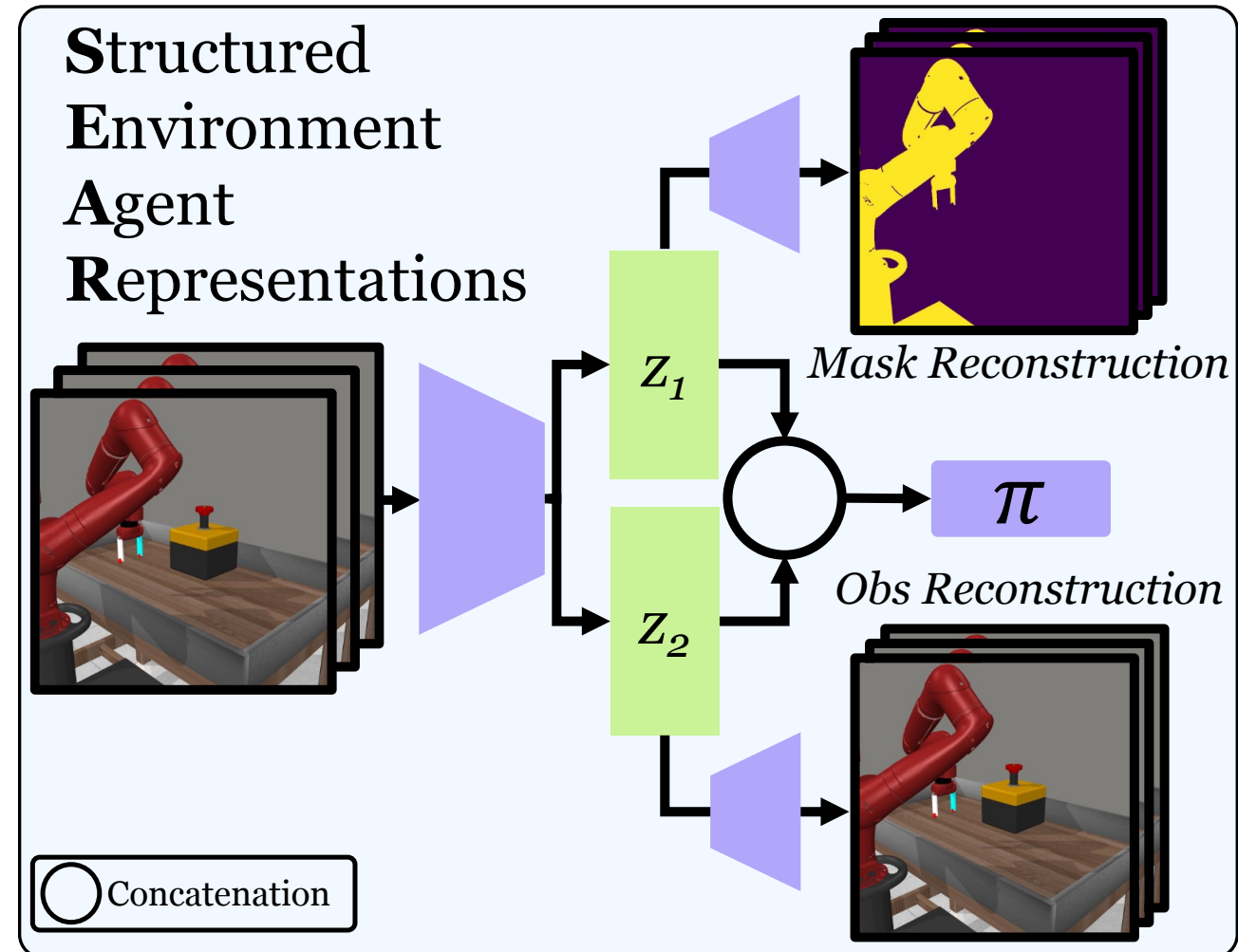


SEAR



# SEAR: Structured Environment-Agent Representations

Structure latent into  
**agent** ( $z_1$ ), **env** ( $z_2$ )

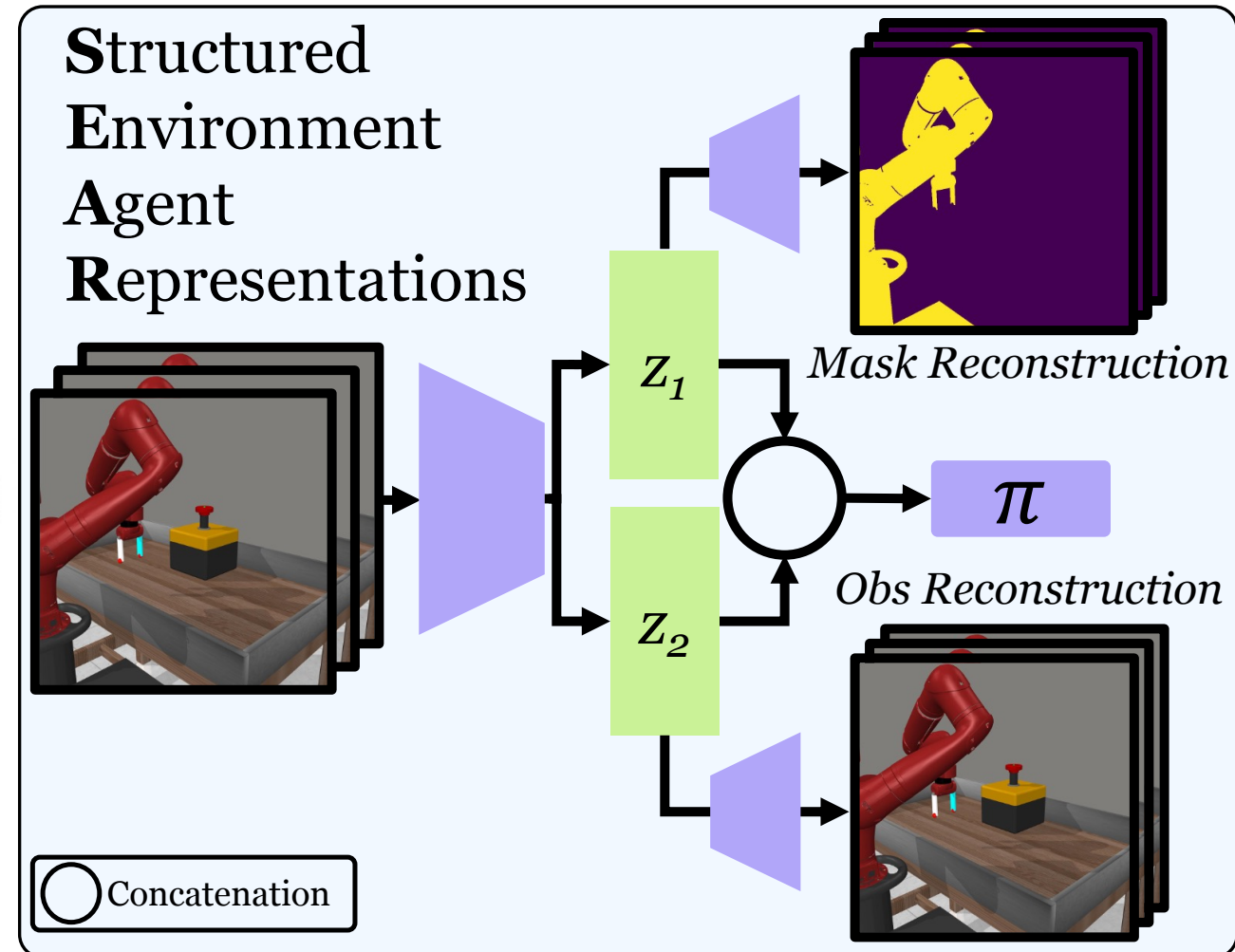


# SEAR: Structured Environment-Agent Representations

Leverage agent **mask** as supervision

$$\mathcal{L}_{mask} = M \log P_{\phi}(M|z_R) + (1-M) \log (1 - P_{\phi}(M|z_R))$$

$$\mathcal{L} = \mathcal{L}_{critic} + c_1 \mathcal{L}_{recon} + c_2 \mathcal{L}_{mask}$$



# Existing Approaches

## Data-Augmentations

- RAD
- DrQ-v2

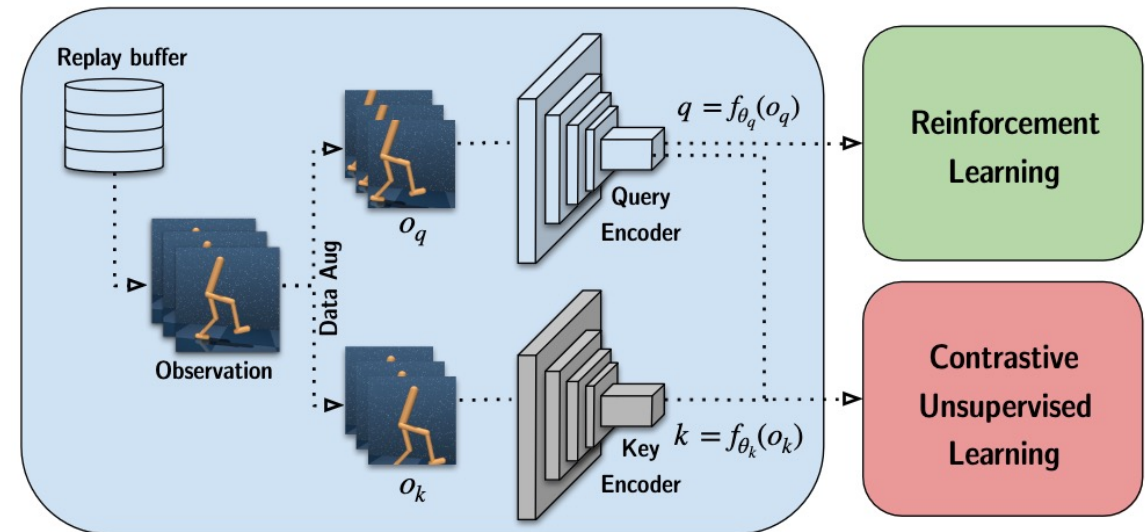
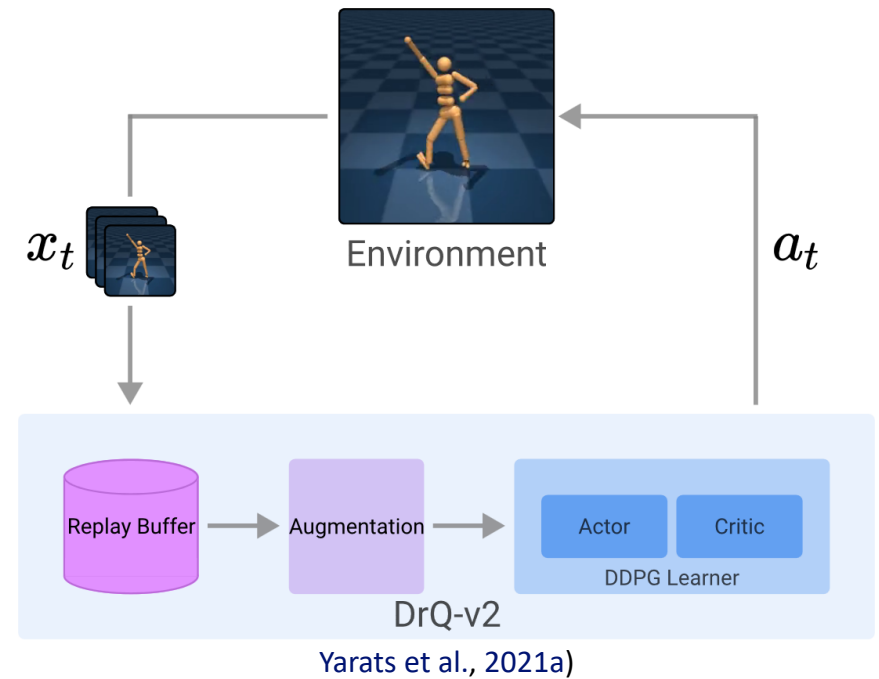
## Auxillary Losses

- CURL
- SAC-AE

## Model-Based

- Dreamer

And many more...





# Existing Approaches

## Data-Augmentations

- RAD
- DrQ-v2

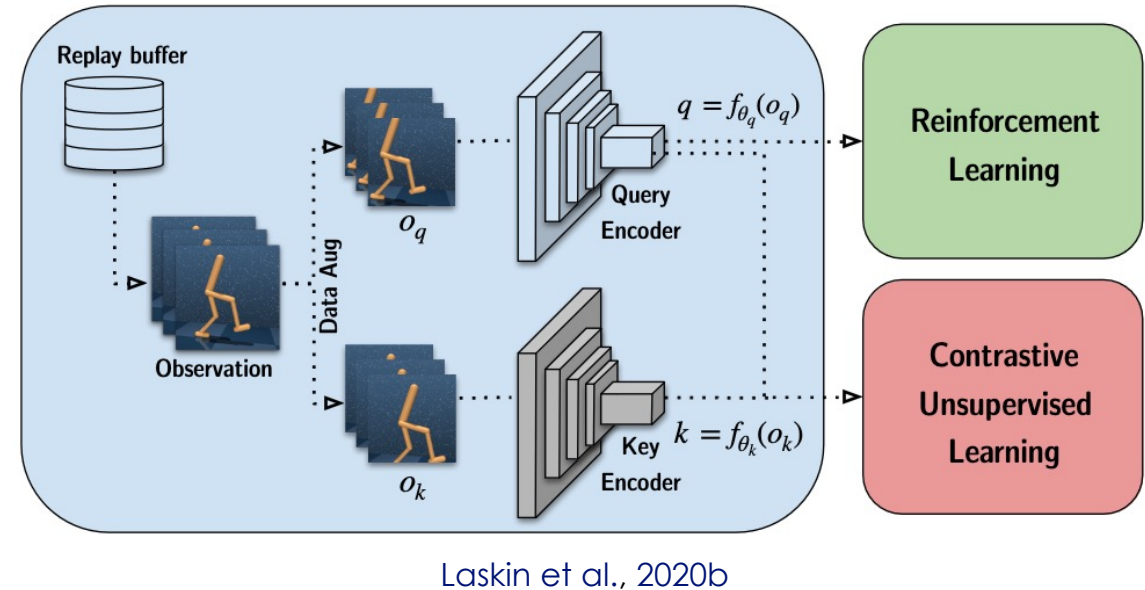
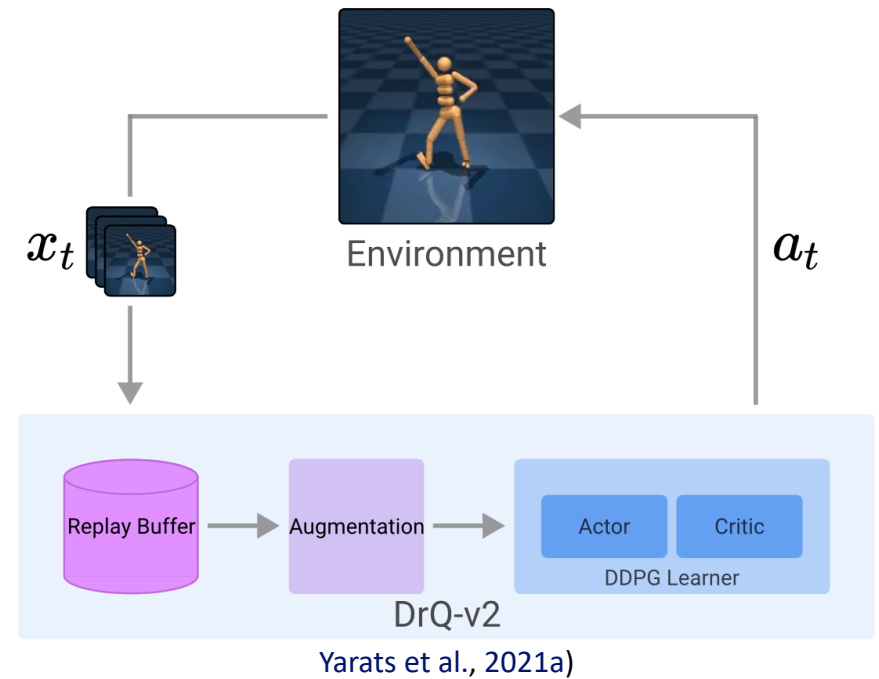
## Auxillary Losses

- CURL
- SAC-AE

## Model-Based

- Dreamer

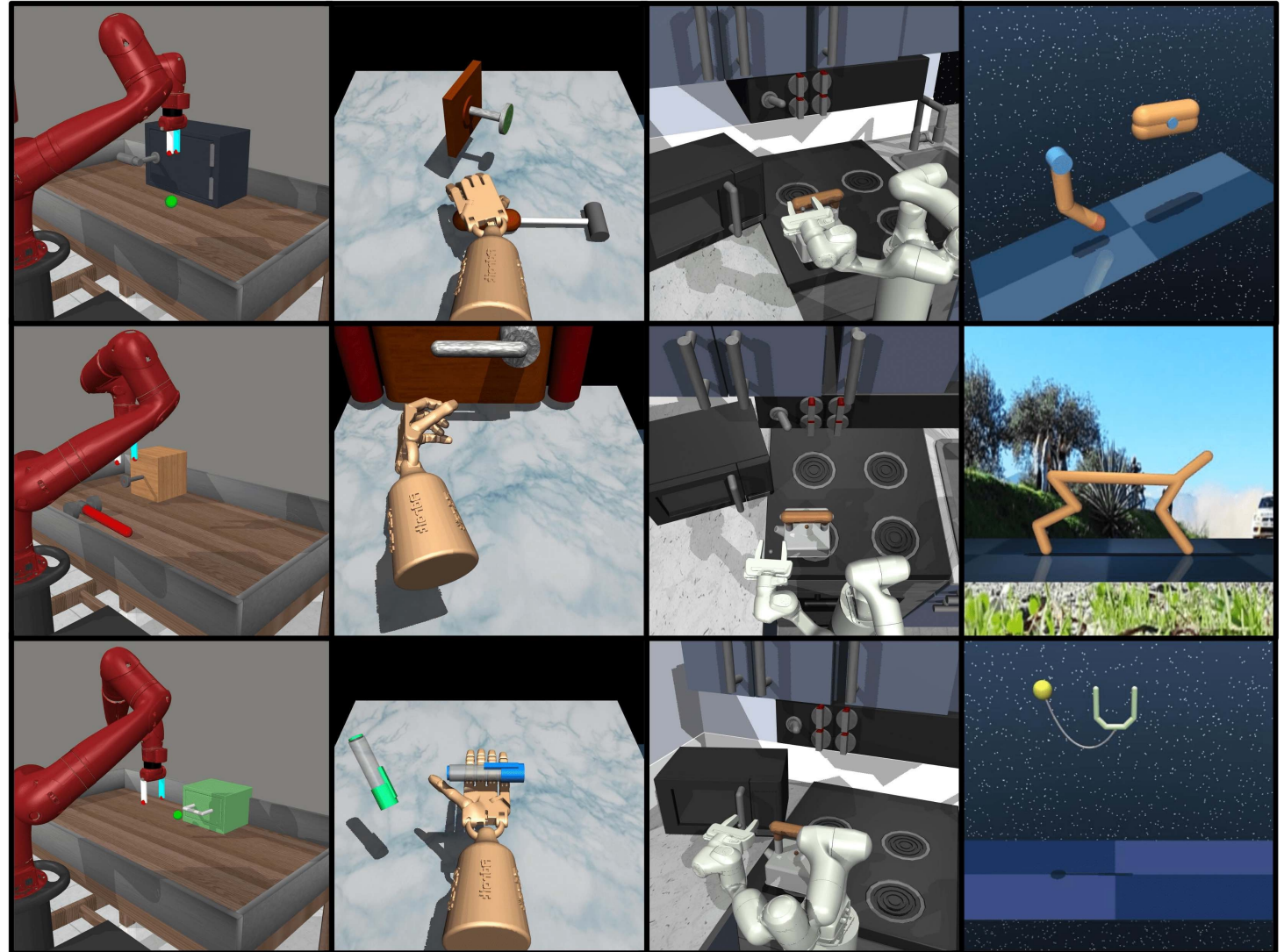
And many more...



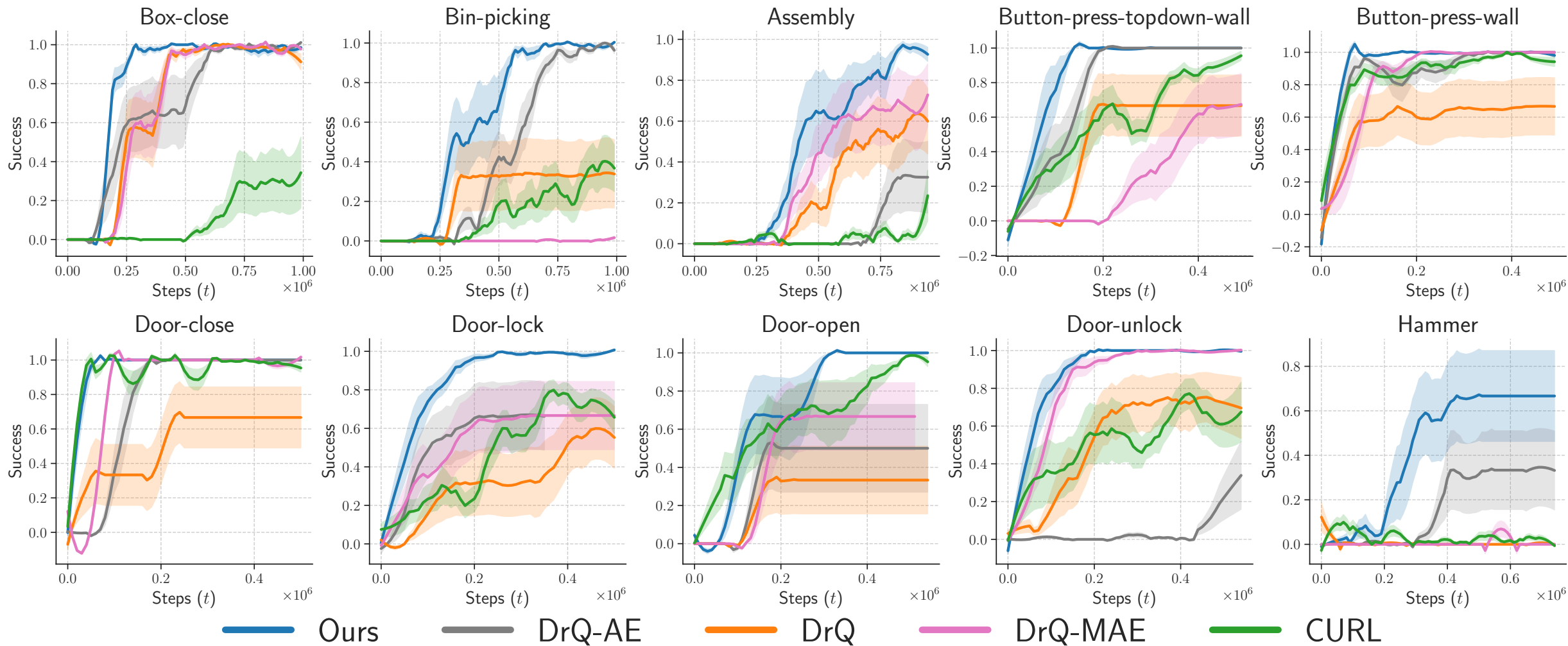
# Environments

---

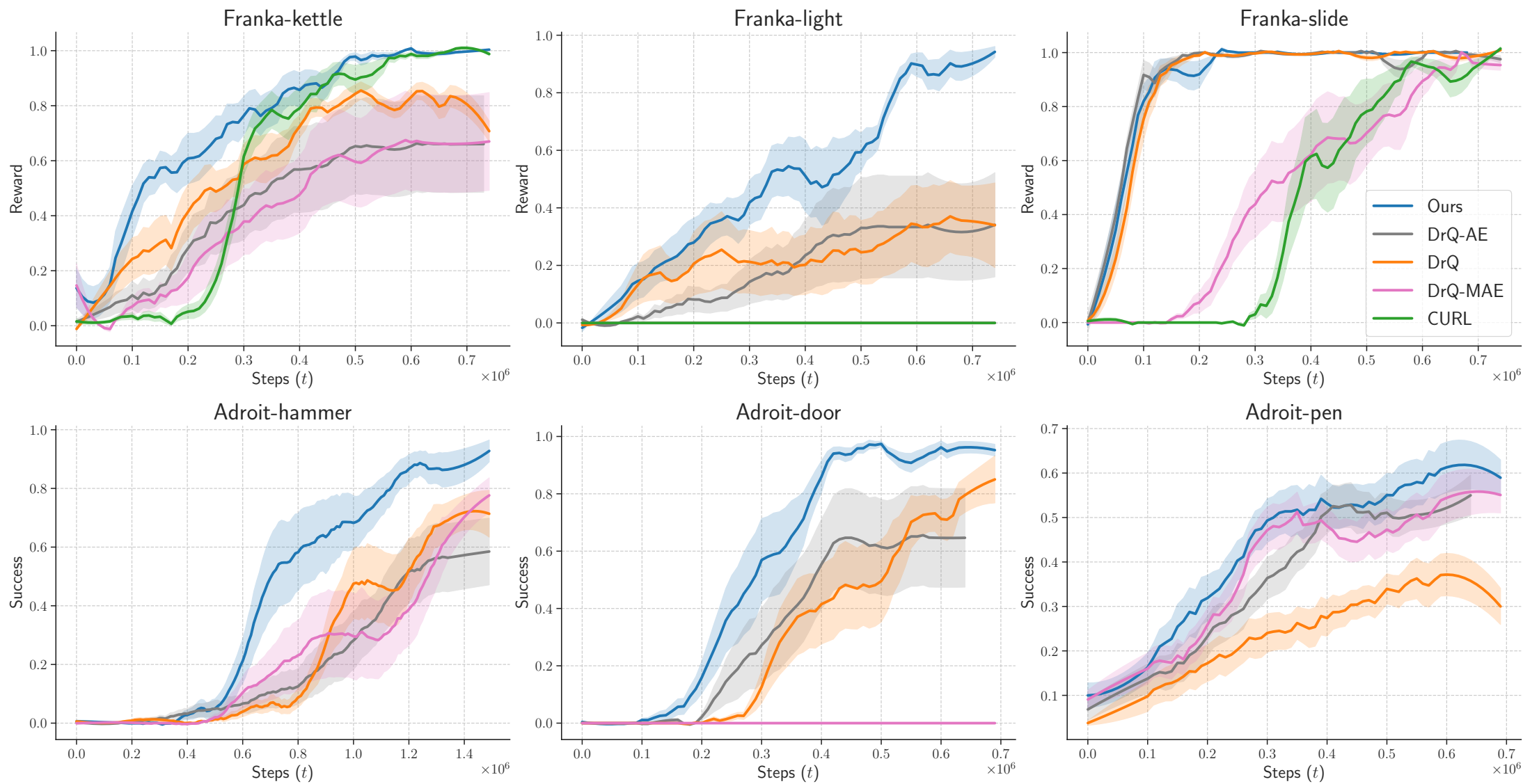
- Meta-World (Yu et al., 2020)
- Hand Manipulation Suite (Rajeswaran et al., 2017)
- Franka Kitchen (Gupta et al., 2019)
- Distracting Control Suite (Stone et al., 2021)



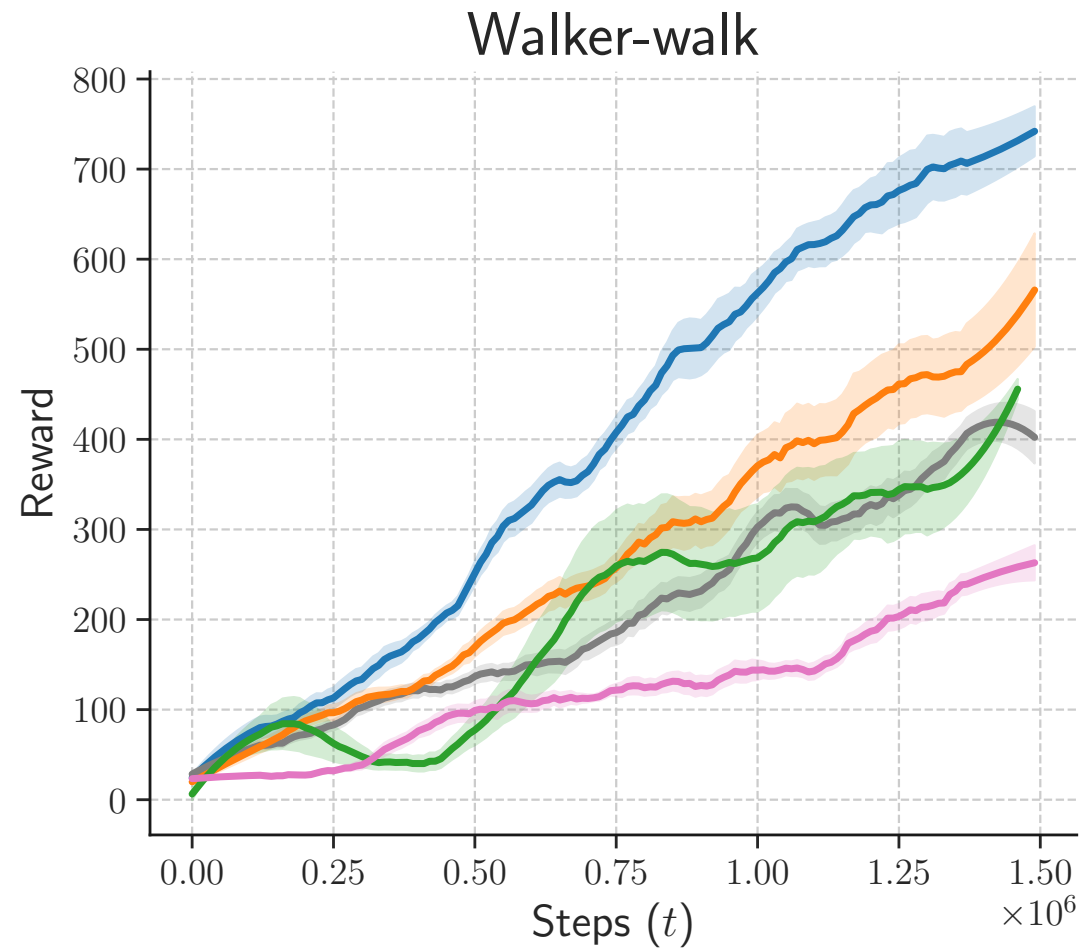
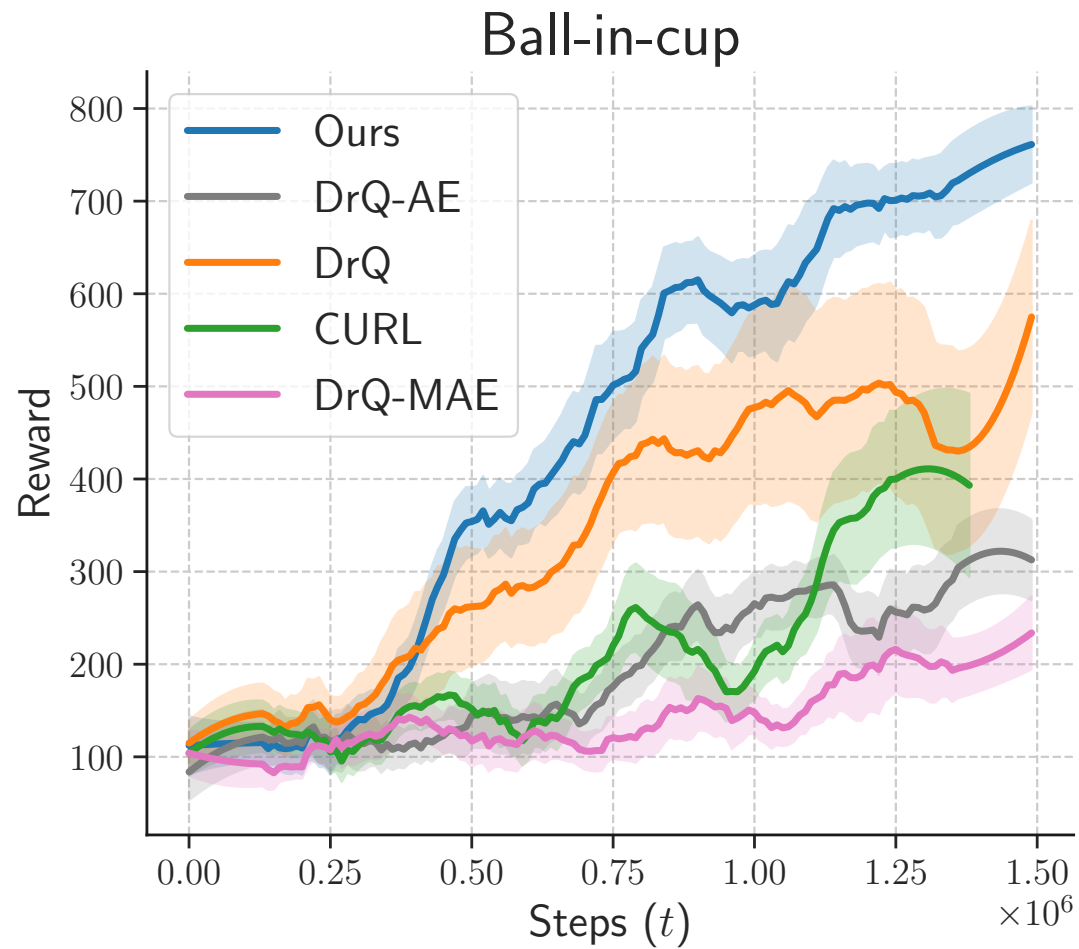
# Meta-World



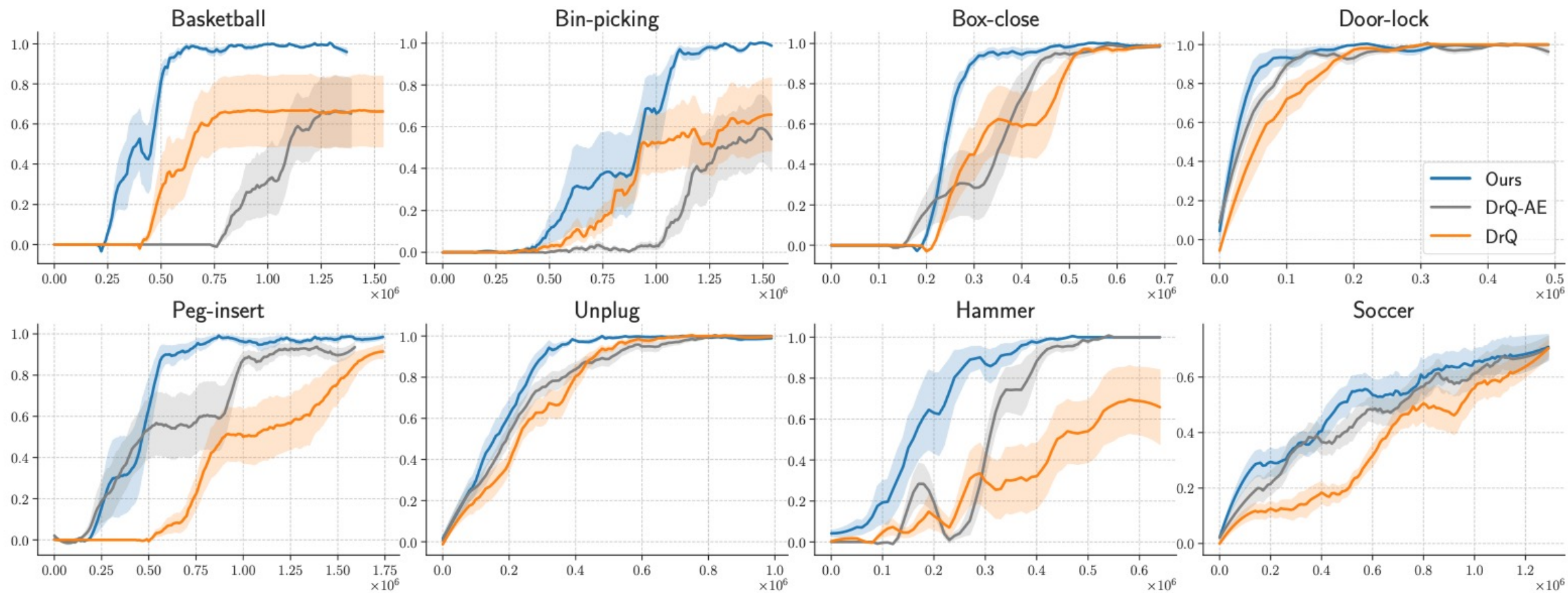
# Franka Kitchen & Adroit



# Distract Gym

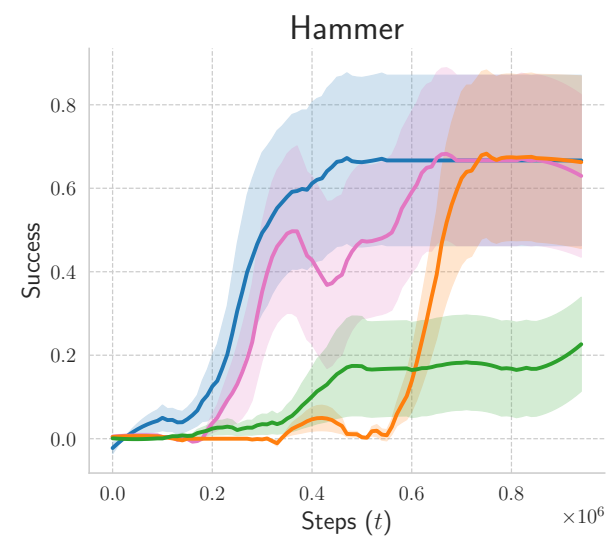
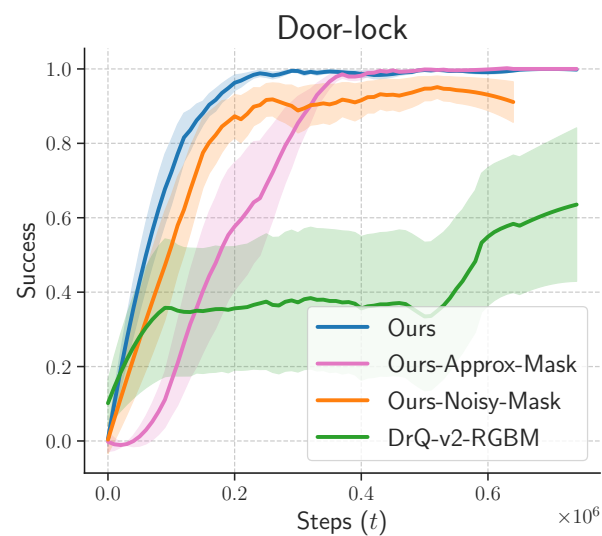


# Transfer to new environments

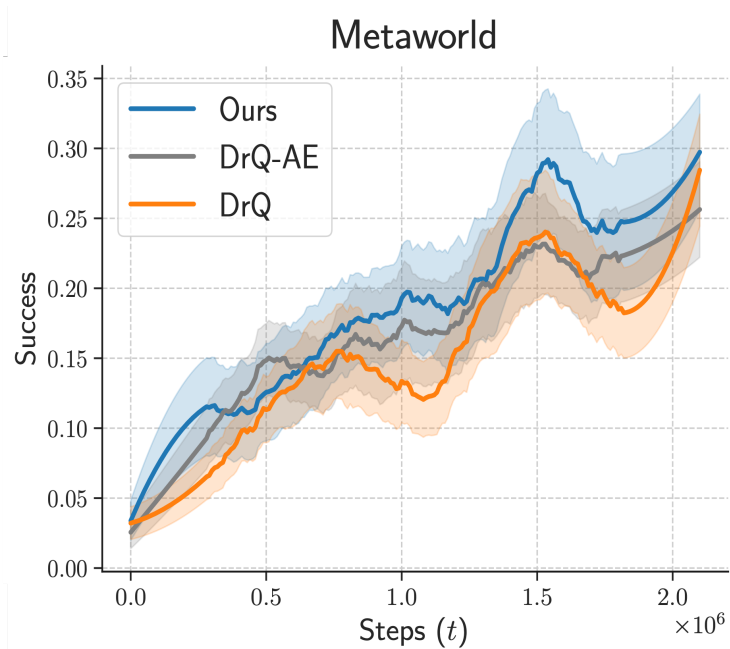


# Robustness to Mask Quality

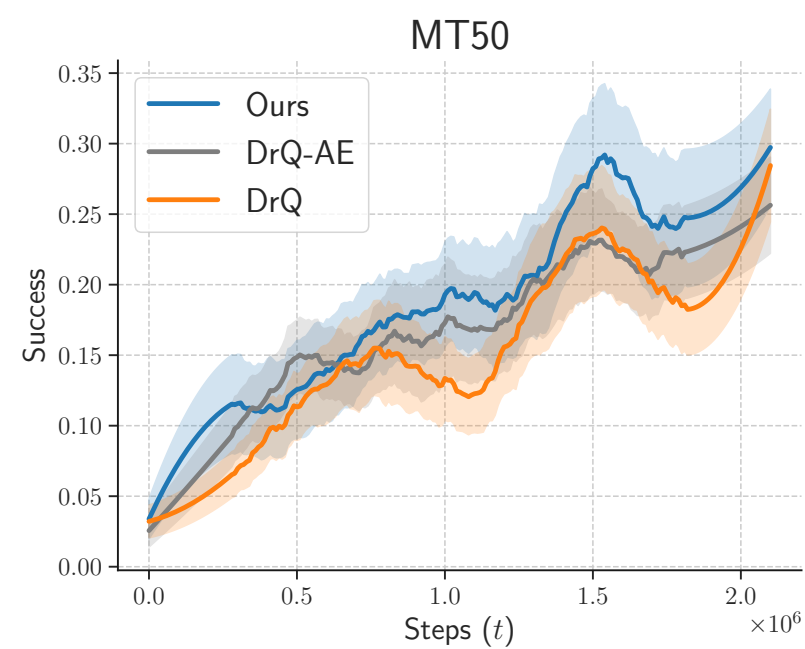
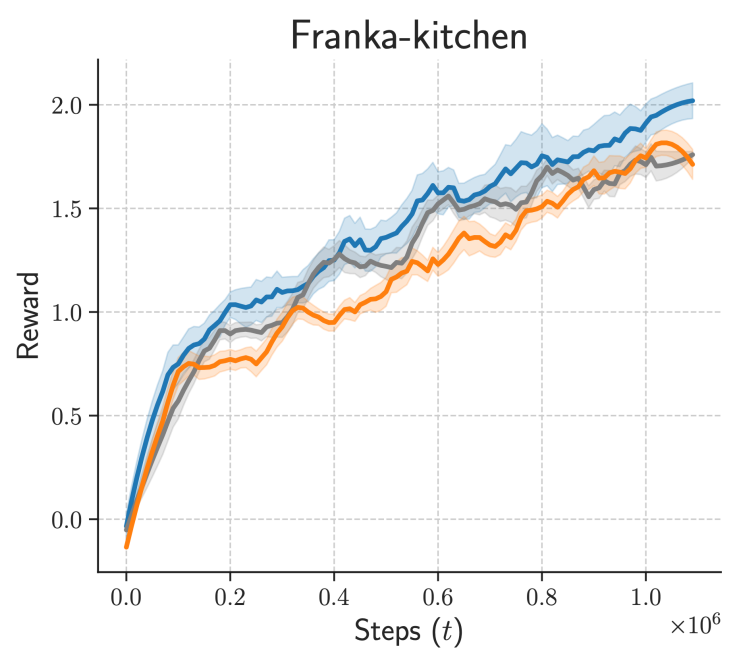
We train with noisy mask (left) and lower resolution, approximate mask (right)



# What about multi-task settings?



Small-Scale Setups



Larger-Scale Setups



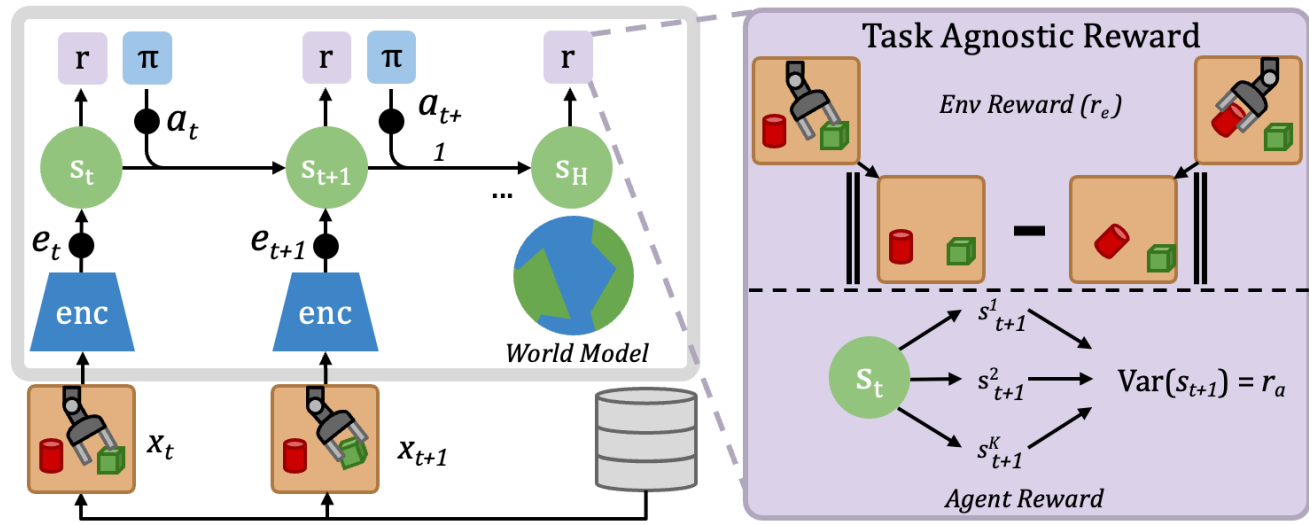
Applications in real robots?

# Environment-centric Exploration



***Env. Chang Exploration Objective***

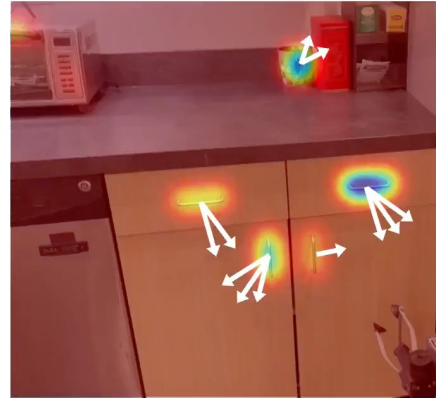
$$c_k = \max_{i,j} \|\Phi_f(R_{k,i}) - \Phi_f(R_{k,j})\|$$



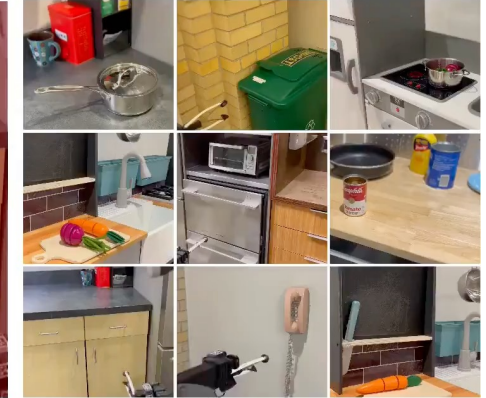
Human Videos



Learned Affordance Model



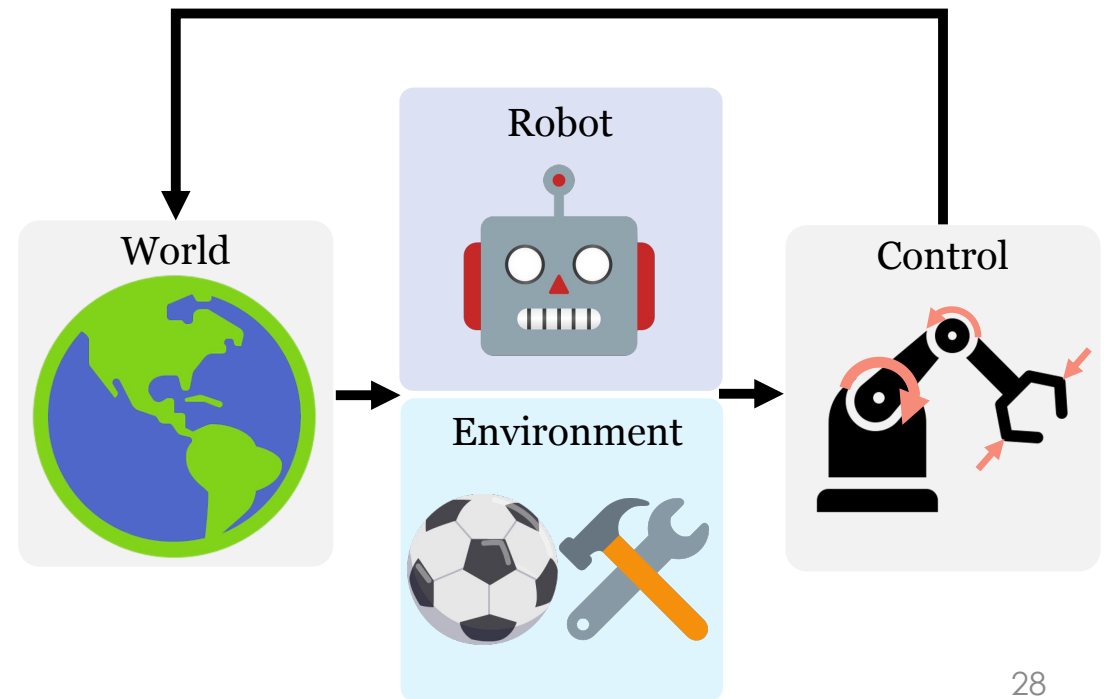
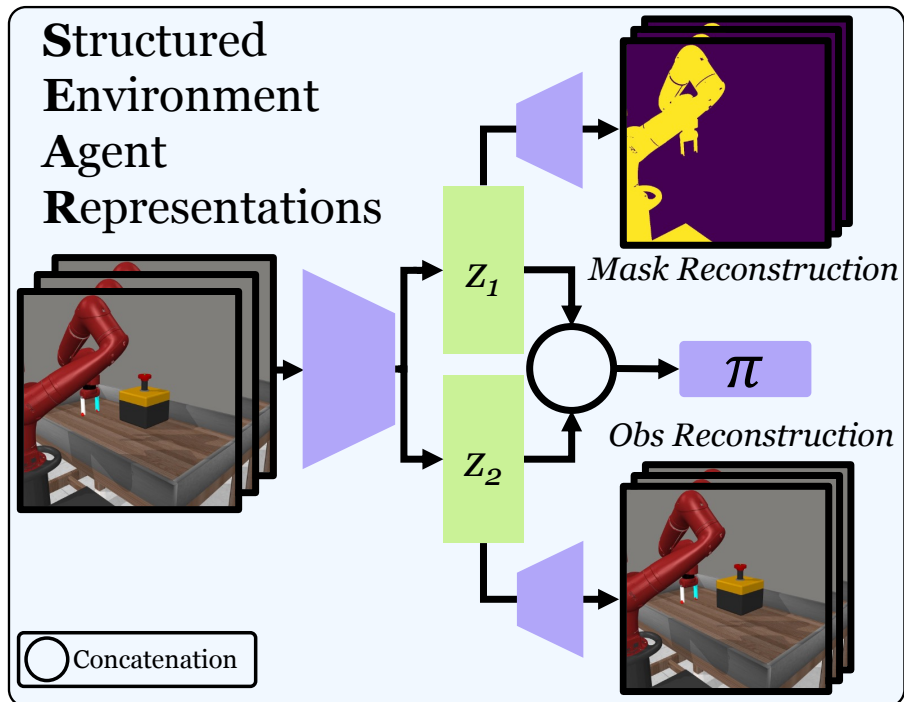
Robot Execution



Agent-Agnostic, change-based exploration can be useful

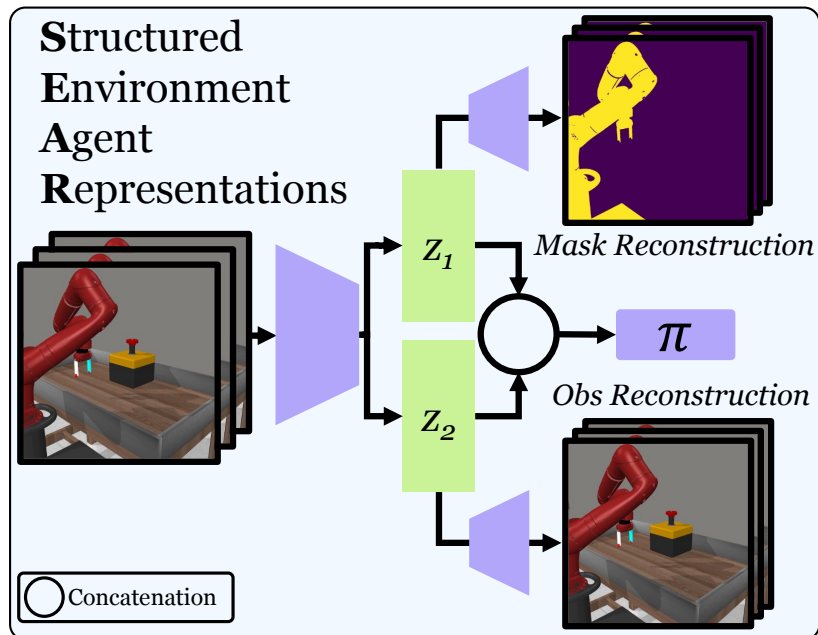
# Limitations

- Further investigation into **multi-task** setting needed
- Robot must be **visible** in image
- Only examined training from **scratch**
- Only added SEAR onto **DrQ-v2**



# Final Thoughts

- Decoupled representation **boosts performance**
- SEAR can help with **transfer**
- Masks are **readily available** from shelf-supervised models
- Can be added to **any visual RL** approach



[sear-rl.github.io](https://sear-rl.github.io)

