



Yanjie Ze  
[yanjieze.com](http://yanjieze.com)



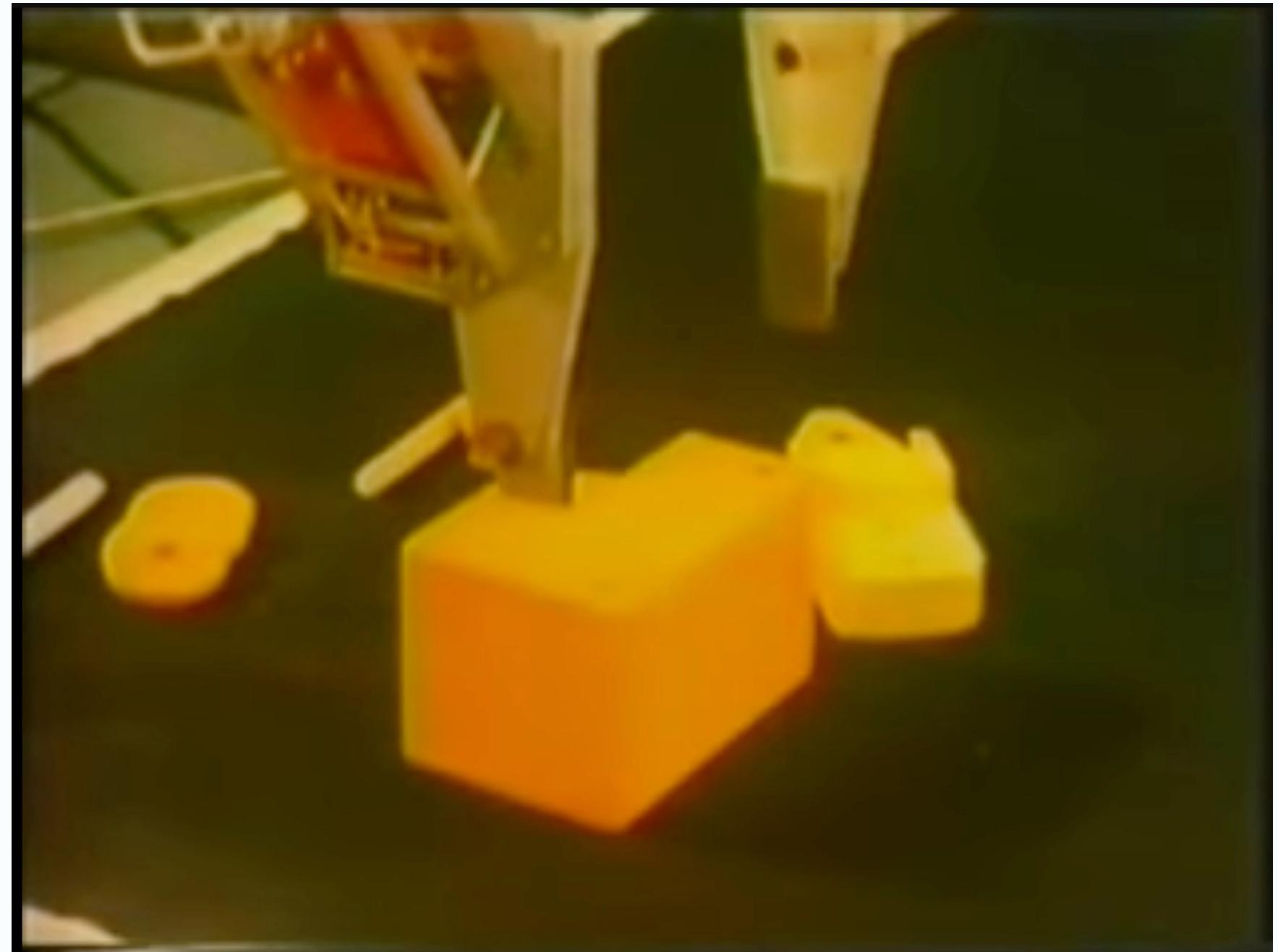
Shanghai Jiao Tong Uni.

# Visual Representations for Generalizable Robotic Manipulation

Oct 27<sup>th</sup>, 2023

- [1] Hansen\*, Yuan\*, **Ze\*** et al., “On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline”, **ICML** 2023.
- [2] **Ze** et al., “Visual Reinforcement Learning with Self-Supervised 3D Representations”, **RA-L** 2023 & **IROS** 2023.
- [3] **Ze** et al., “GNFactor: Multi-Task Real Robot Learning with Generalizable Neural Feature Fields”, **CoRL** 2023 **Oral**.
- [4] **Ze** et al., “H-InDex: Visual Reinforcement Learning with Hand-Informed Representations for Dexterous Manipulation”, **NeurIPS** 2023.
- [5] Yang\*, **Ze\*** et al., “MoVie: Visual Model-Based Policy Adaptation for View Generalization”, **NeurIPS** 2023.

# Generalization in Robotics



**NOT JUST:**

**1 task**

**1 scene**

**1 object**

The 1973 Lighthill debate on Artificial Intelligence:  
"The general purpose robot is a **mirage**"

# Generalization in Computer Vision

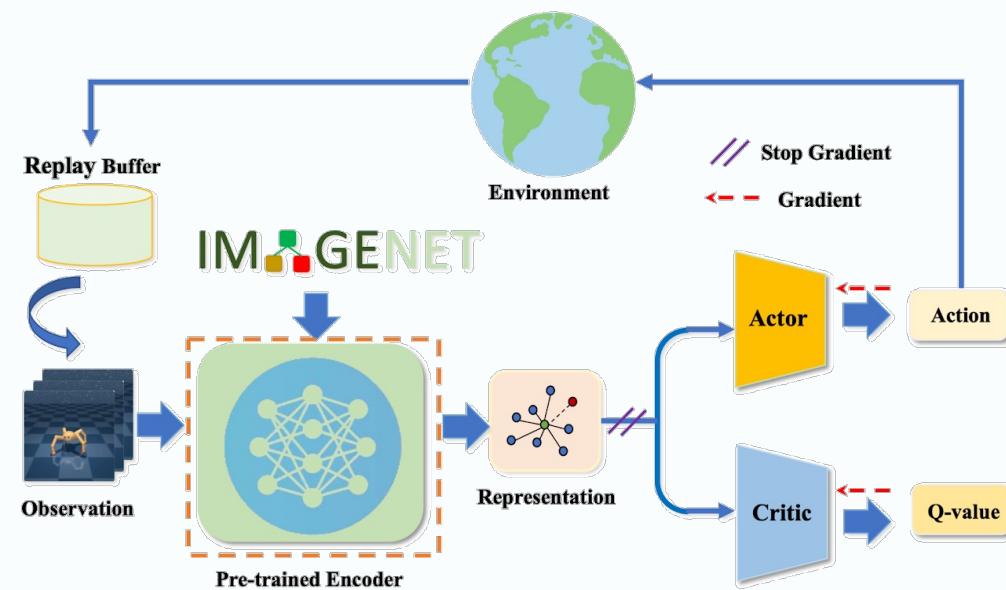


**Segment Anything Model**  
Meta, 2023

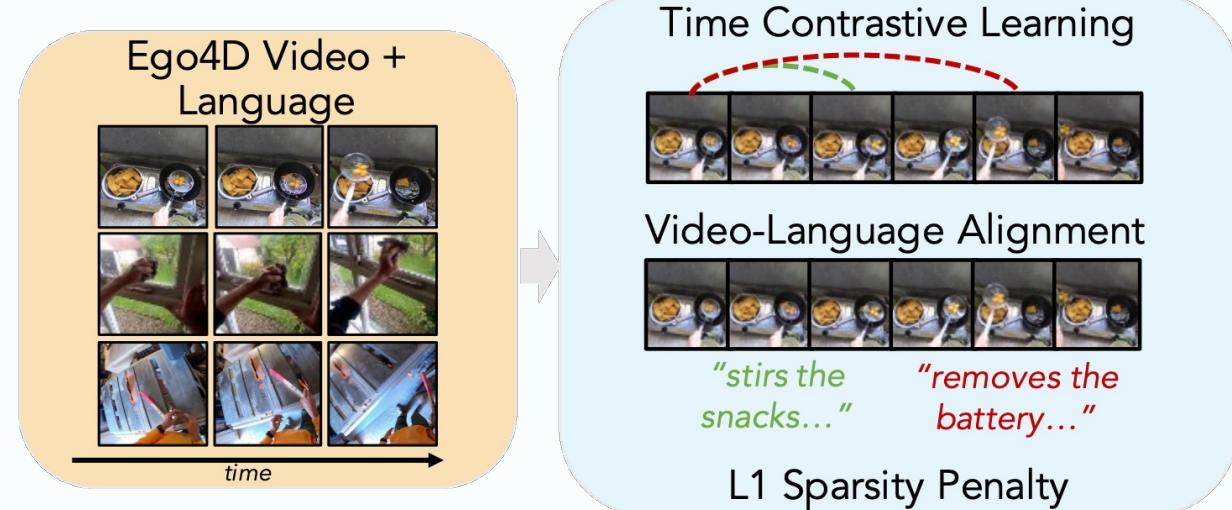


**DALL·E 3**  
OpenAI, 2023

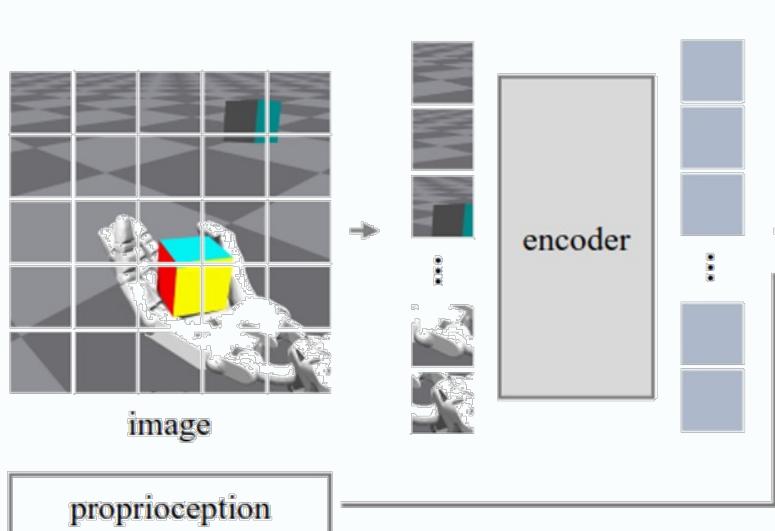
# 2D Visual Representations for Robotics



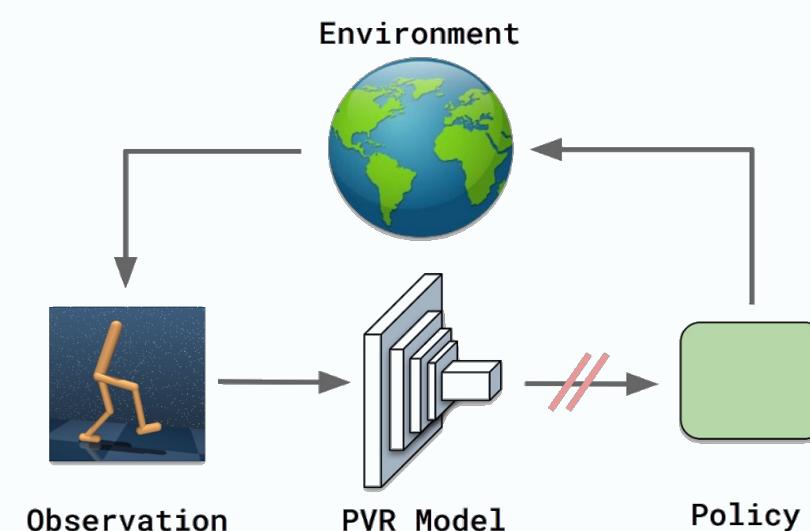
Yuan et al., 2022



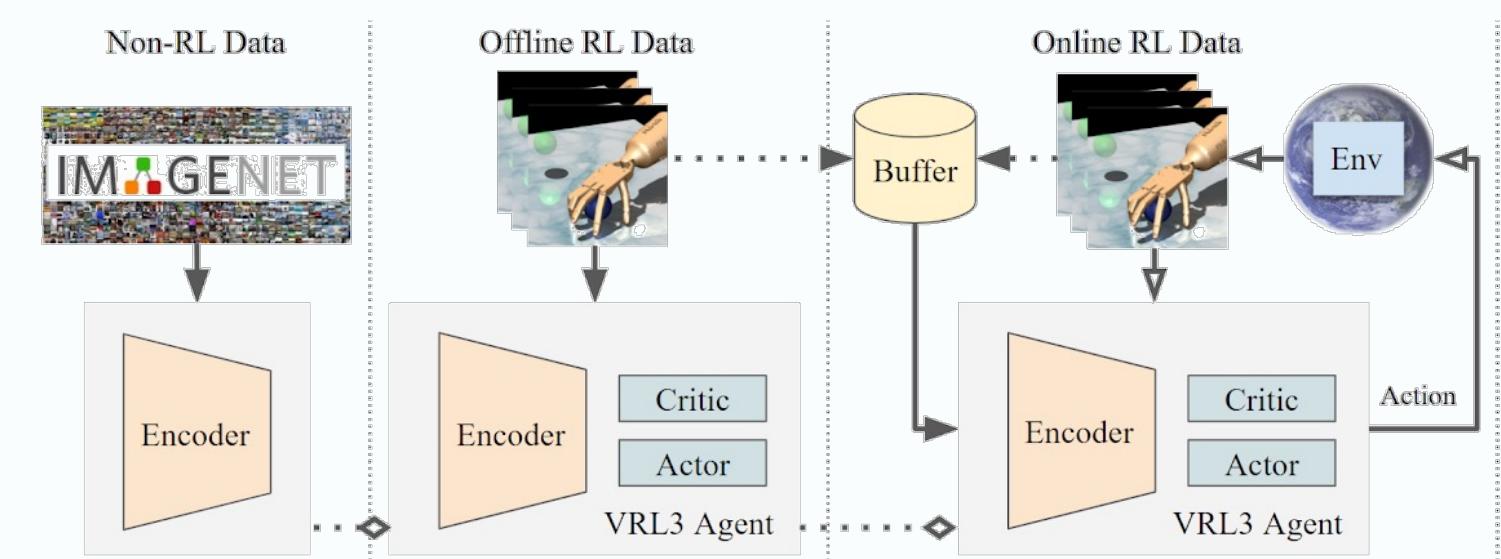
New Environment, New Tasks



Xiao et al., 2022

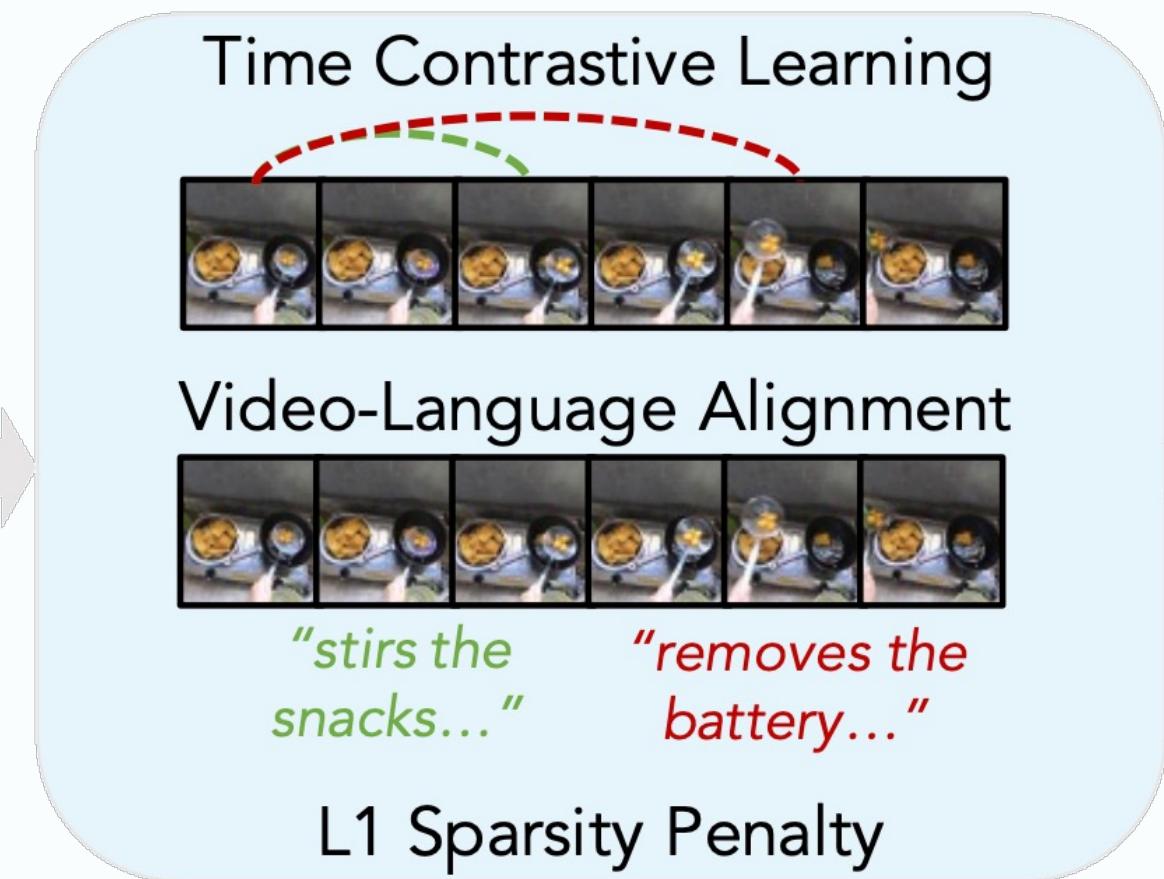


Parisi et al., 2022



Wang et al., 2022

# 2D Visual Representations for Robotics



Pre-Trained R3M Representation

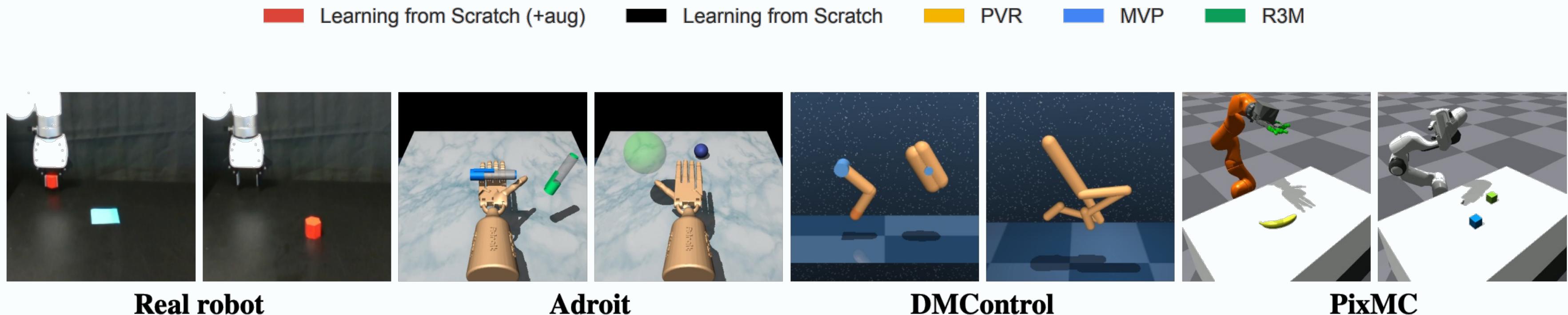


Nair et al., 2022

- How much are they helping manipulation?
- Actually, **not too much.**

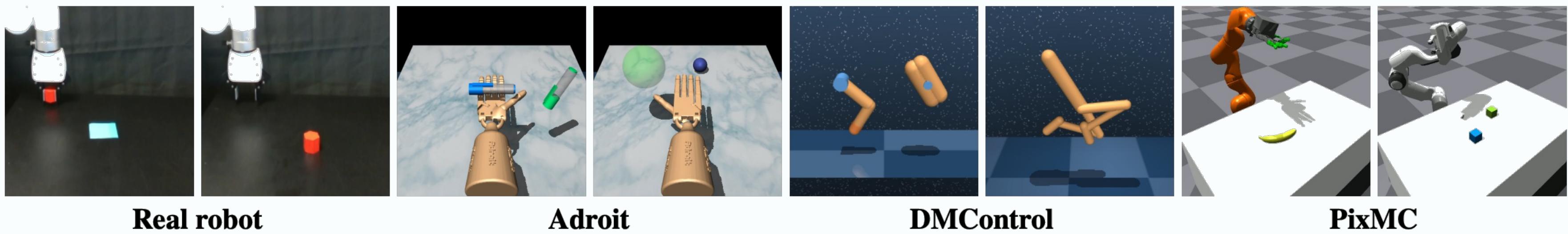
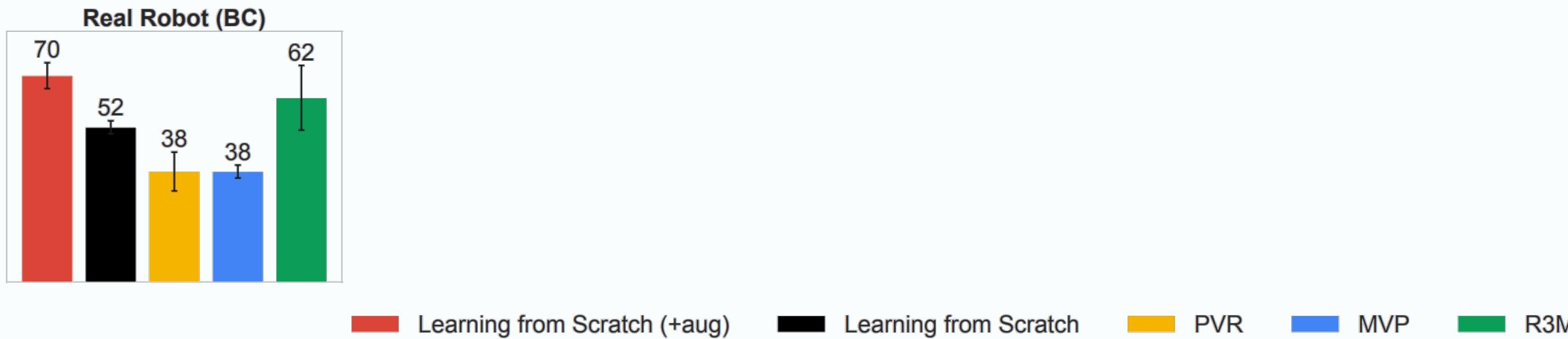
# Revisiting Visual Representations for Visuomotor Control

[1] Hansen\*, Yuan\*, Ze\* et al., On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline, ICML 2023.



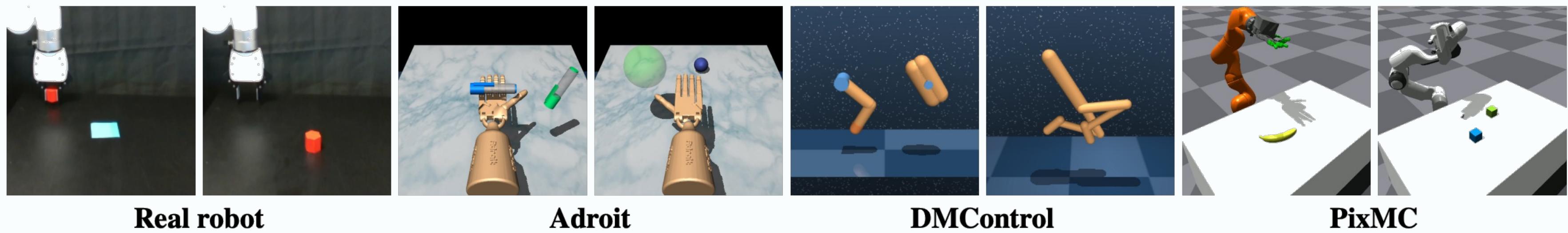
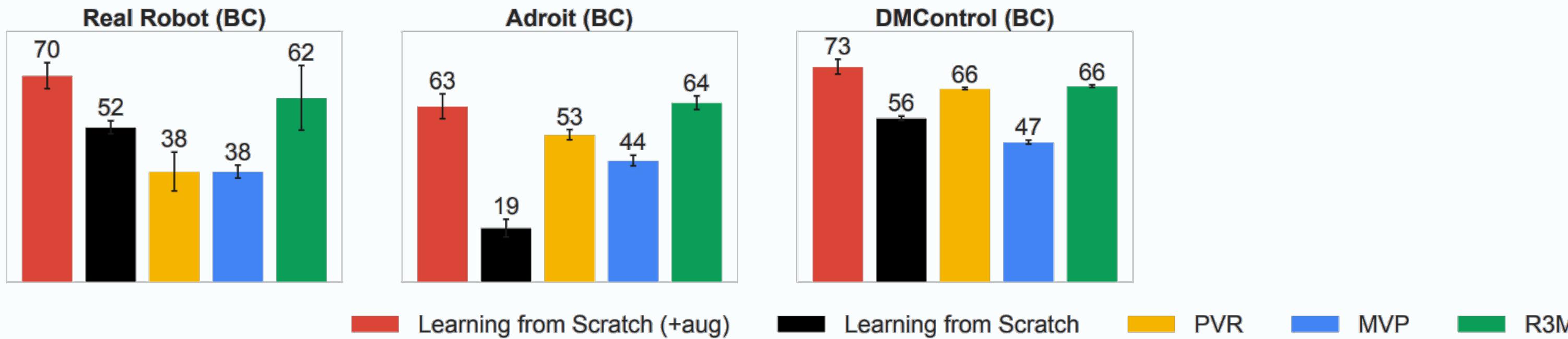
# Revisiting Visual Representations for Visuomotor Control

[1] Hansen\*, Yuan\*, Ze\* et al., On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline, ICML 2023.



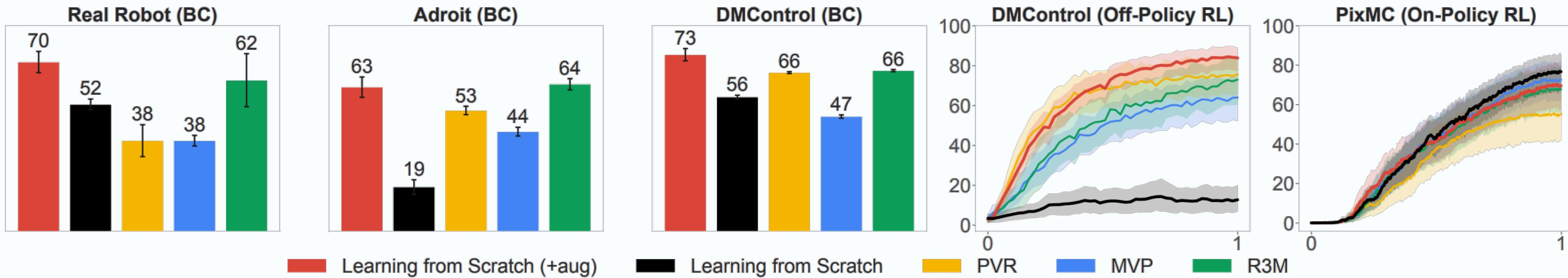
# Revisiting Visual Representations for Visuomotor Control

[1] Hansen\*, Yuan\*, Ze\* et al., On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline, ICML 2023.



# Revisiting Visual Representations for Visuomotor Control

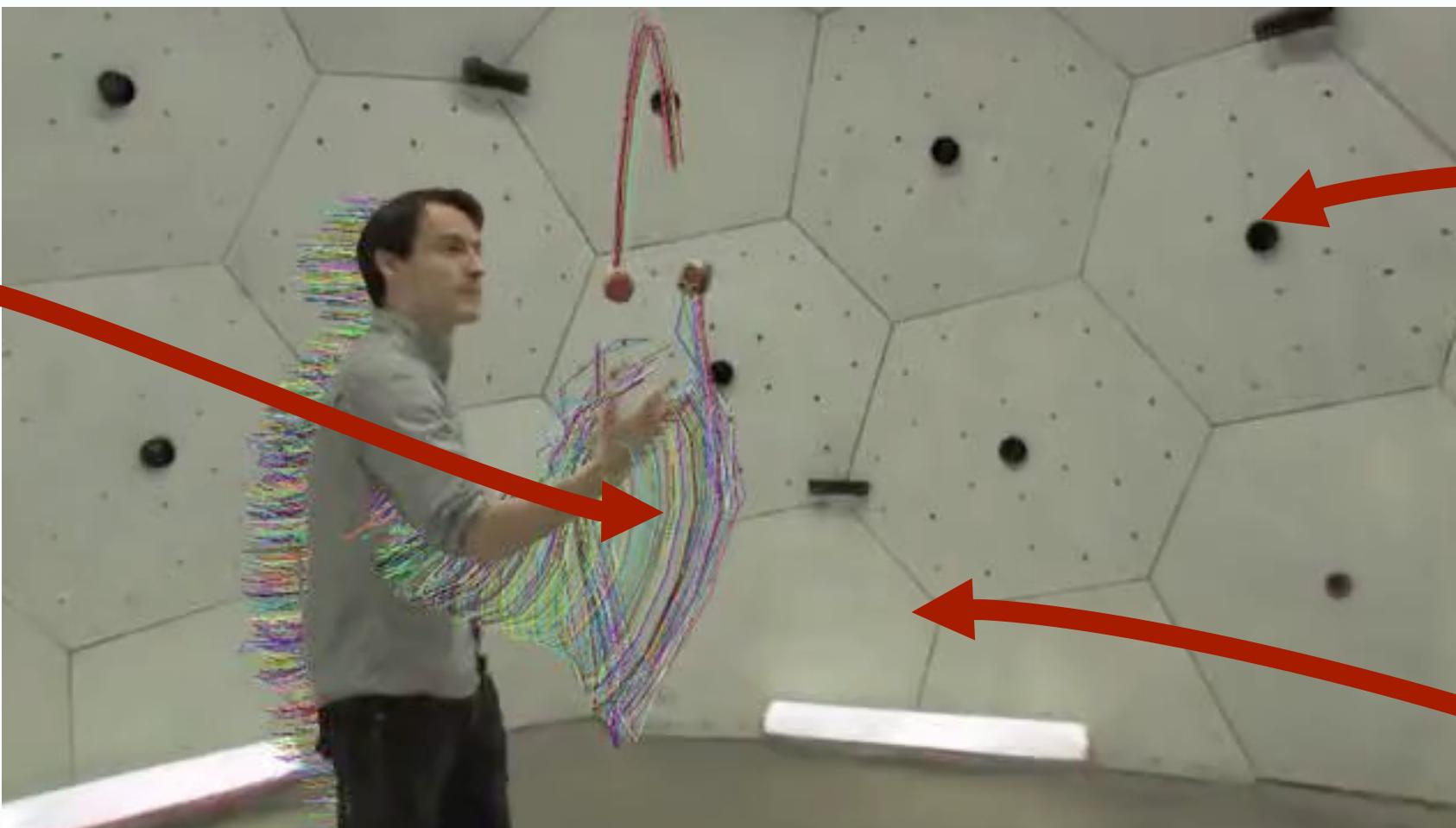
[1] Hansen\*, Yuan\*, Ze\* et al., On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline, ICML 2023.



# What is the Missing Point?

The world itself contains rich **prior**.

**Dynamics Prior**  
How environment transits



**Geometric Prior**  
3D scene, object, ...

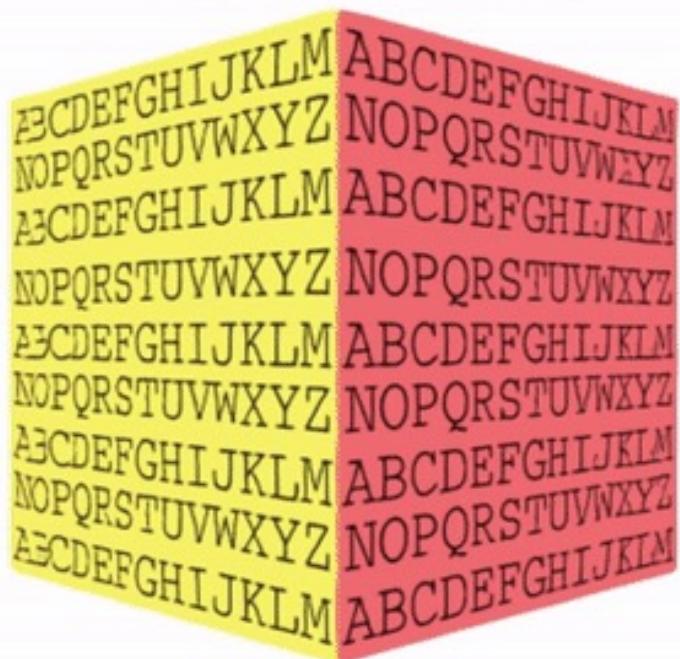
**Human Prior**  
body pose, hand pose, ...

Luiten et al., 2023

# Rich Visual Prior from the World

**Geometric Prior**

**Voxel**



**Human Prior**

**Dynamics Prior**

Sitzmann et al., 2019

# Rich Visual Prior from the World

**Geometric Prior**

**NeRF**



**Human Prior**

**Dynamics Prior**

Barron et al., 2022

# Rich Visual Prior from the World

## Geometric Prior

NeRF



Barron et al., 2022

## Human Prior

Human Pose



Cai et al., 2023

## Dynamics Prior

# Rich Visual Prior from the World

## Geometric Prior

NeRF



Barron et al., 2022

## Human Prior

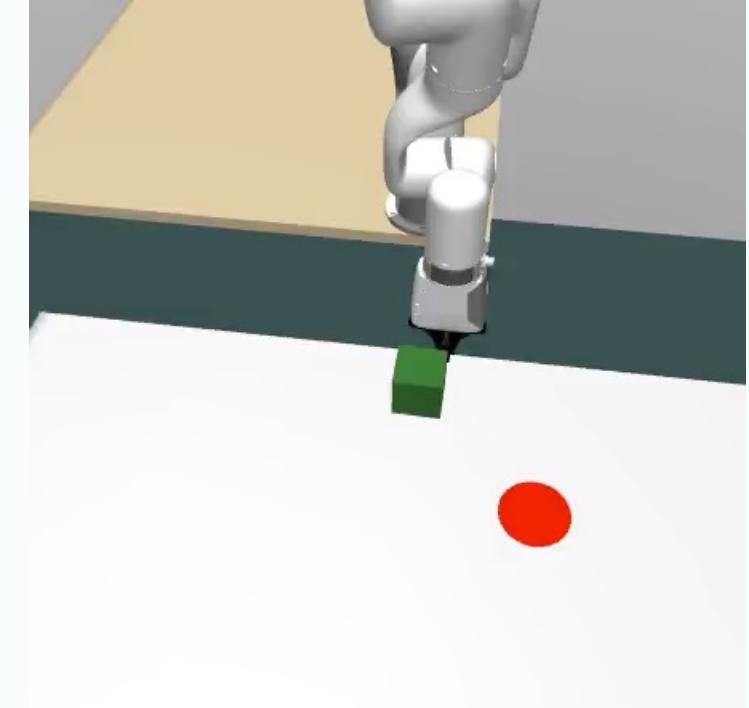
Human Pose



Cai et al., 2023

## Dynamics Prior

Environment Transition



Yang et al., 2023

# Visual Representations for Generalizable Robotic Manipulation

1

## Geometric Prior

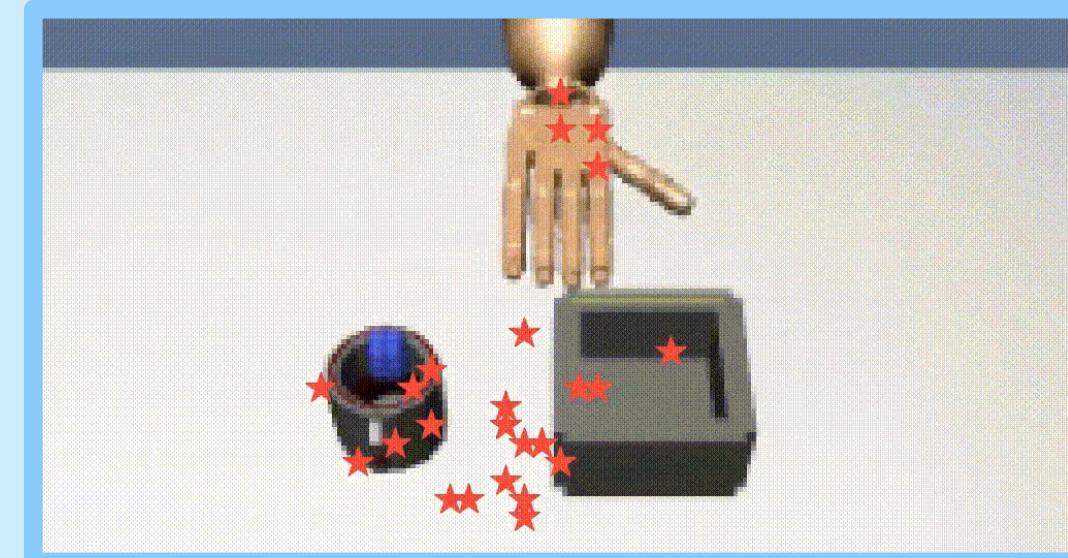


[2] Ze et al., “Visual Reinforcement Learning with Self-Supervised 3D Representations”, **RA-L** 2023 & **IROS** 2023.

[3] Ze et al., “GNFactor: Multi-Task Real Robot Learning with Generalizable Neural Feature Fields”, **CoRL** 2023 Oral.

2

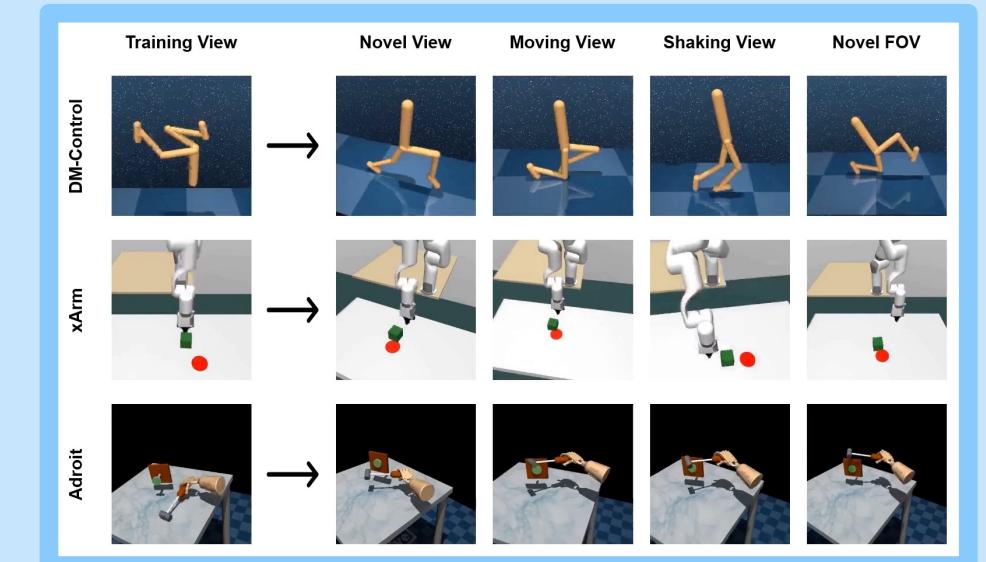
## Human Prior



[4] Ze et al., “H-InDex: Visual Reinforcement Learning with Hand-Informed Representations for Dexterous Manipulation”, **NeurIPS** 2023.

3

## Dynamics Prior

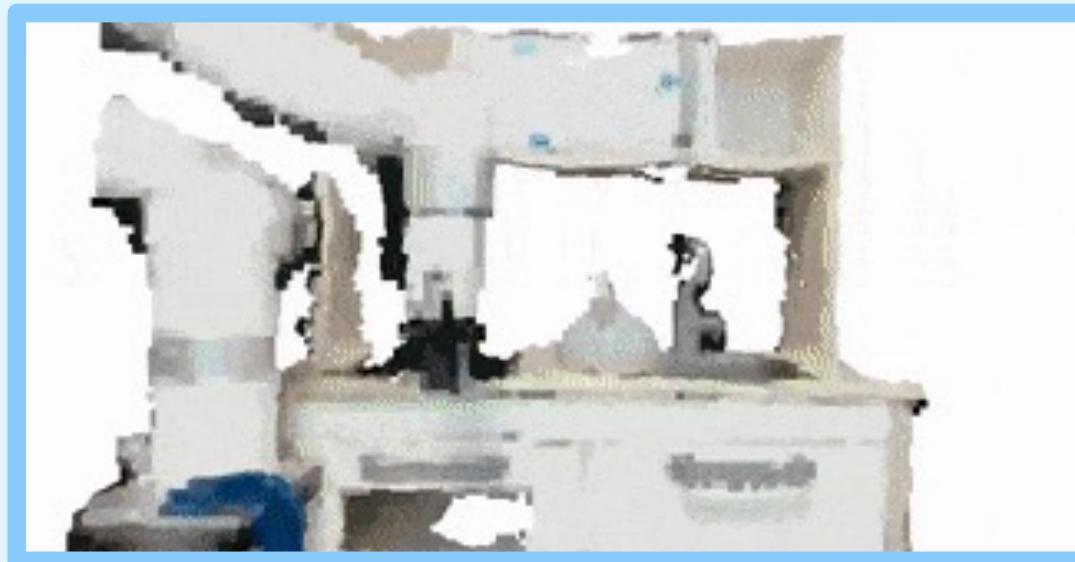


[5] Yang\*, Ze\* et al., “MoVie: Visual Model-Based Policy Adaptation for View Generalization”, **NeurIPS** 2023.

# Visual Representations for Generalizable Robotic Manipulation

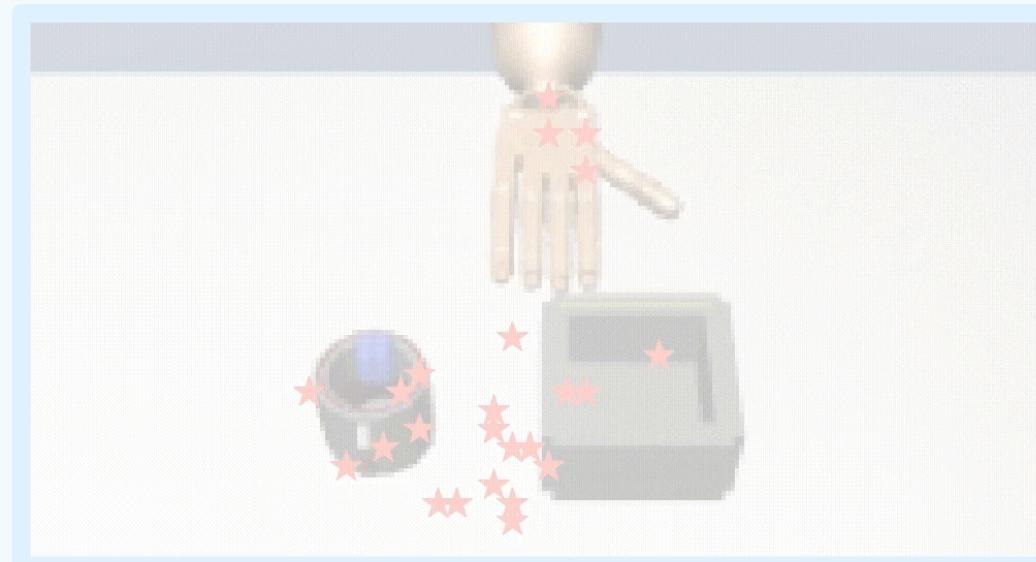
1

## Geometric Prior



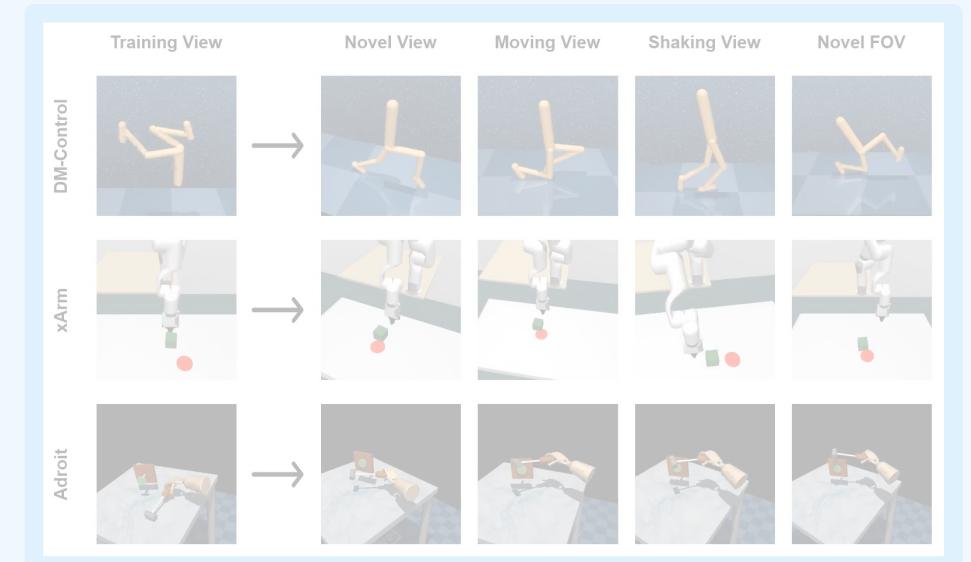
2

## Human Prior



3

## Dynamics Prior



[2] Ze et al., “Visual Reinforcement Learning with Self-Supervised 3D Representations”, **RA-L** 2023 & **IROS** 2023.

[3] Ze et al., “GNFactor: Multi-Task Real Robot Learning with Generalizable Neural Feature Fields”, **CoRL** 2023 Oral.

[4] Ze et al., “H-InDex: Visual Reinforcement Learning with Hand-Informed Representations for Dexterous Manipulation”, **NeurIPS** 2023.

[5] Yang\*, Ze\* et al., “MoVie: Visual Model-Based Policy Adaptation for View Generalization”, **NeurIPS** 2023.

# Visual Reinforcement Learning with Self-Supervised 3D Representations (**RL3D**)

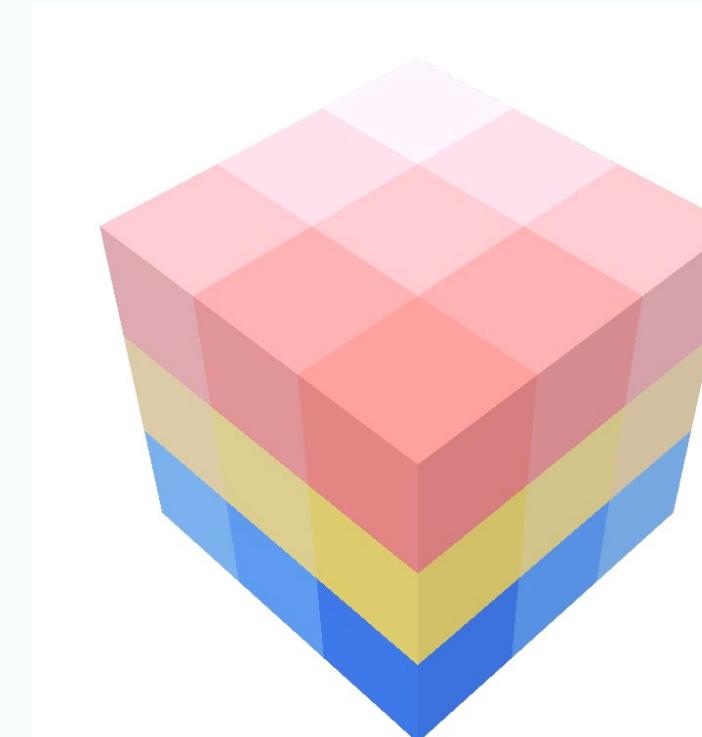
Yanjie Ze<sup>1\*</sup> Nicklas Hansen<sup>2\*</sup> Yinbo Chen<sup>2</sup> Mohit Jain<sup>2</sup> Xiaolong Wang<sup>2</sup>

<sup>1</sup>Shanghai Jiao Tong University <sup>2</sup>UC San Diego

RA-L 2023 & IROS 2023



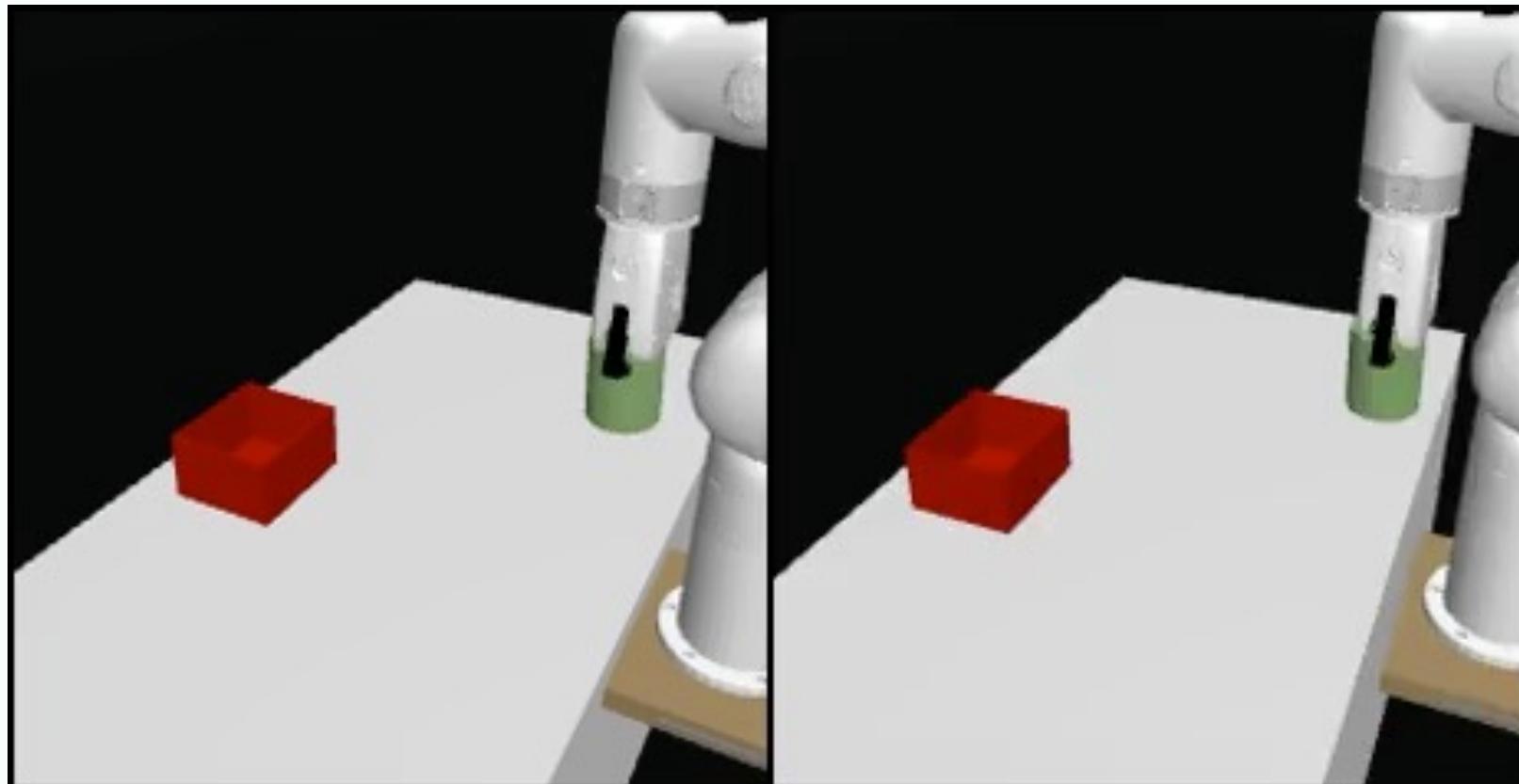
# How about 3D pre-training for visuomotor control?



# One Key Design

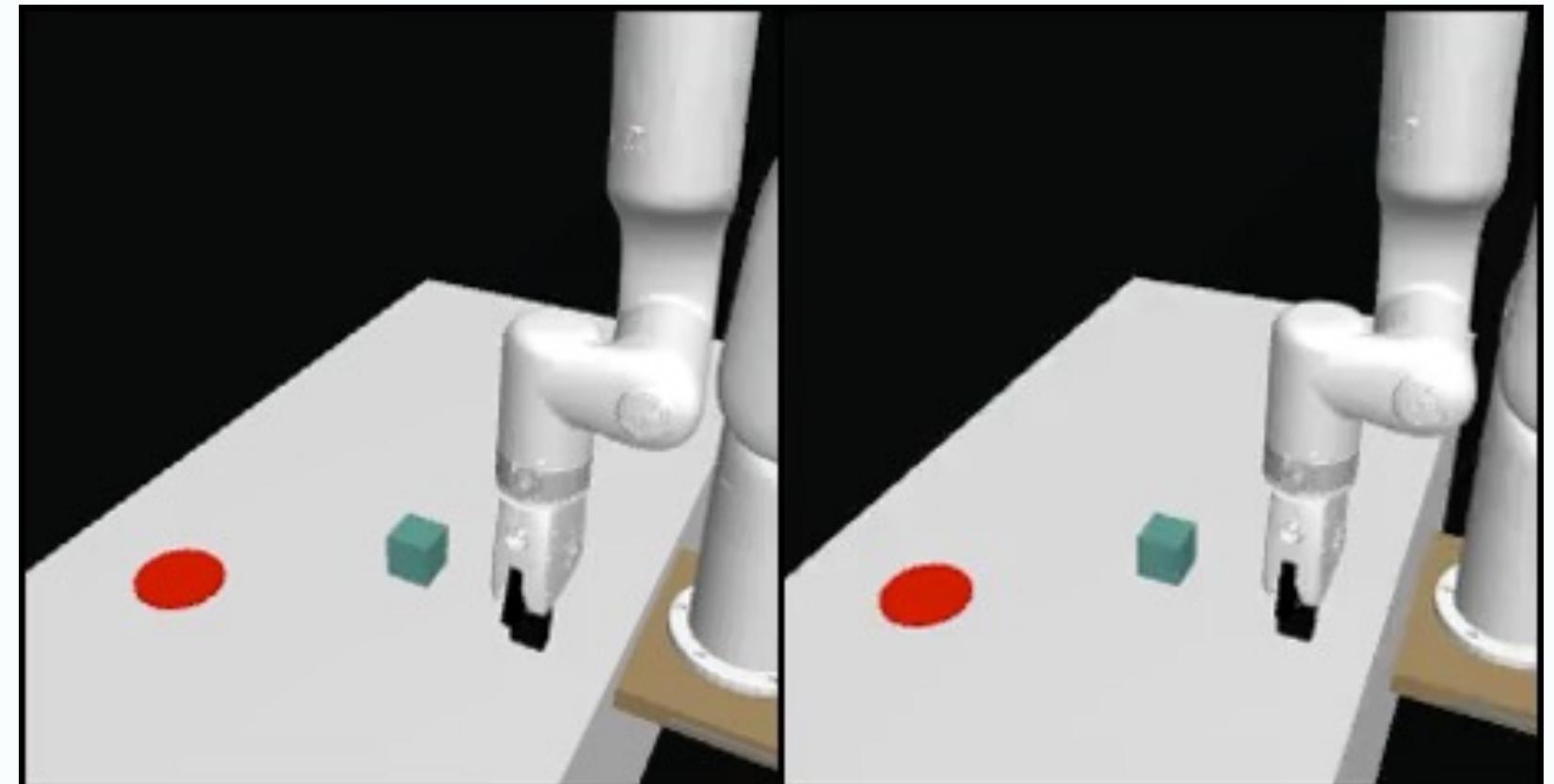
We design **two** cameras:  
one **fixed** camera and one **dynamic** camera.

Static

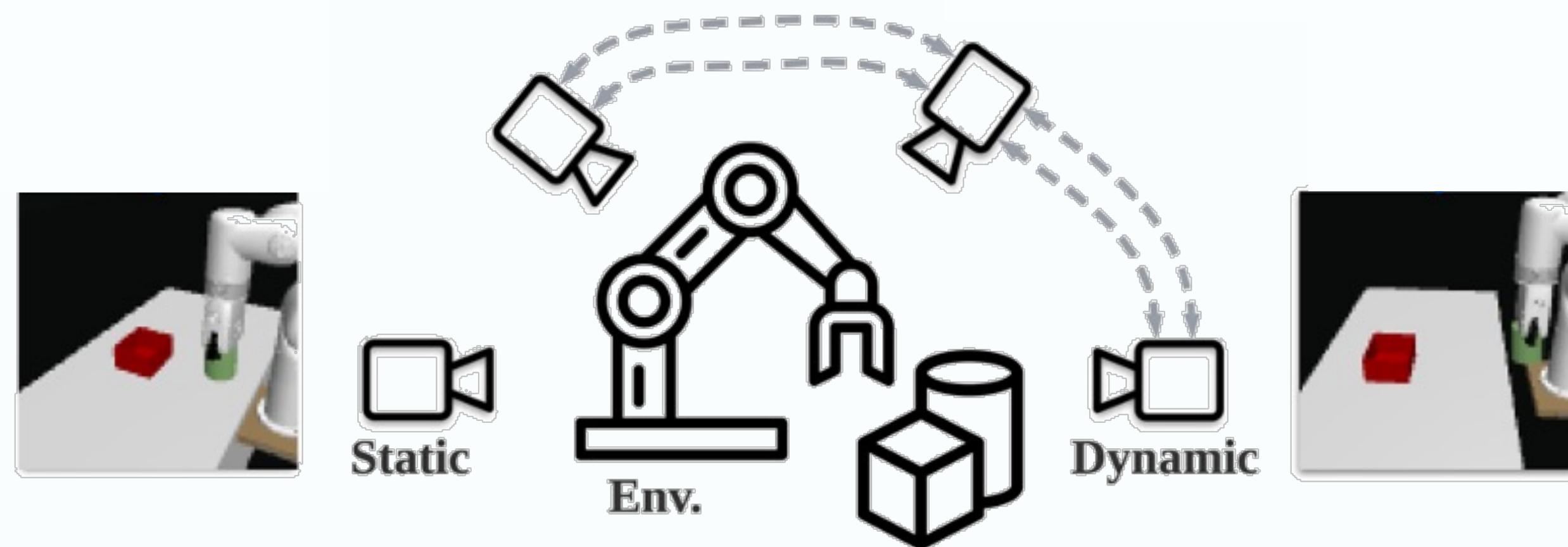


Dynamic

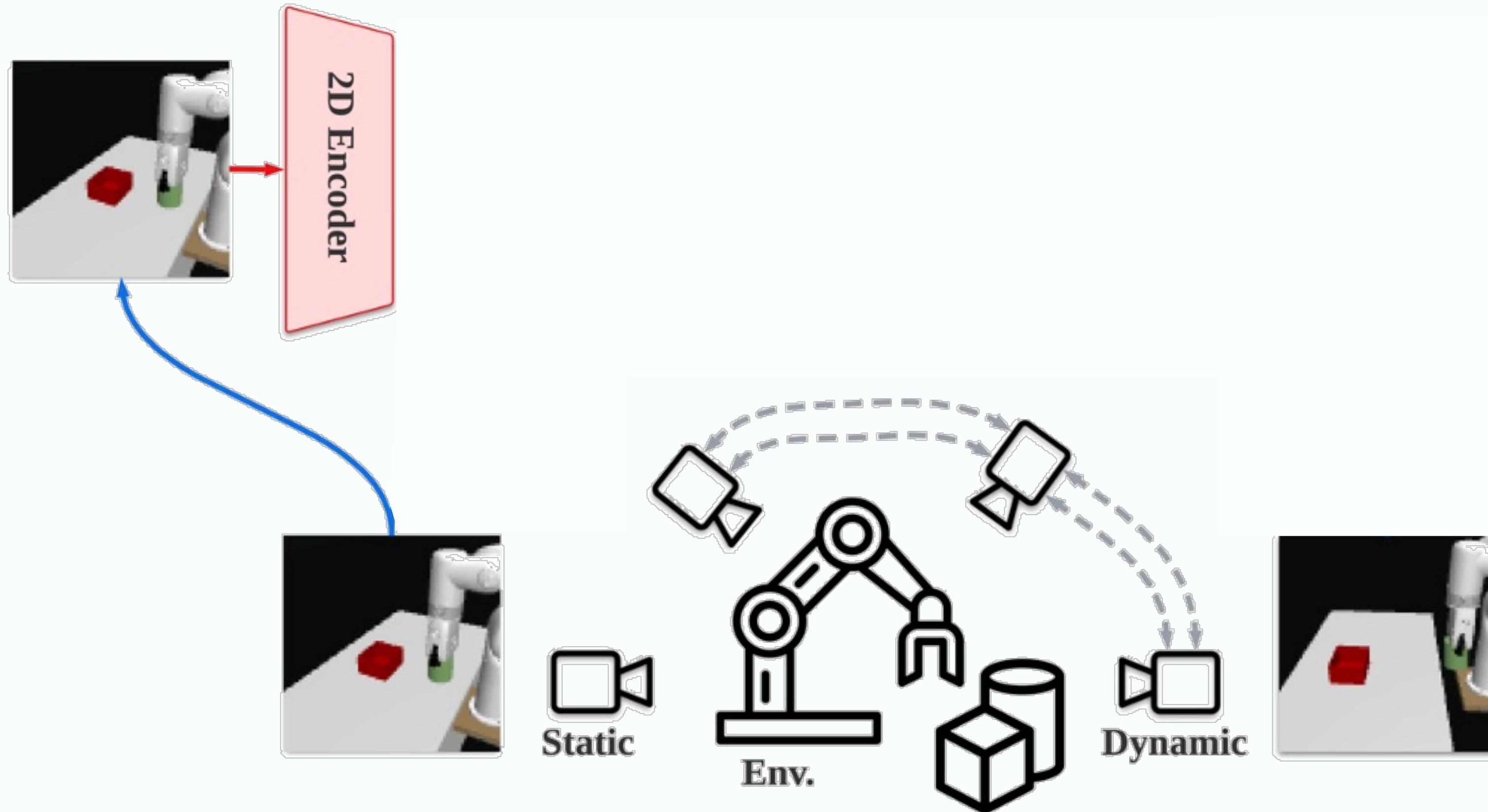
Static



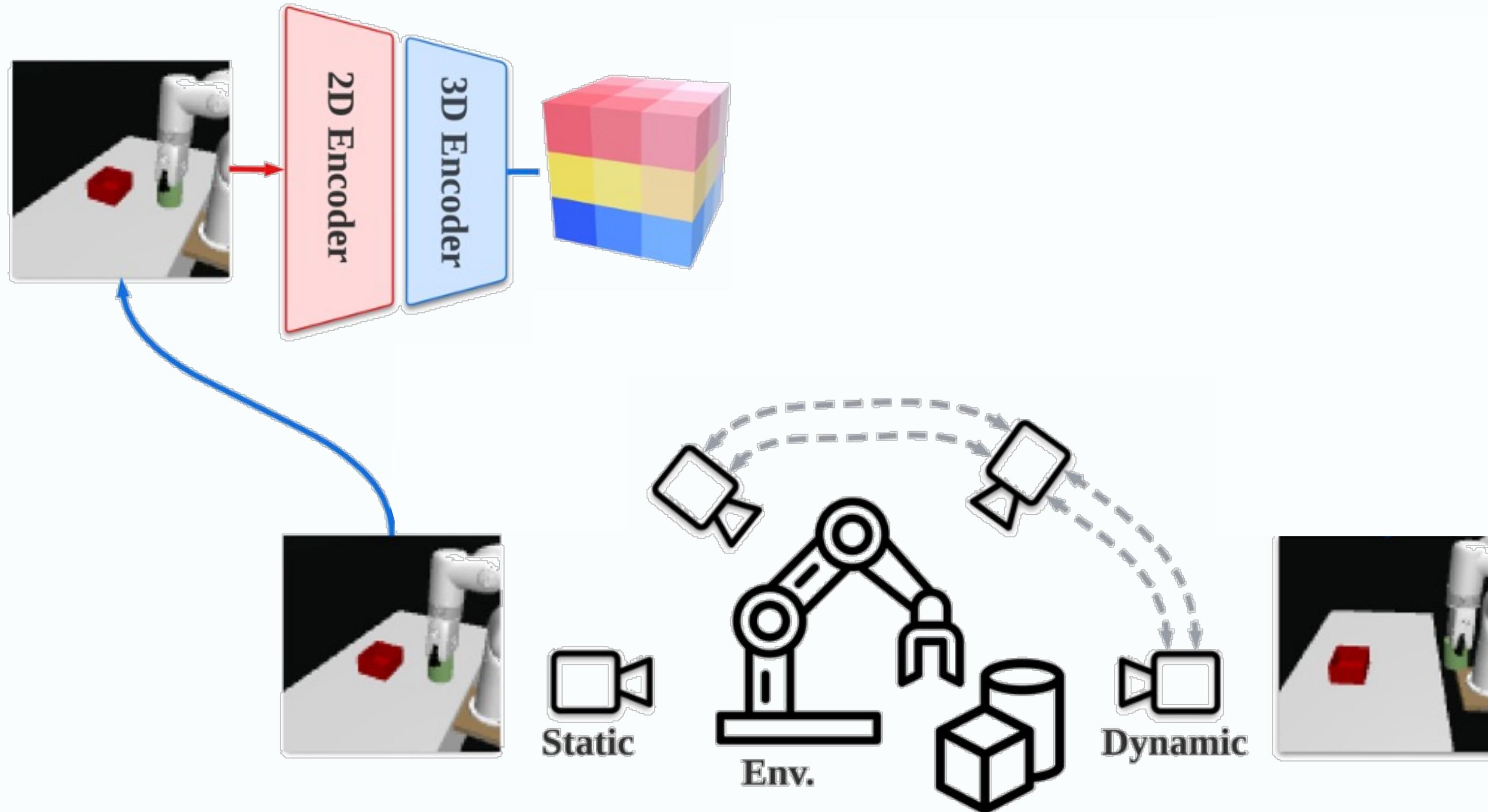
# Algorithm Foundation: Video Auto-encoder



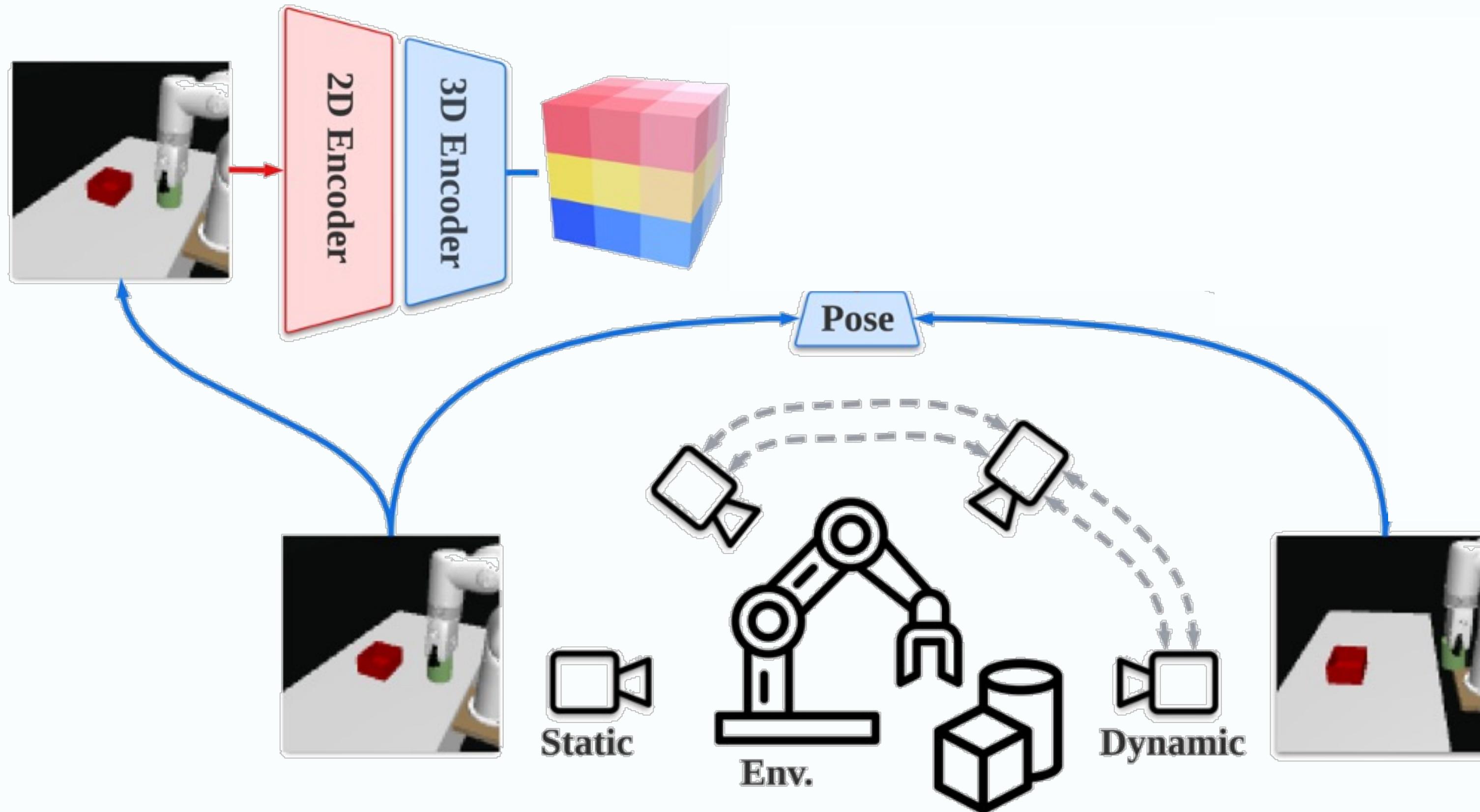
# Algorithm Foundation: Video Auto-encoder



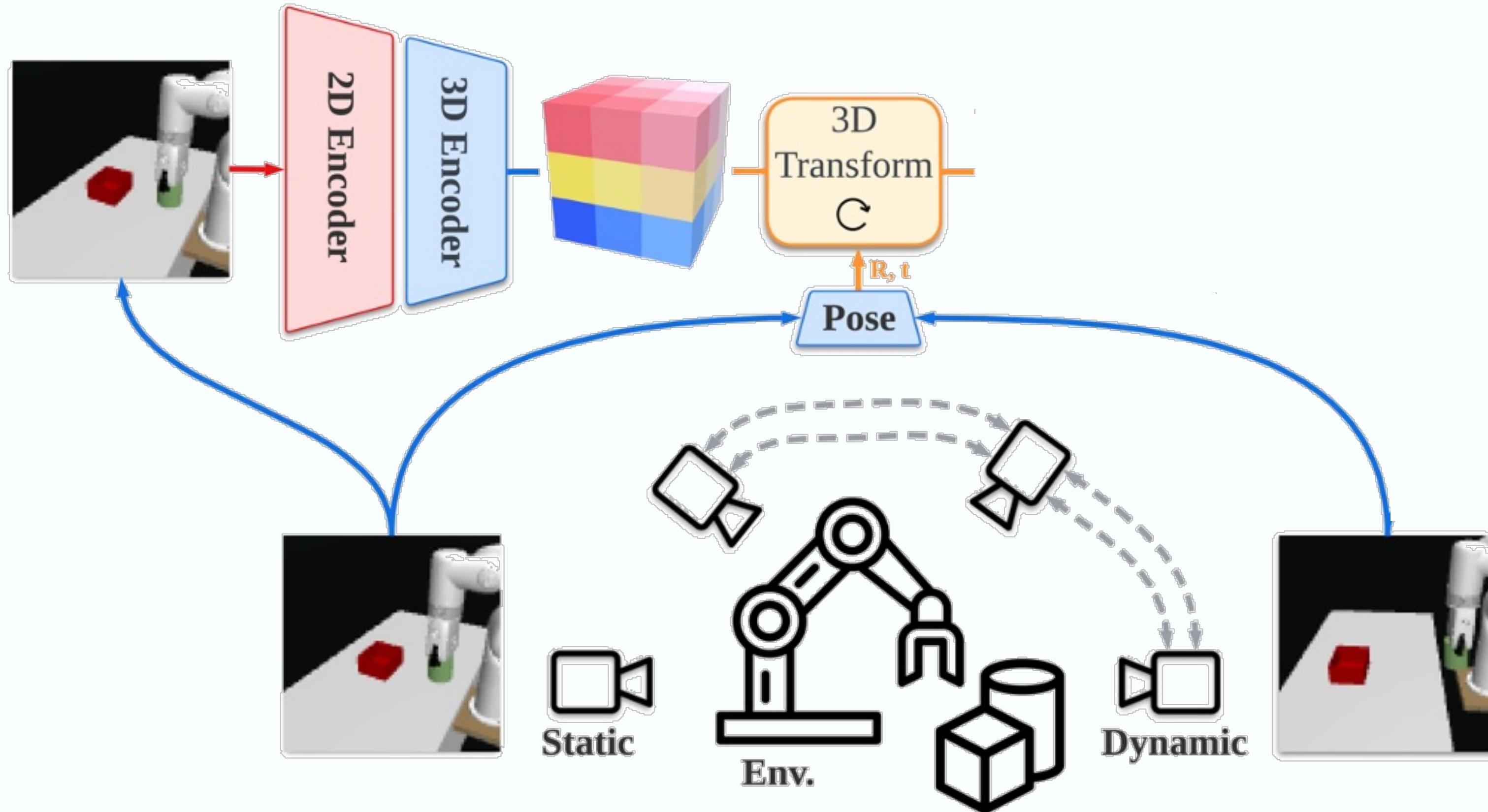
# Algorithm Foundation: Video Auto-encoder



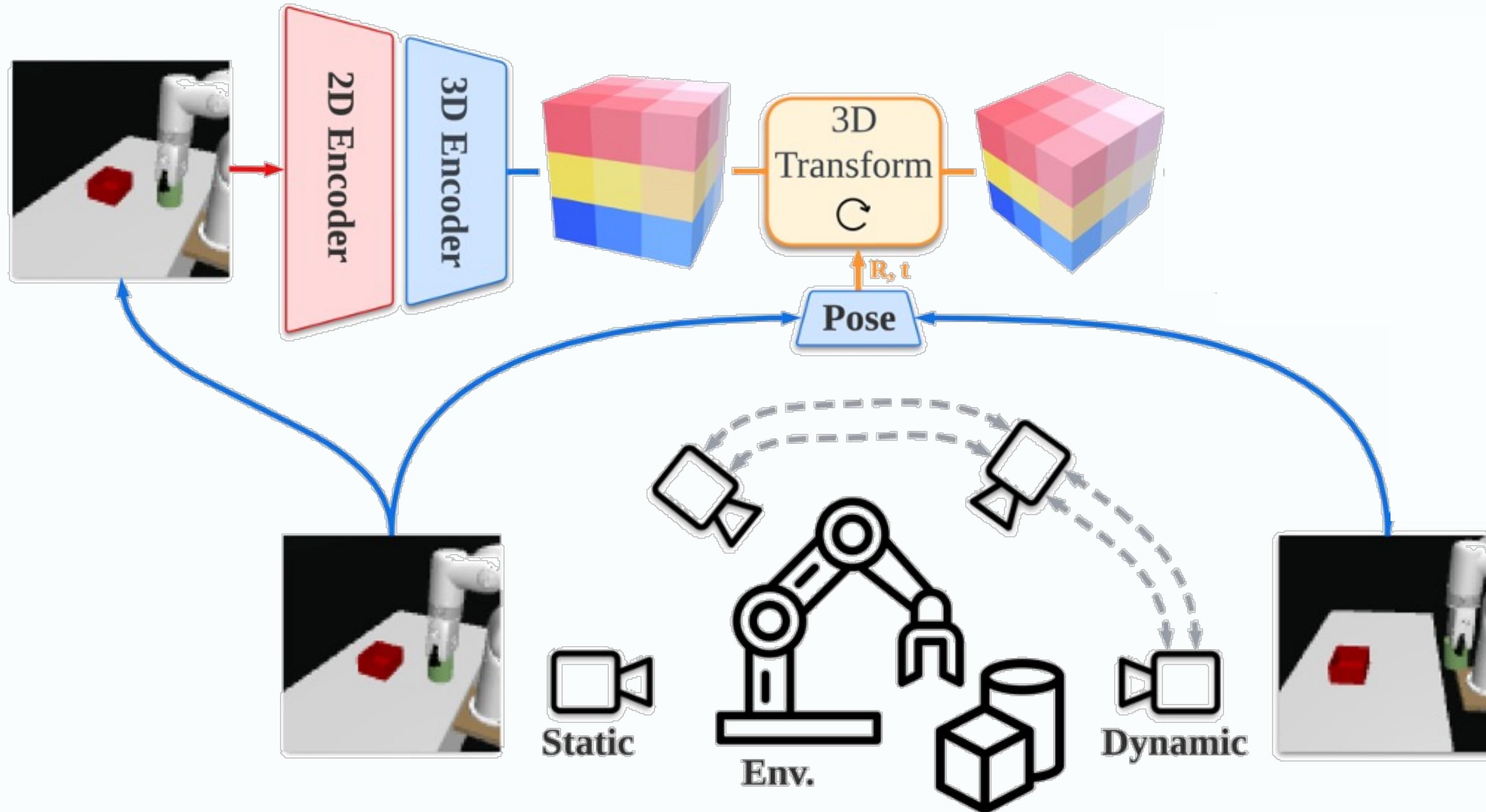
# Algorithm Foundation: Video Auto-encoder



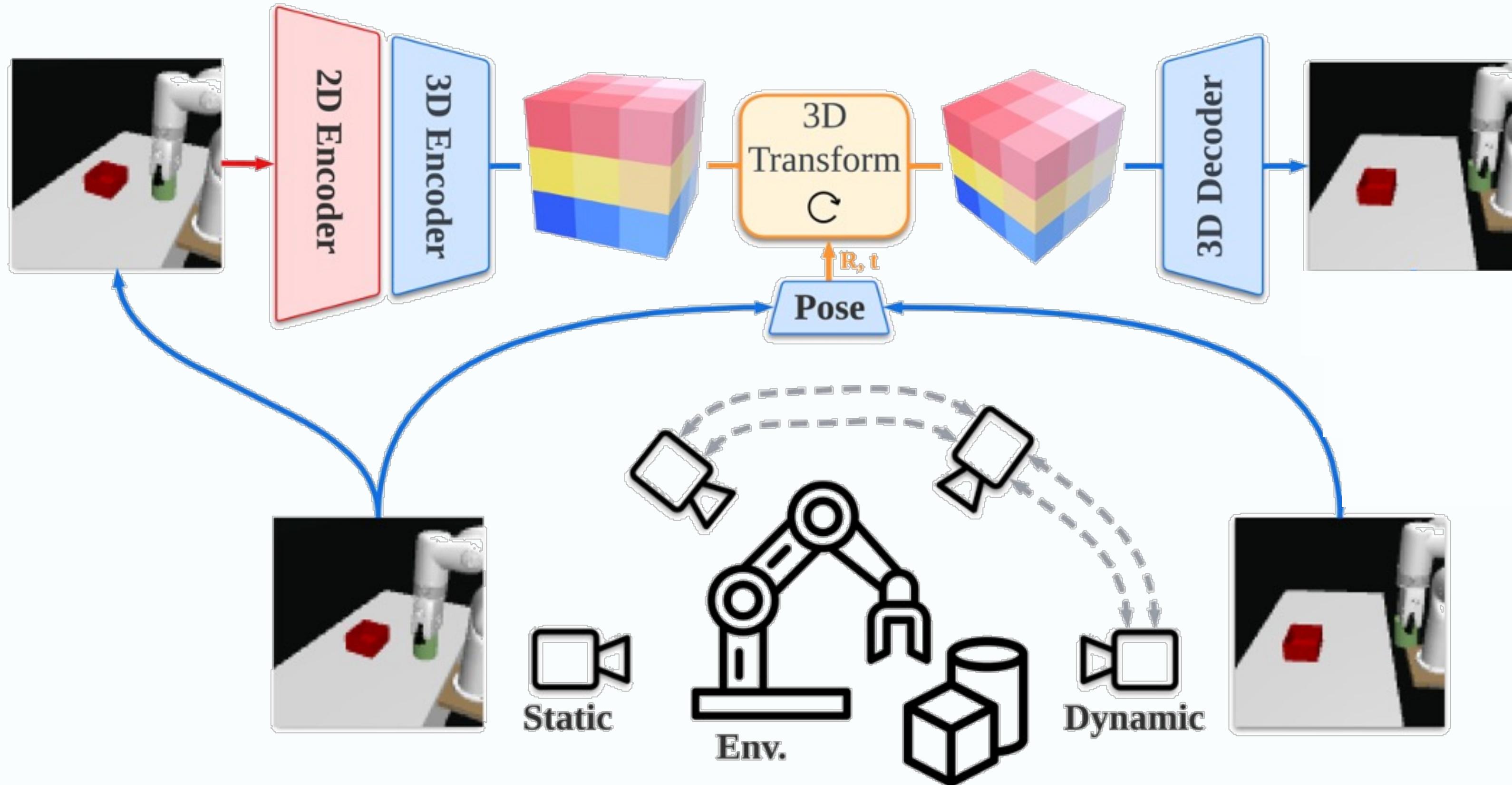
# Algorithm Foundation: Video Auto-encoder



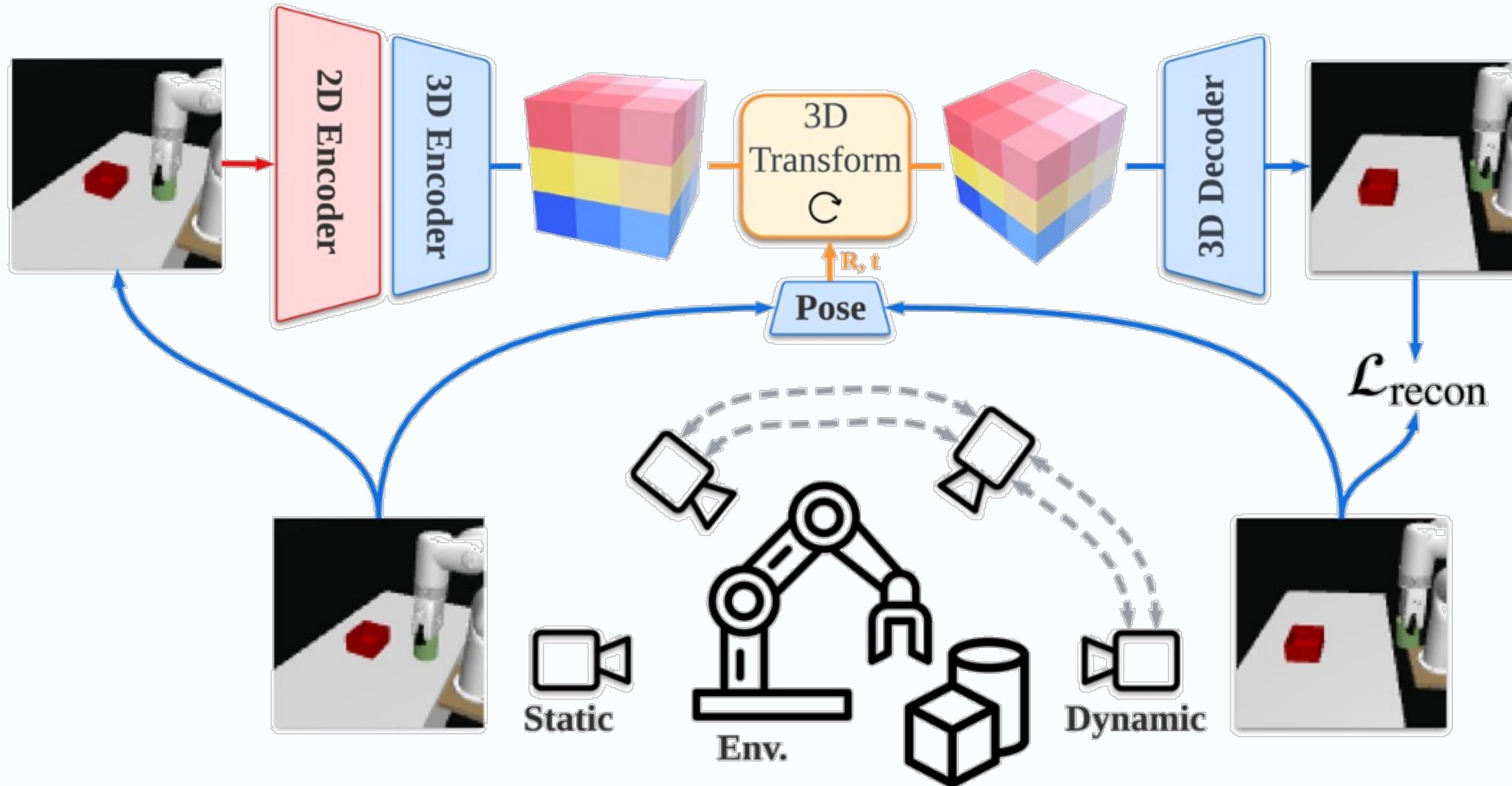
# Algorithm Foundation: Video Auto-encoder



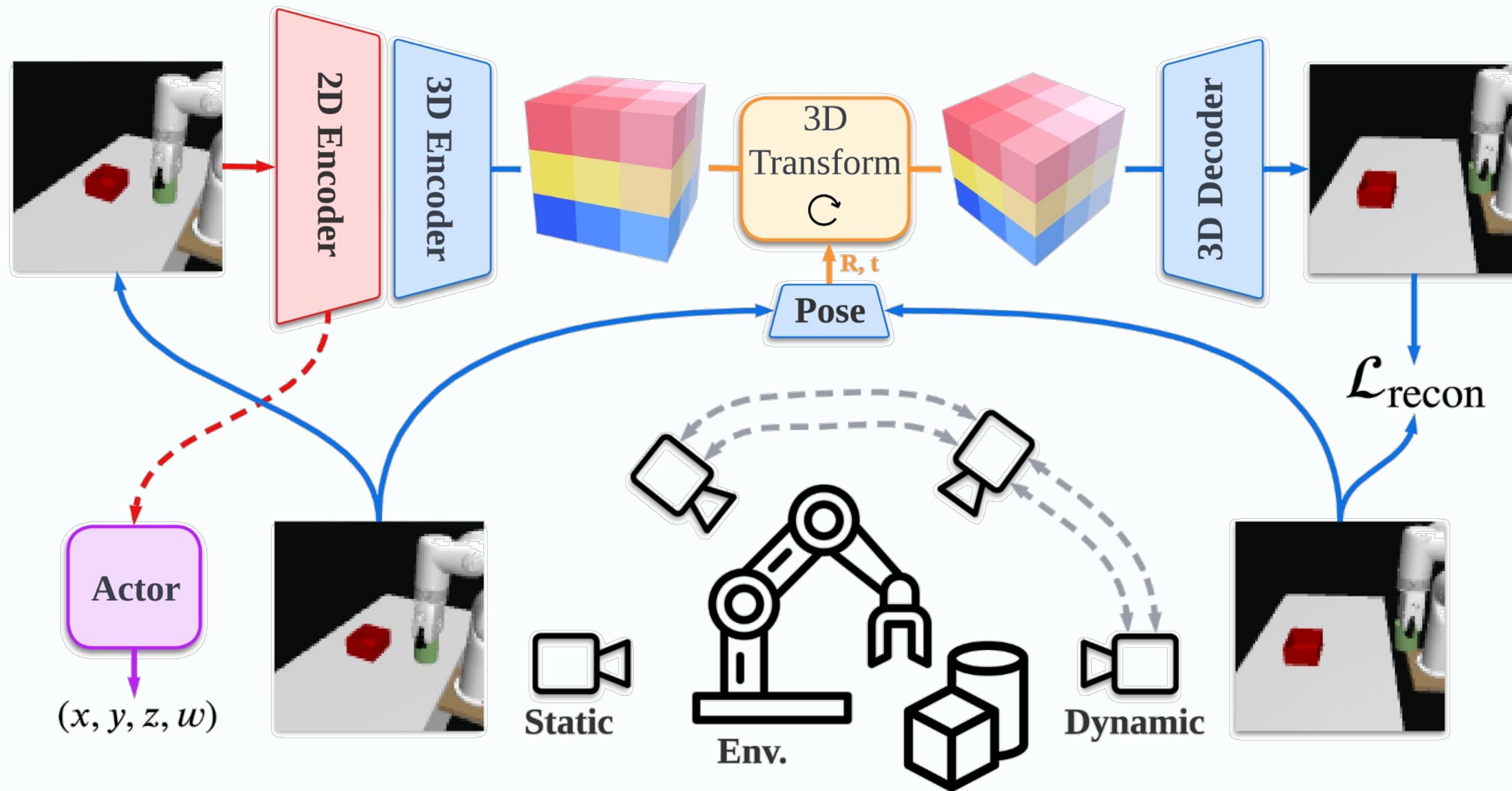
# Algorithm Foundation: Video Auto-encoder



# Algorithm Foundation: Video Auto-encoder



# Algorithm Foundation: Video Auto-encoder



# Framework of RL3D

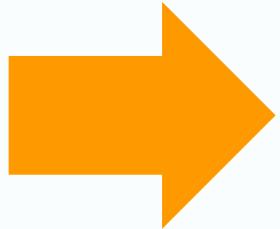
Pretrain on CO3D



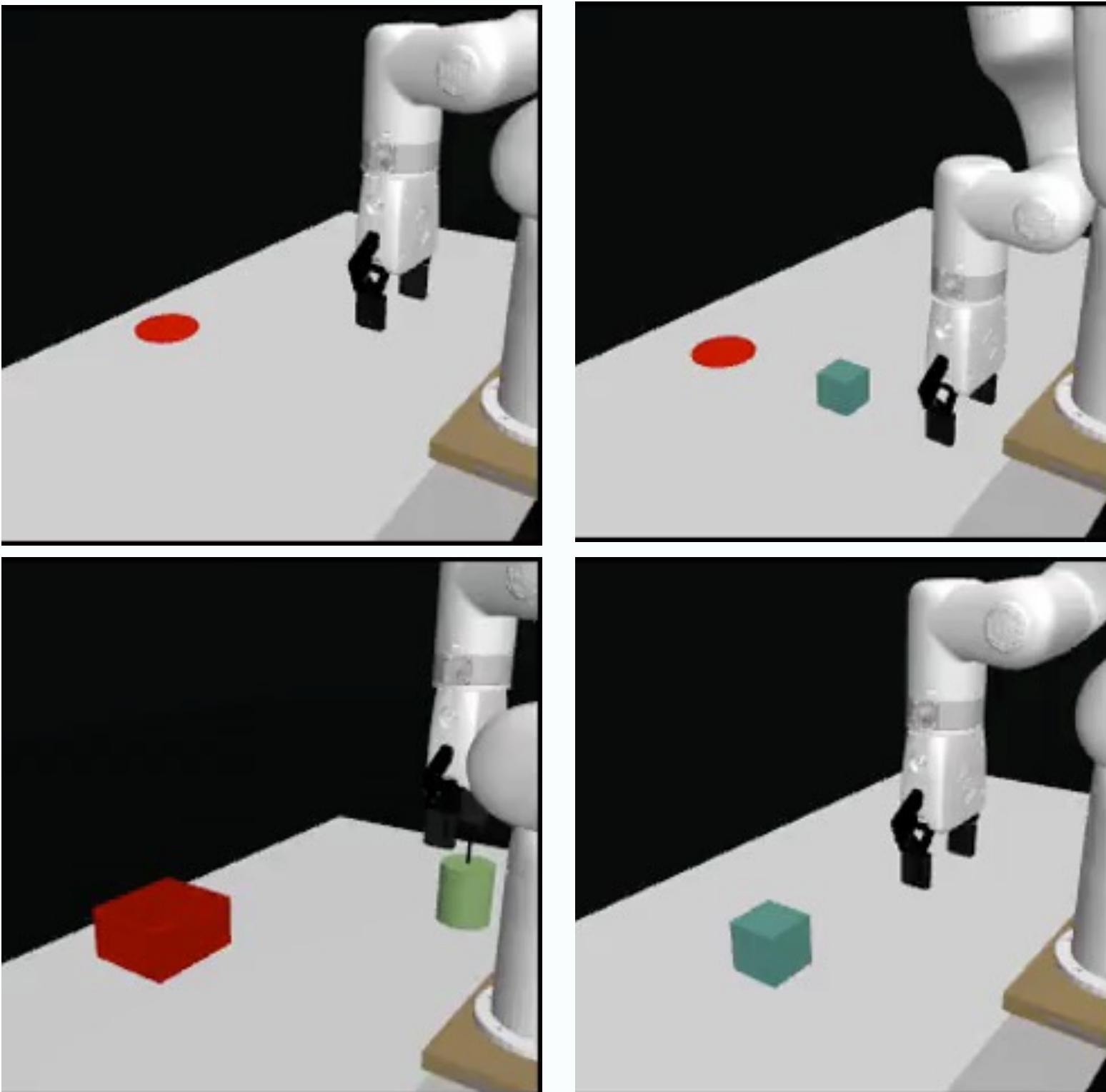
Reizenstein et al, Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction, ICCV 2021.

# Framework of RL3D

Pretrain on CO3D



Finetune on Robot Env

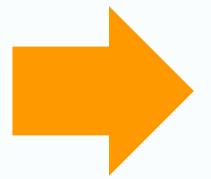
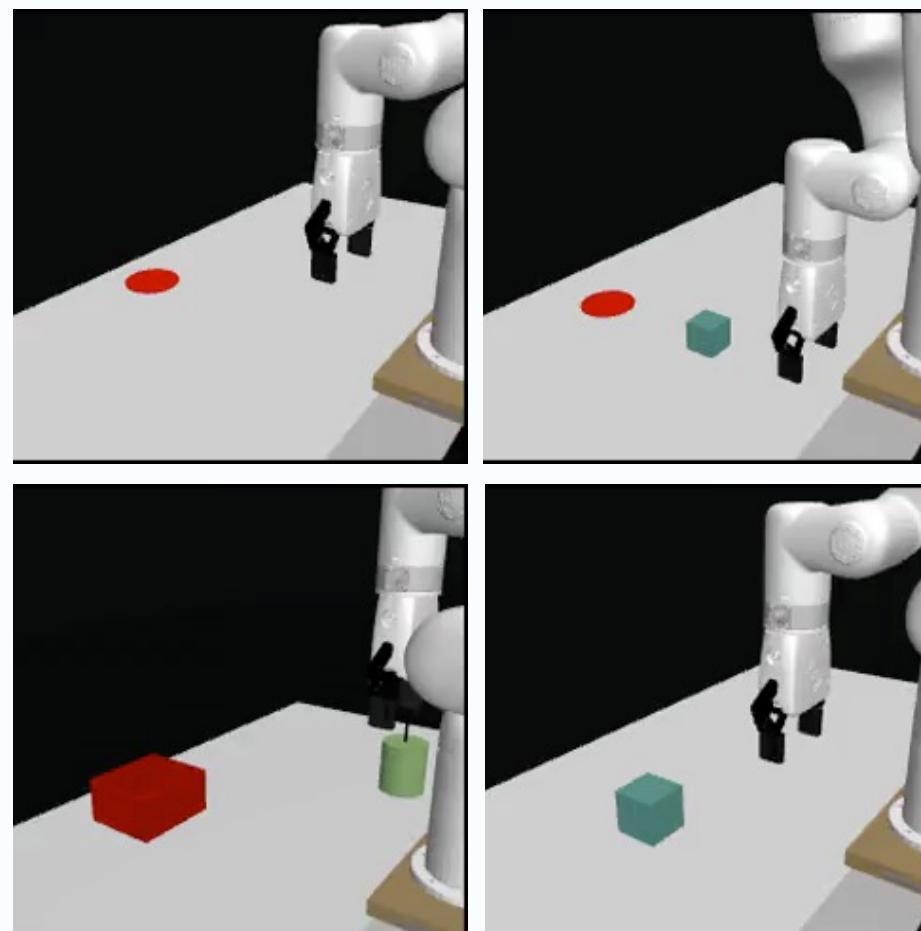


# Framework of RL3D

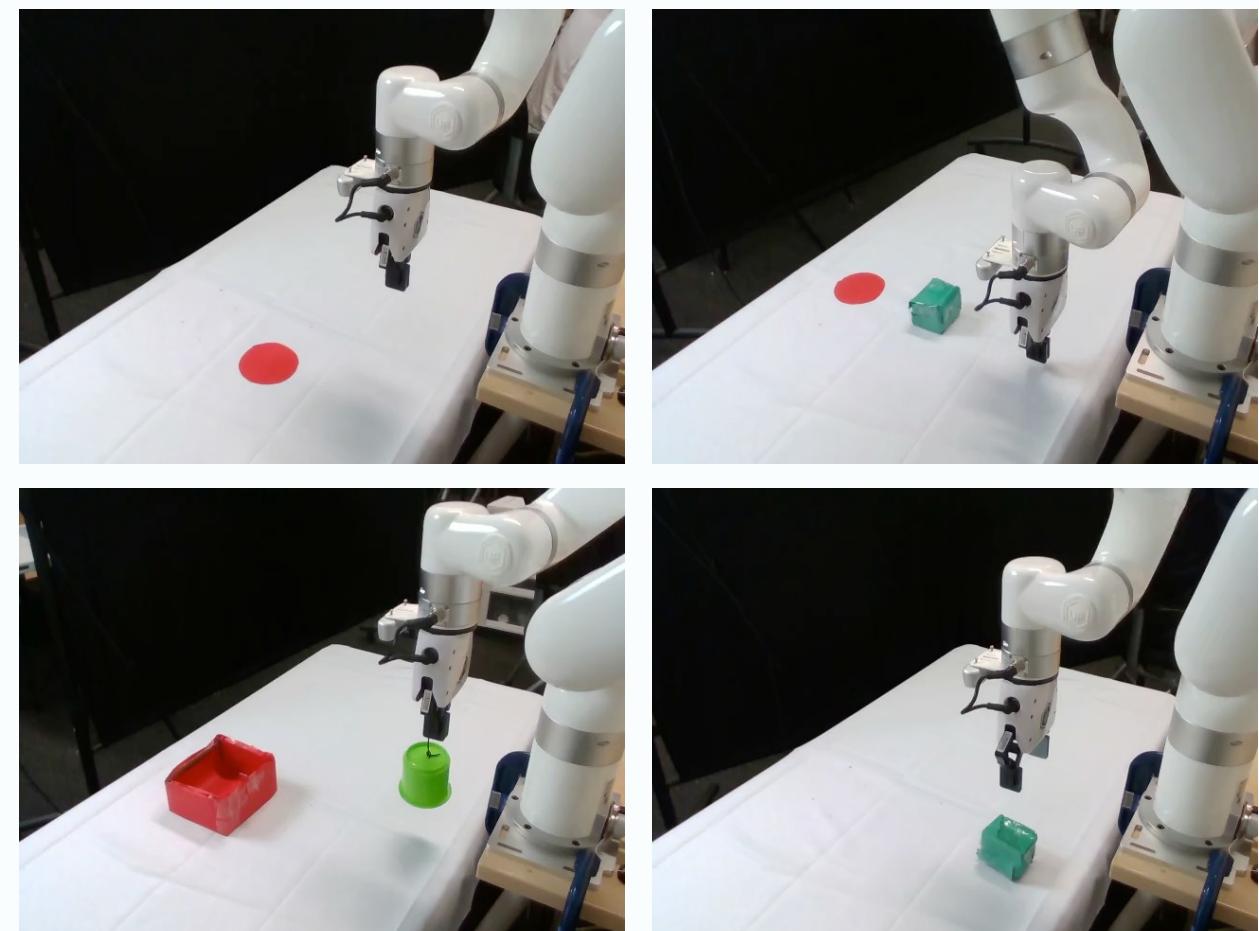
Pretrain on **CO3D**



Finetune on **Robot Env**

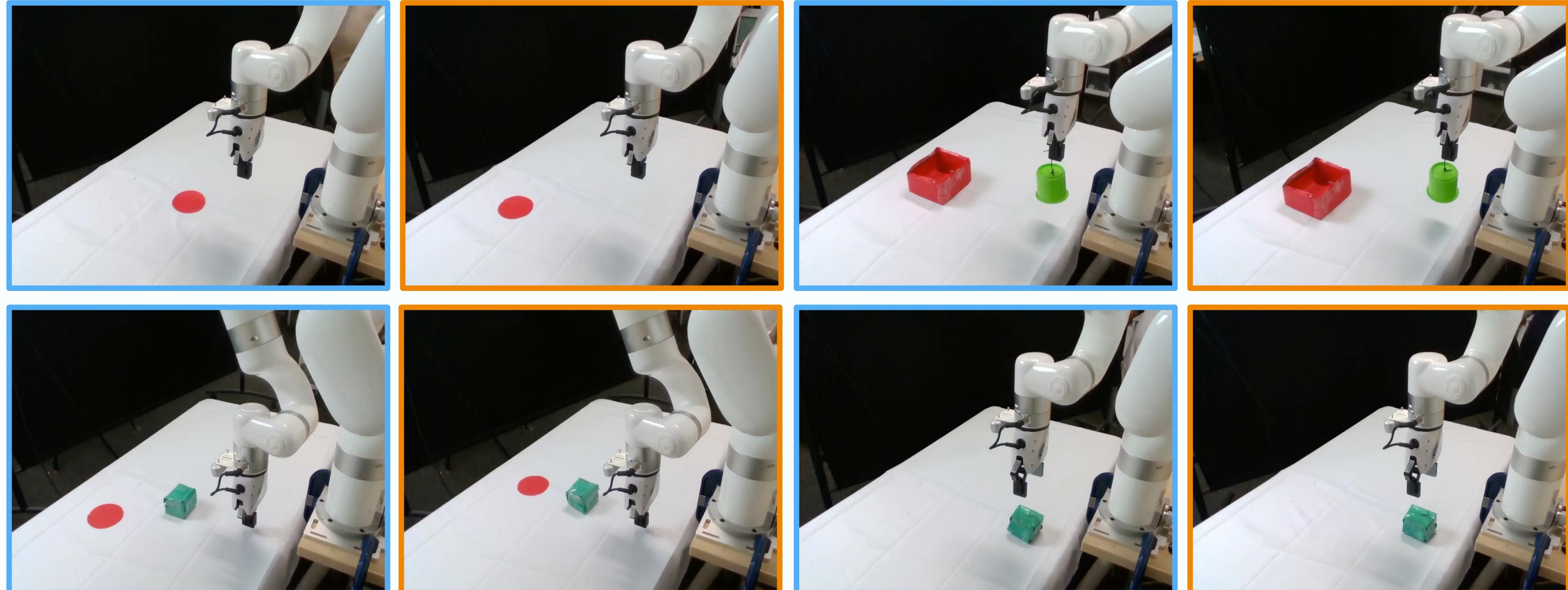


**Sim-to-Real**

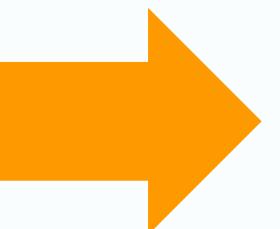
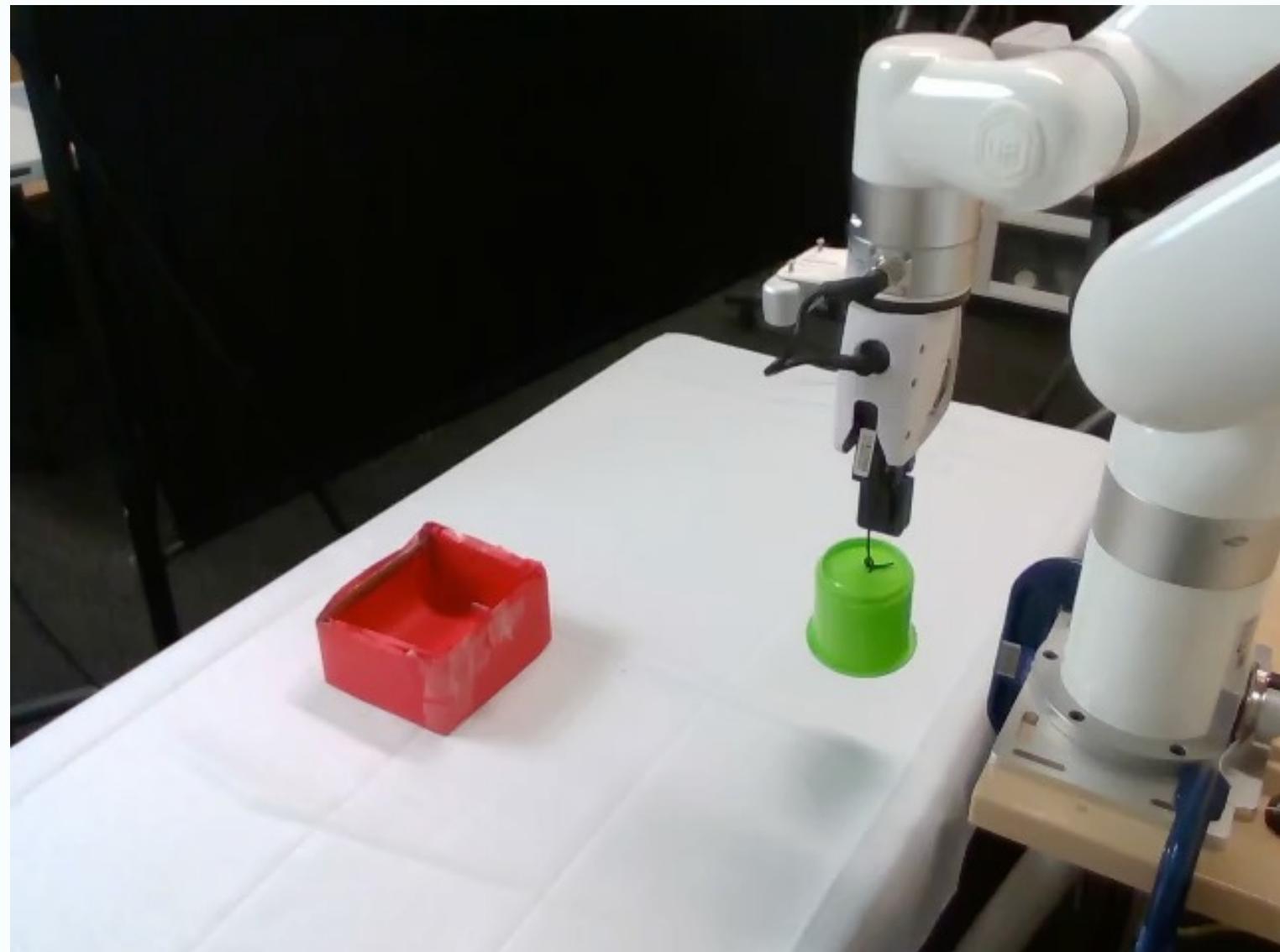


# Sim-to-Real

- Baseline (MoCo)
- RL3D



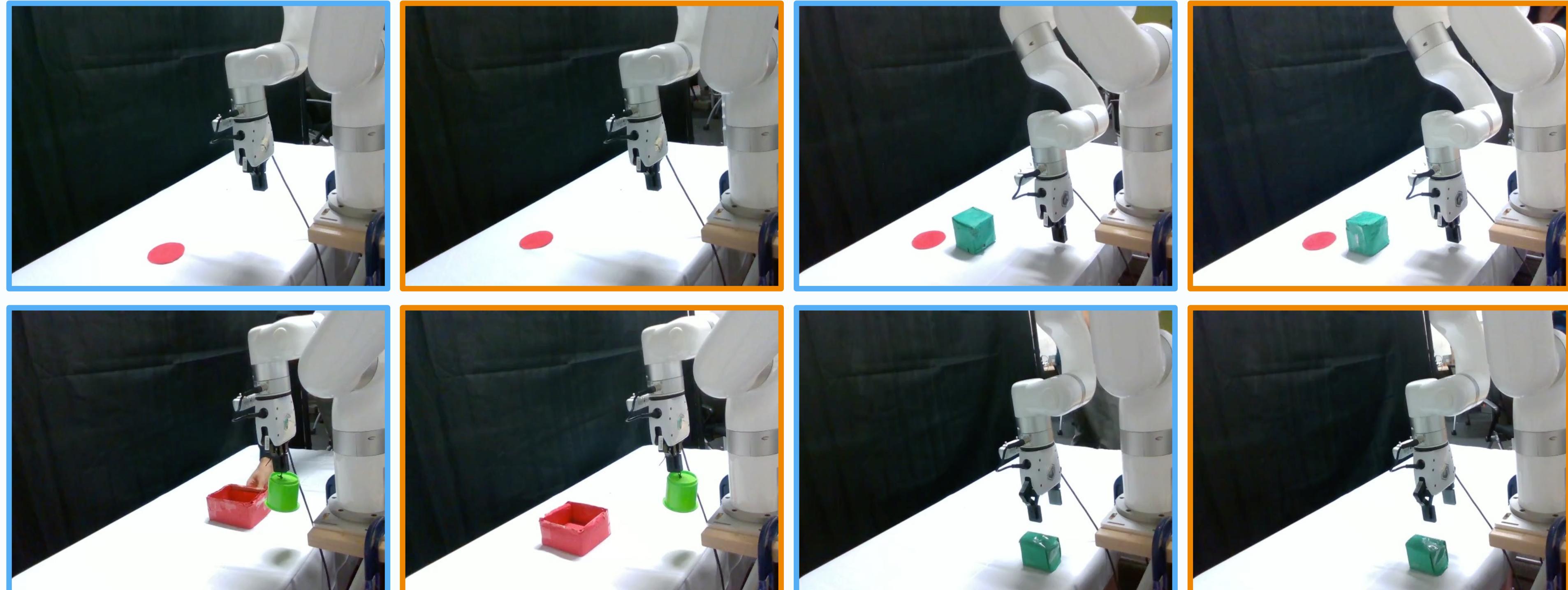
# Generalize to New View



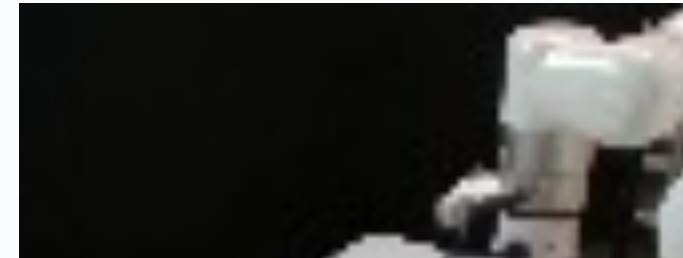
# Generalize to New View

- Baseline (MoCo)
- RL3D

3D representations are more view-invariant!

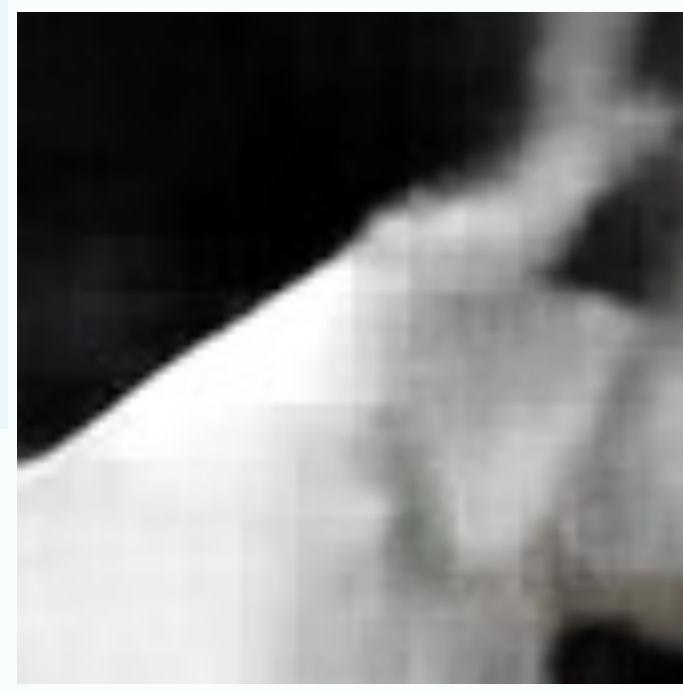
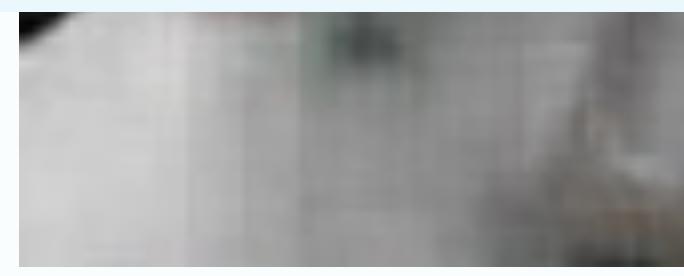


# Real World View Synthesis

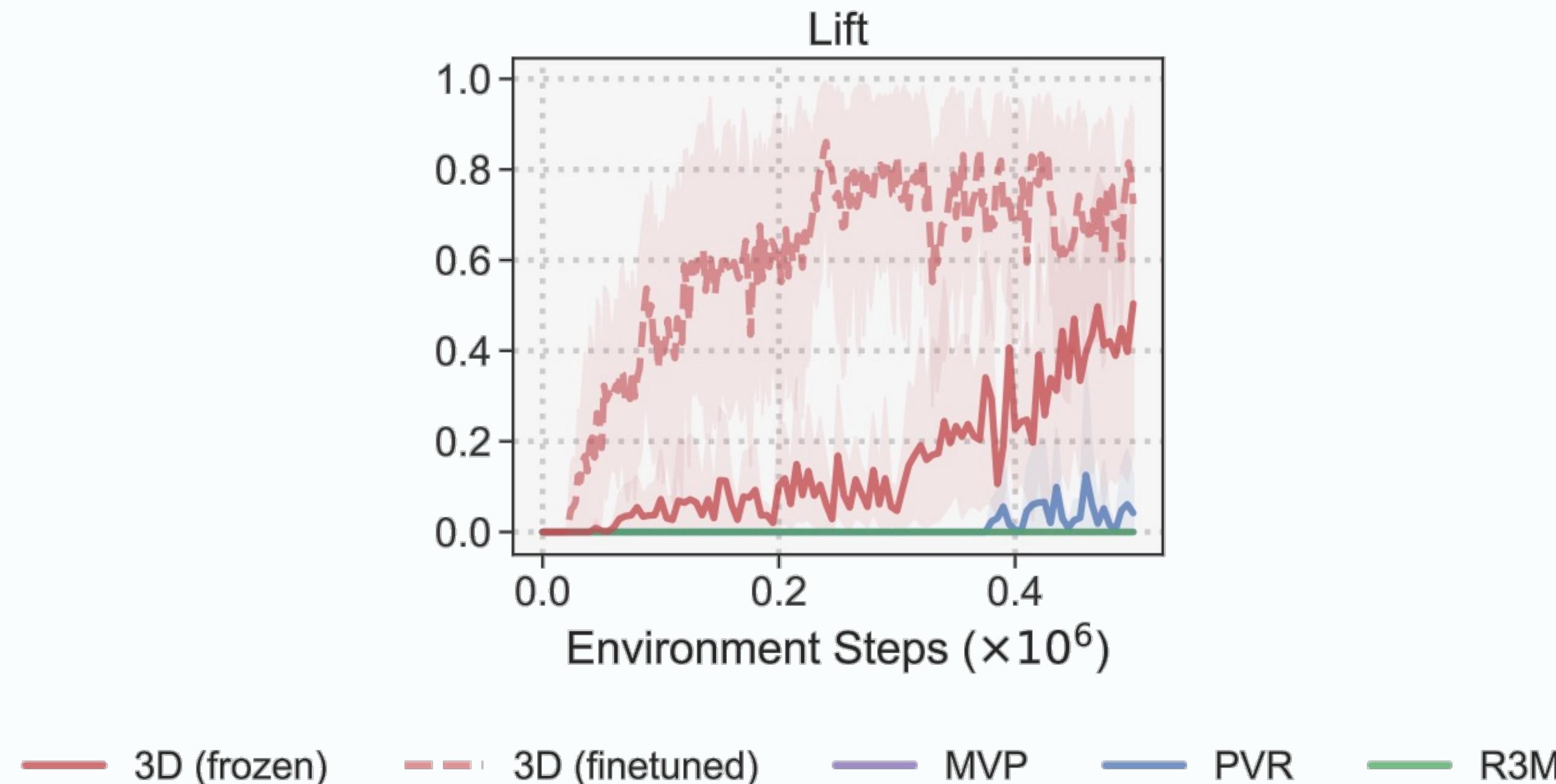


**Reasons for low-quality synthesis:**

- Sim-to-real
- Joint training
- Voxel-based representations

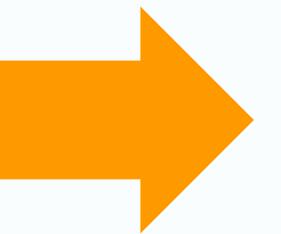


# Finetuning Pre-trained Model is Critical !



# Limitations of RL3D

- Only single-task agent
- Poor view synthesis
- Only easy tasks
- Only simplified scene structure
- Require millions of interactions with env



GNFactor!

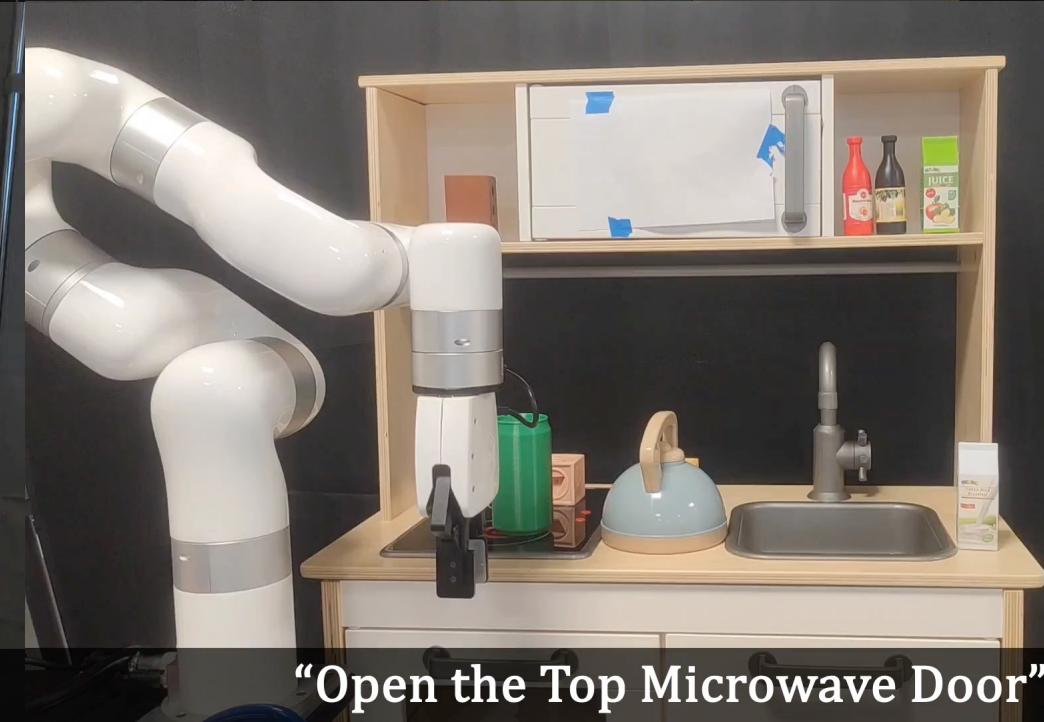
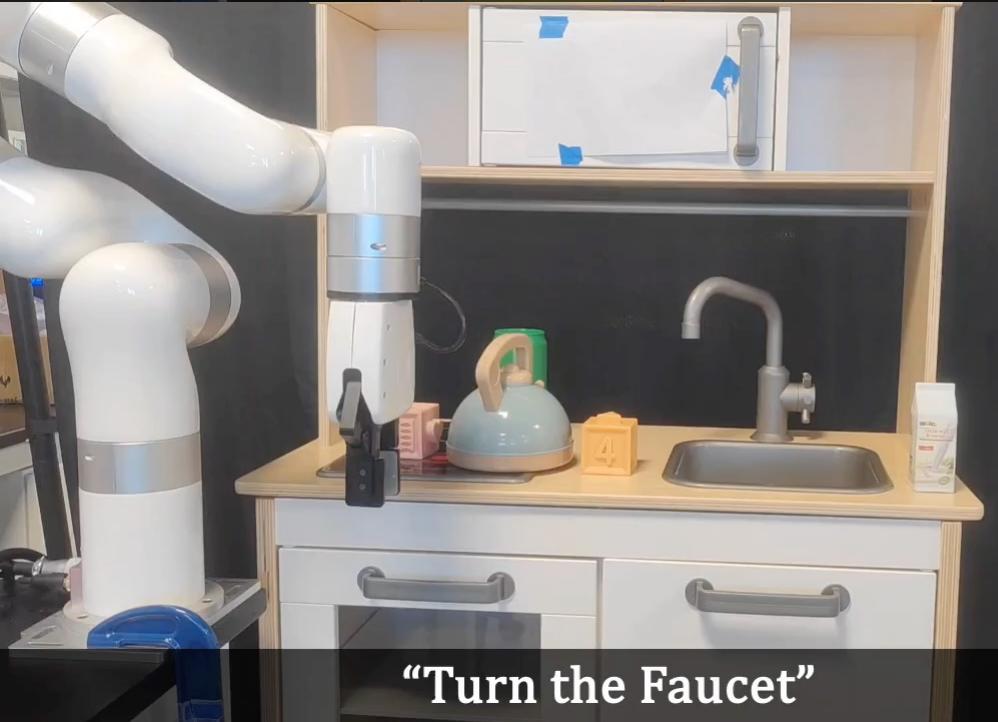
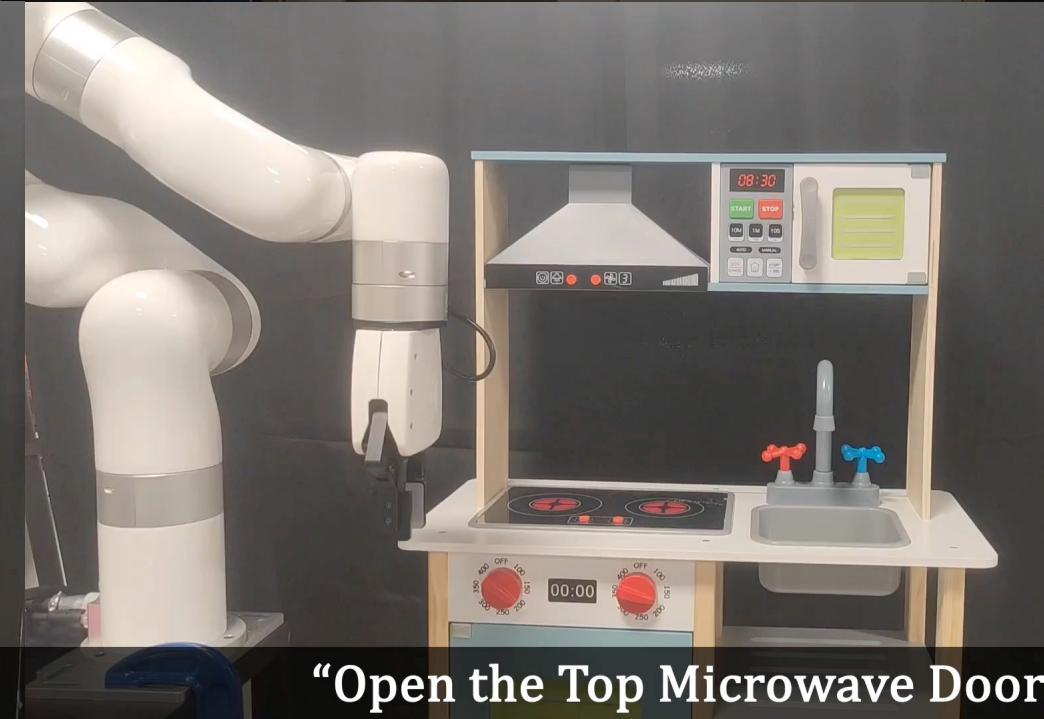
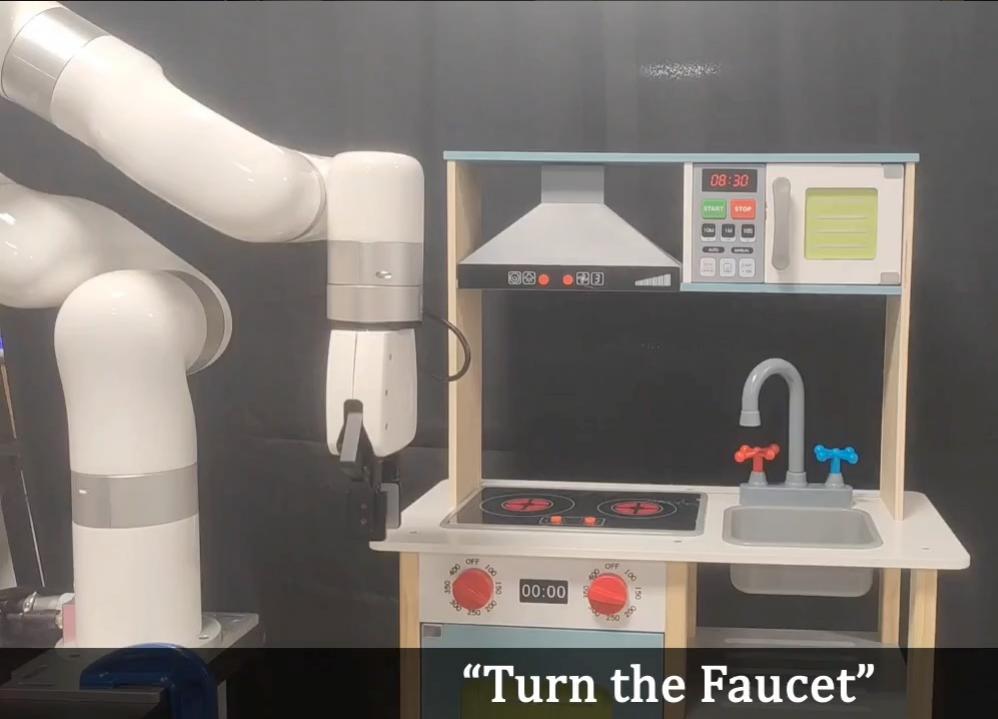
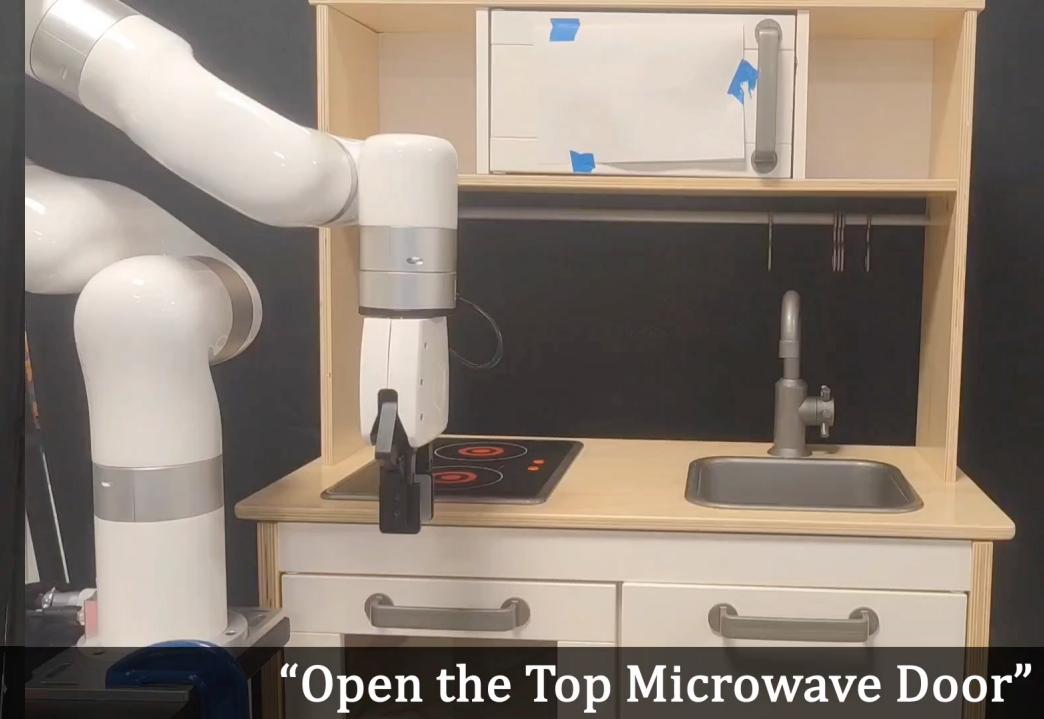
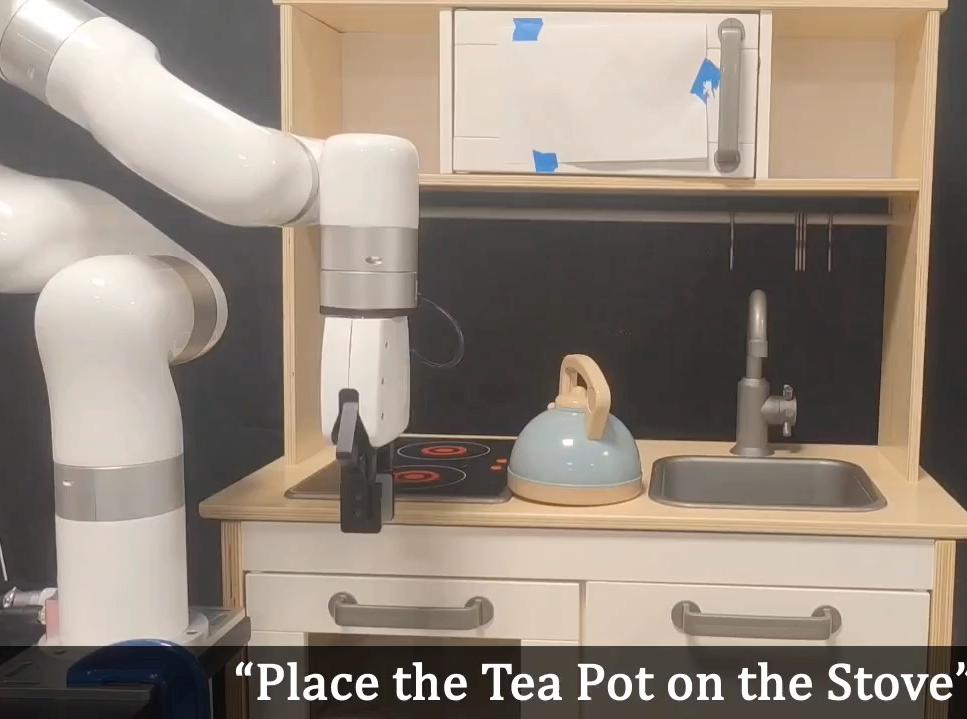
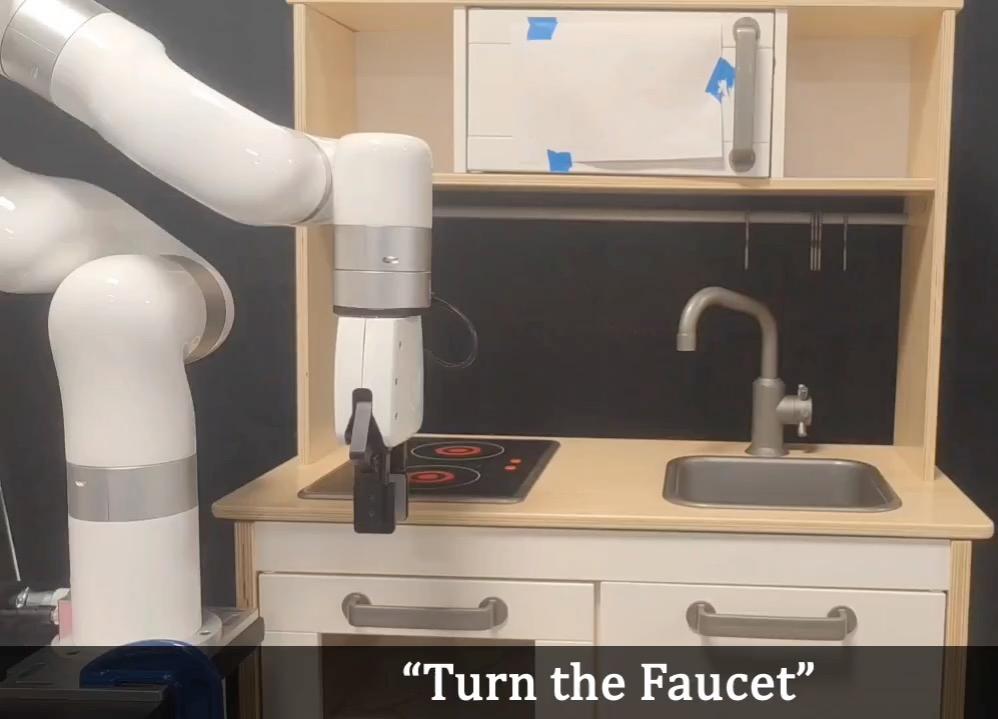
# GNFactor: Multi-Task Real Robot Learning with Generalizable Neural Feature Fields

Yanjie Ze<sup>1\*</sup> Ge Yan<sup>2\*</sup> Yueh-Hua Wu<sup>2\*</sup> Annabella Macaluso<sup>2</sup>  
Yuying Ge<sup>3</sup> Jianglong Ye<sup>2</sup> Nicklas Hansen<sup>2</sup> Li Erran Li<sup>4</sup> Xiaolong Wang<sup>2</sup>

<sup>1</sup>Shanghai Jiao Tong University <sup>2</sup>UC San Diego <sup>3</sup>University of Hong Kong <sup>4</sup>AWS AI, Amazon

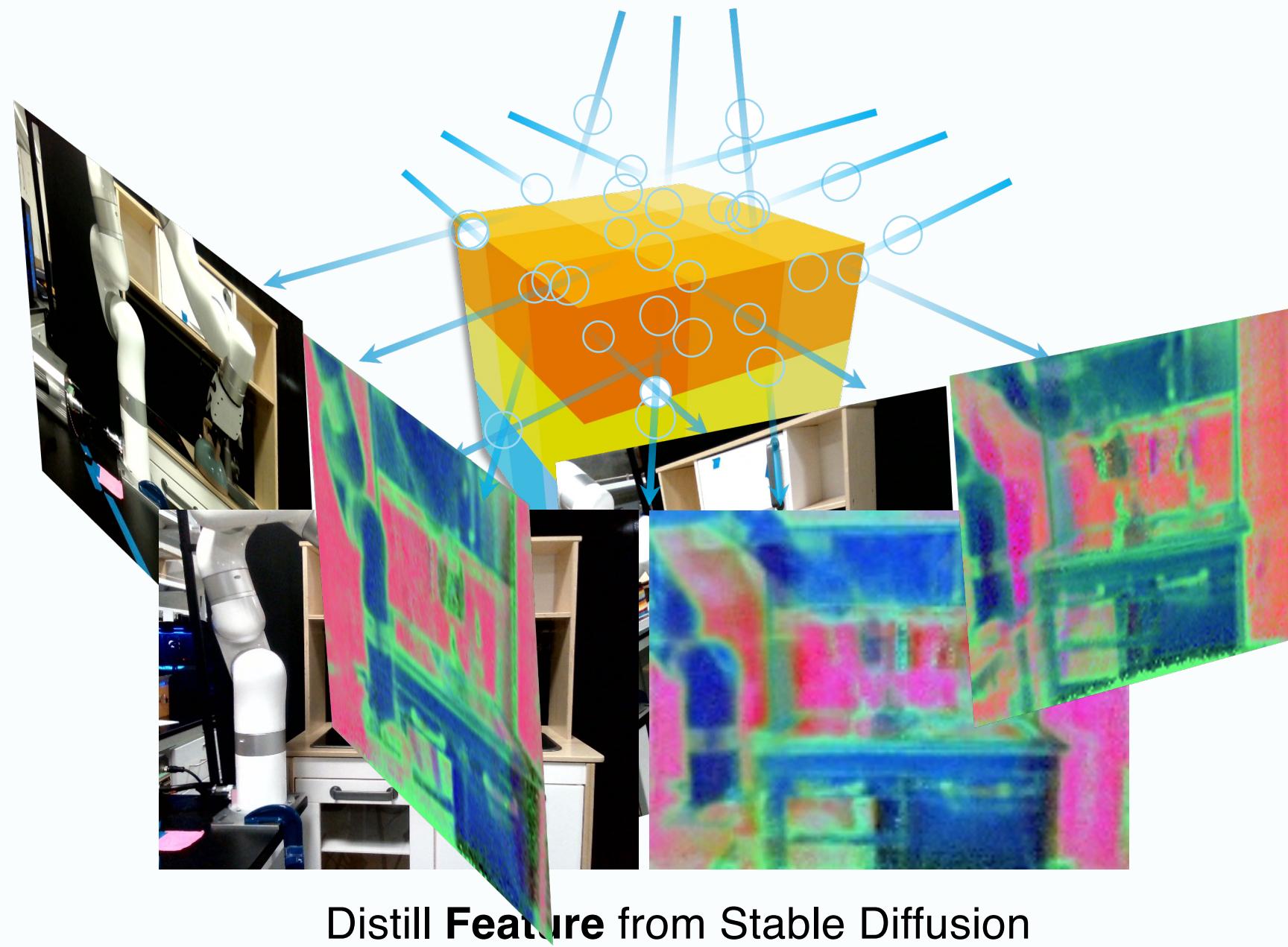
CoRL 2023 Oral





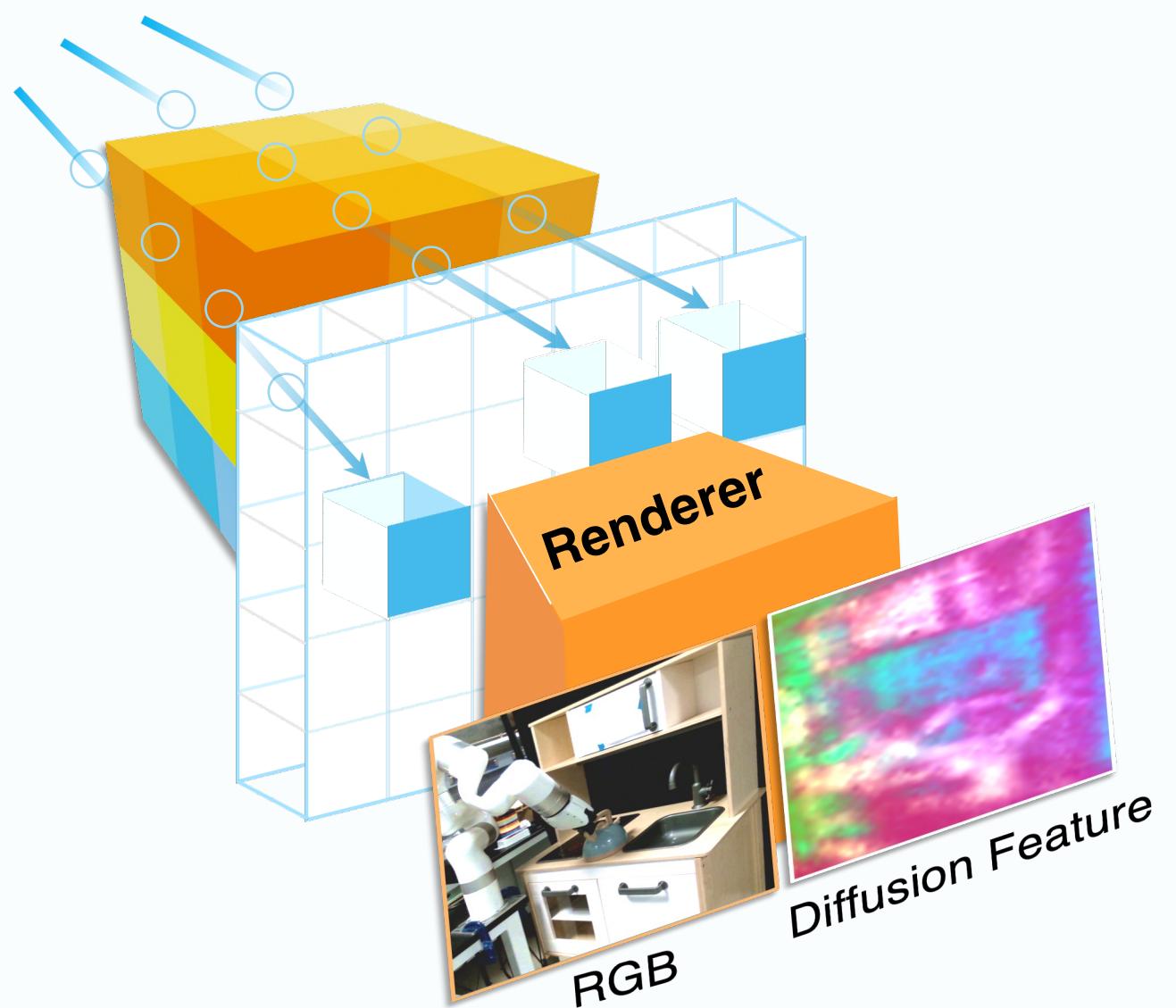
# Overview of GNFactor

## 1. Train a Generalizable Neural Feature Field

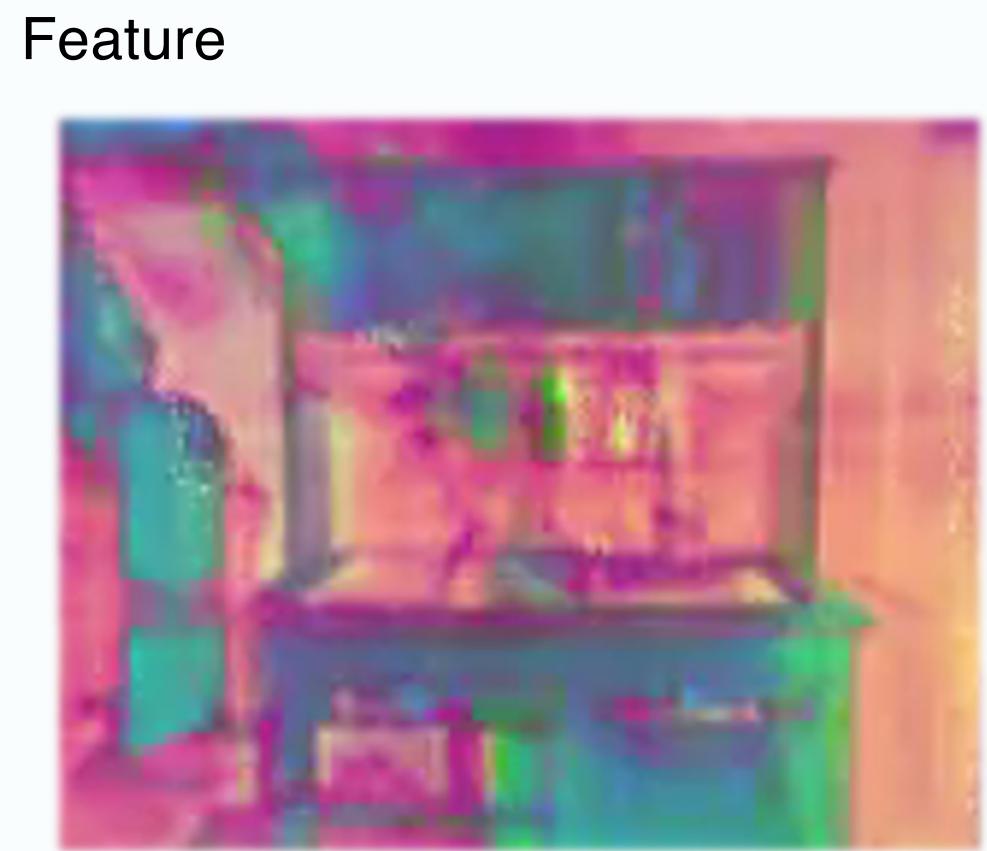


# Overview of GNFactor

## 1. Train a Generalizable Neural Feature Field

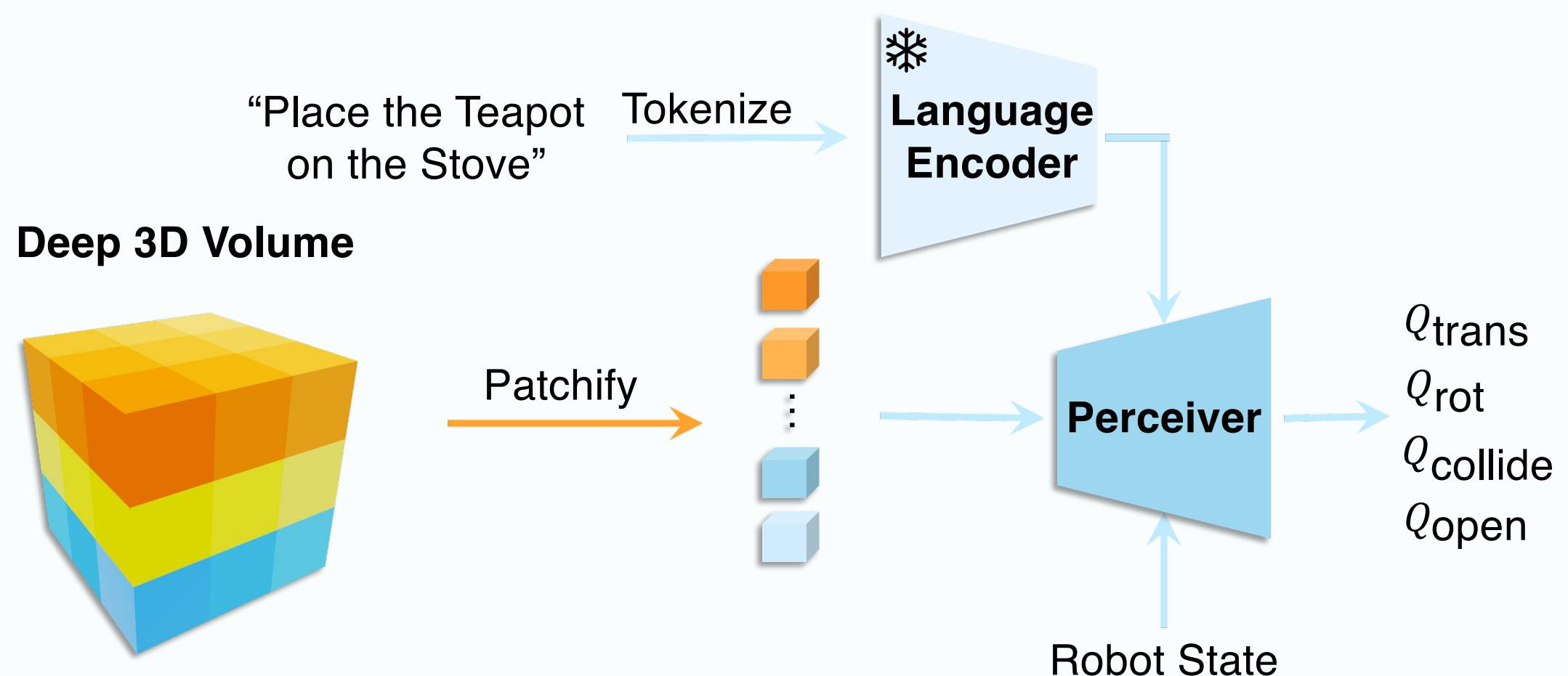


## Novel View Synthesis



# Overview of GNFactor

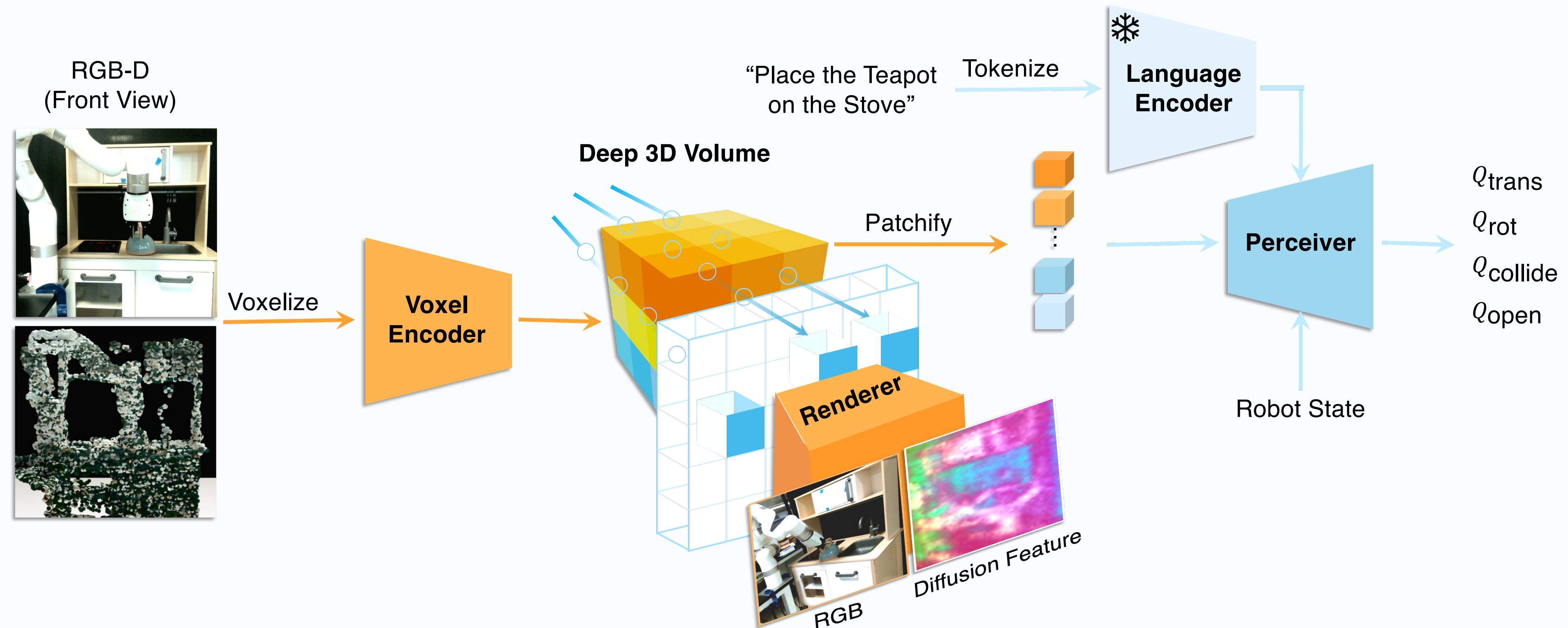
## 2. Train a Multi-Task Policy Jointly



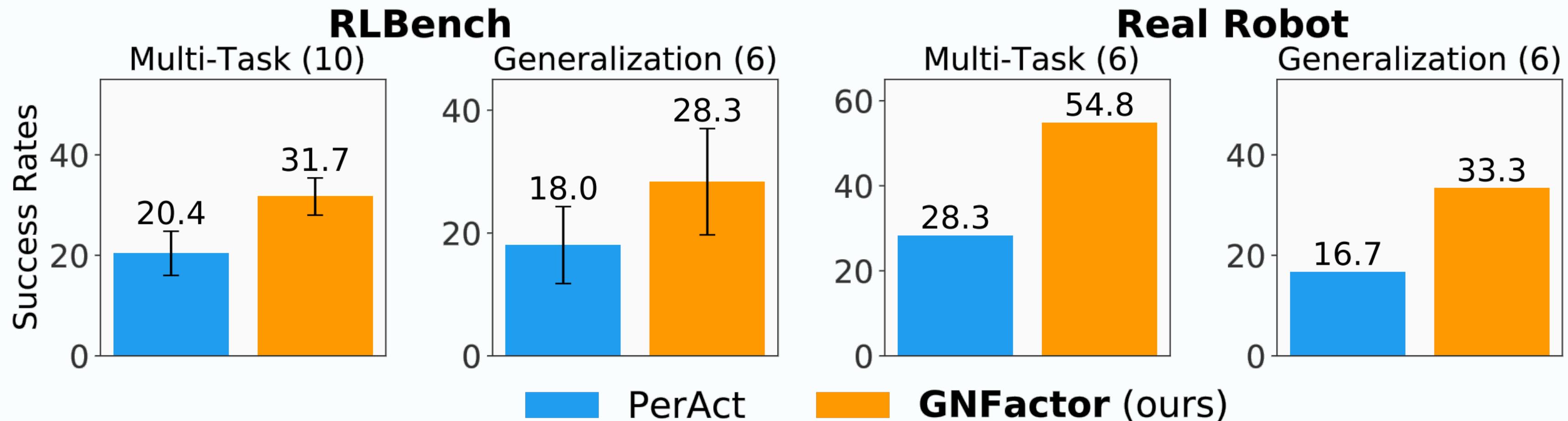
Grad-CAM in 3D



# Overview of GNFactor



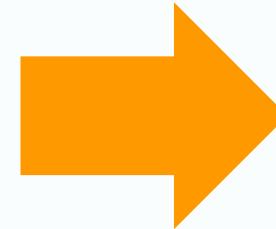
# Strong Results in both Sim and Real



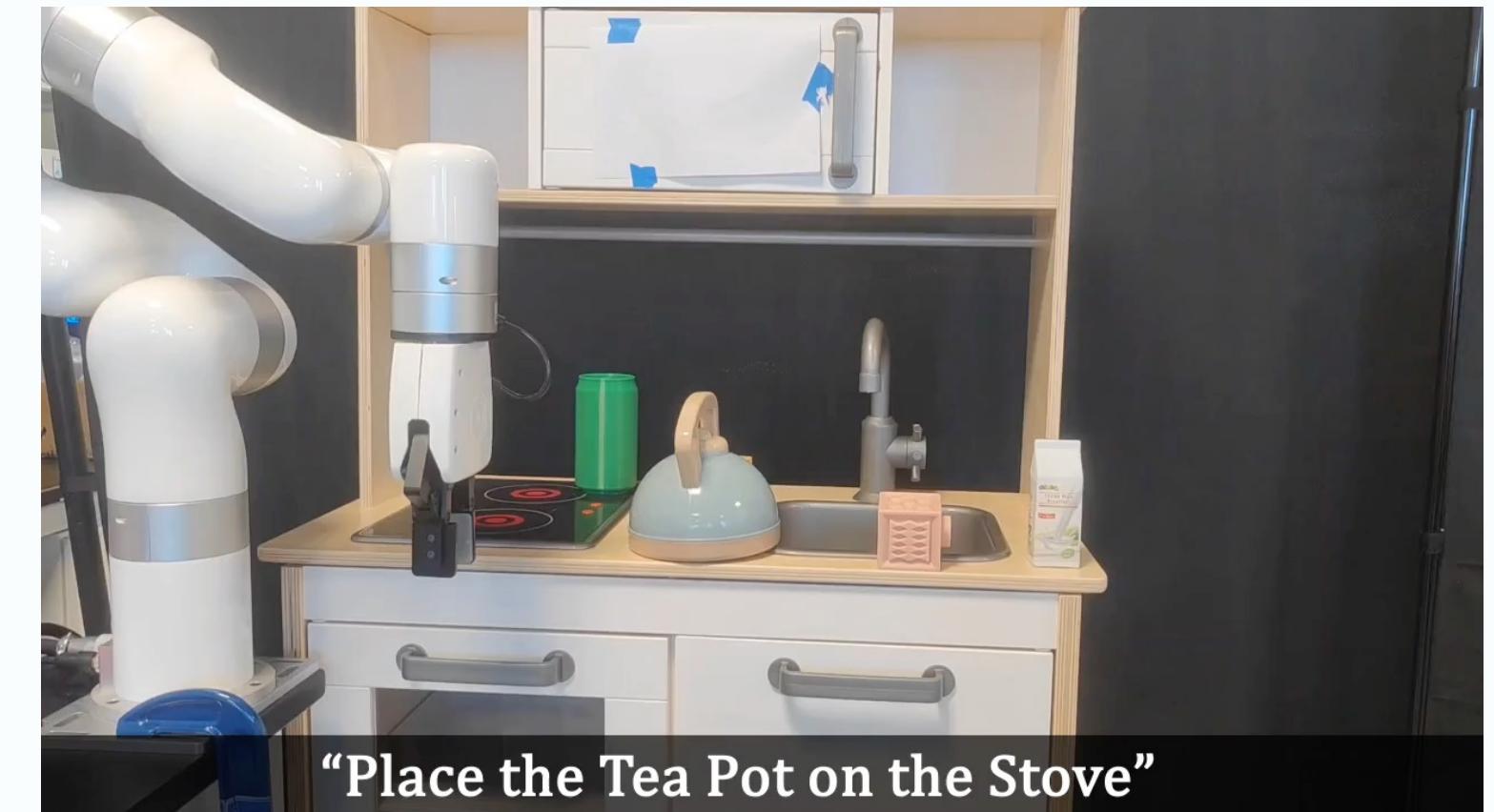
Shridhar et al, Perceiver-Actor: A Multi-Task Transformer for Robotic Manipulation, CoRL 2022.

# Generalize to Cluttered Scenes

Training Scene



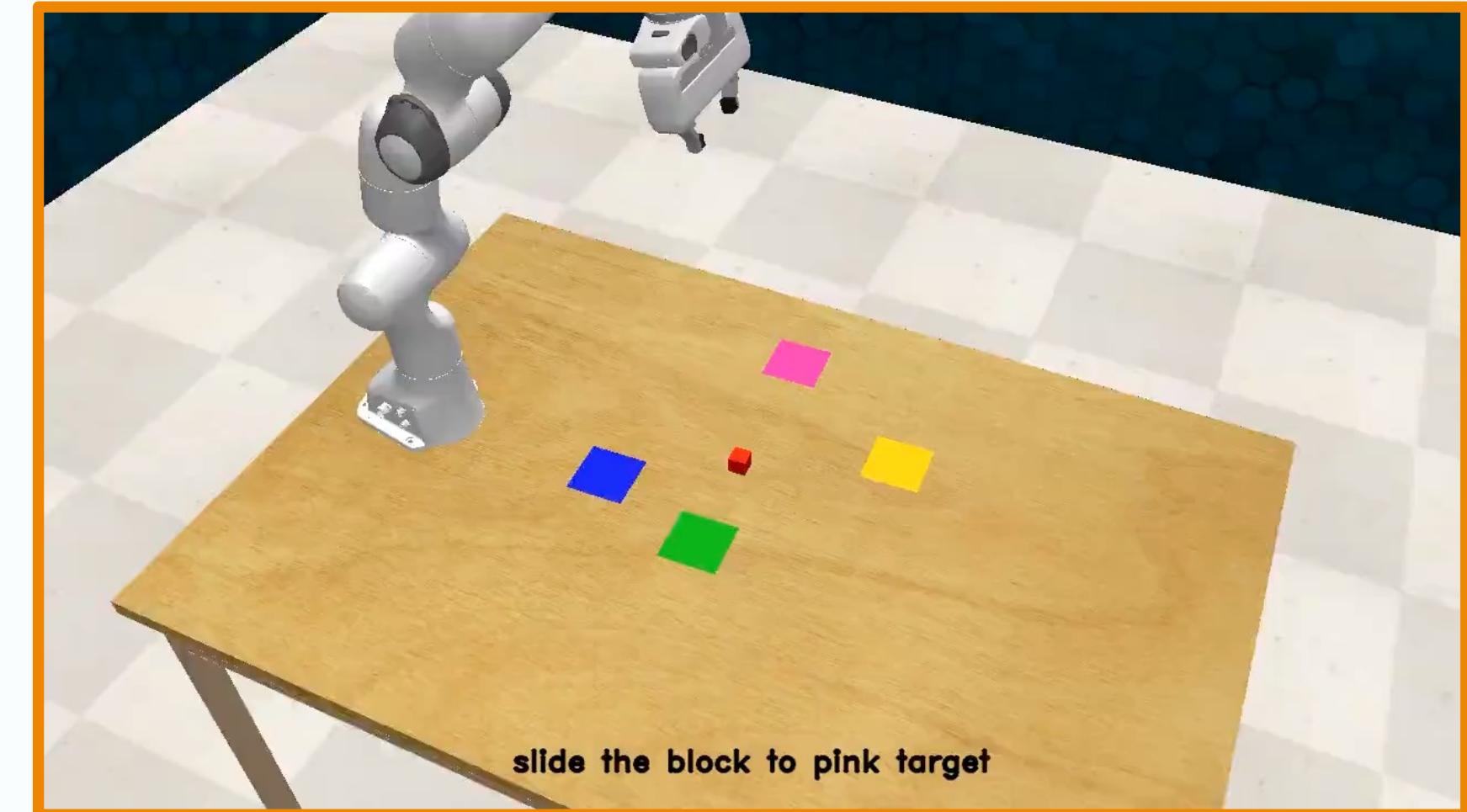
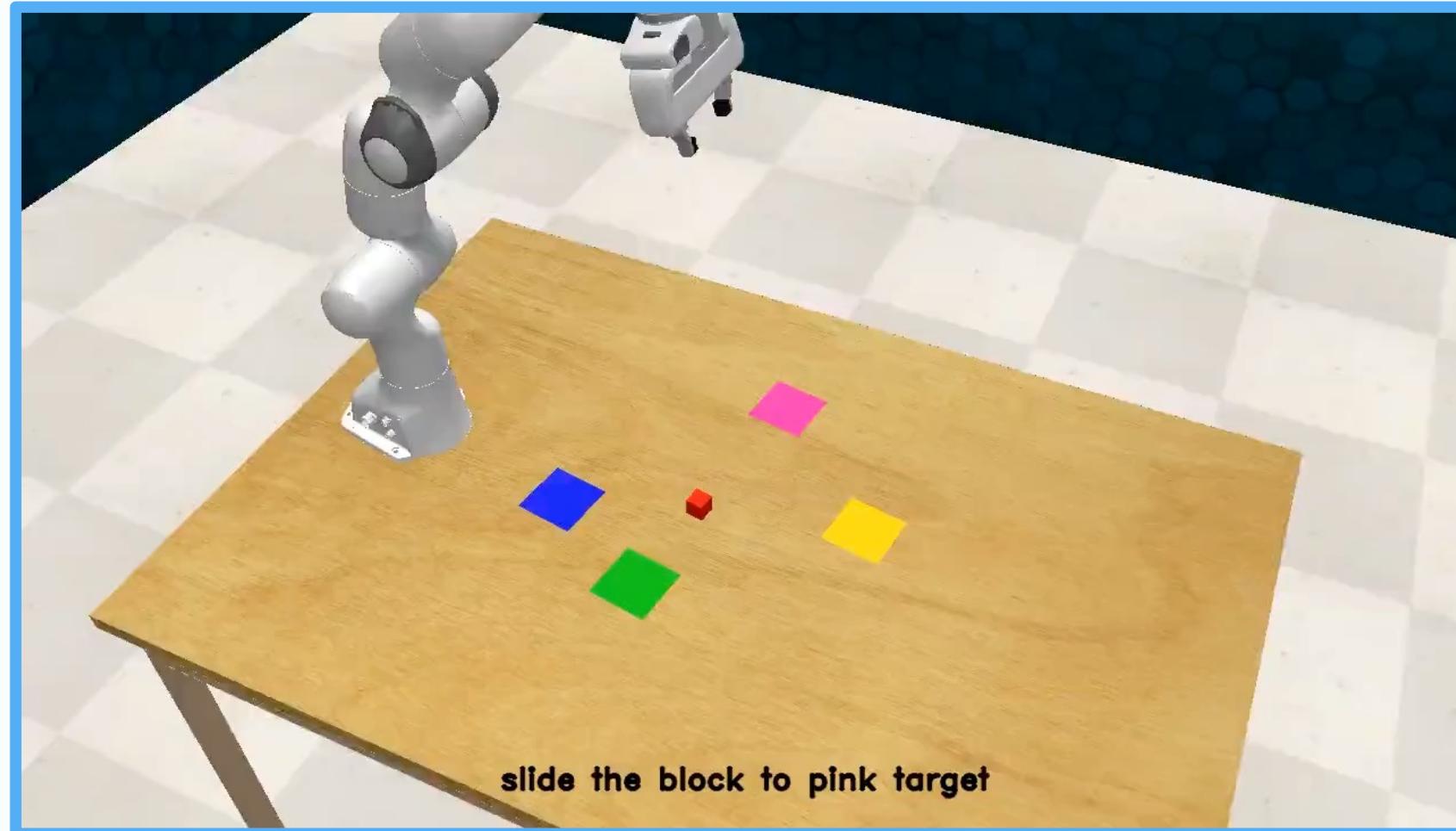
Cluttered Scene



# Generalize to Unseen Scenarios

- PerAct
- GNFactor

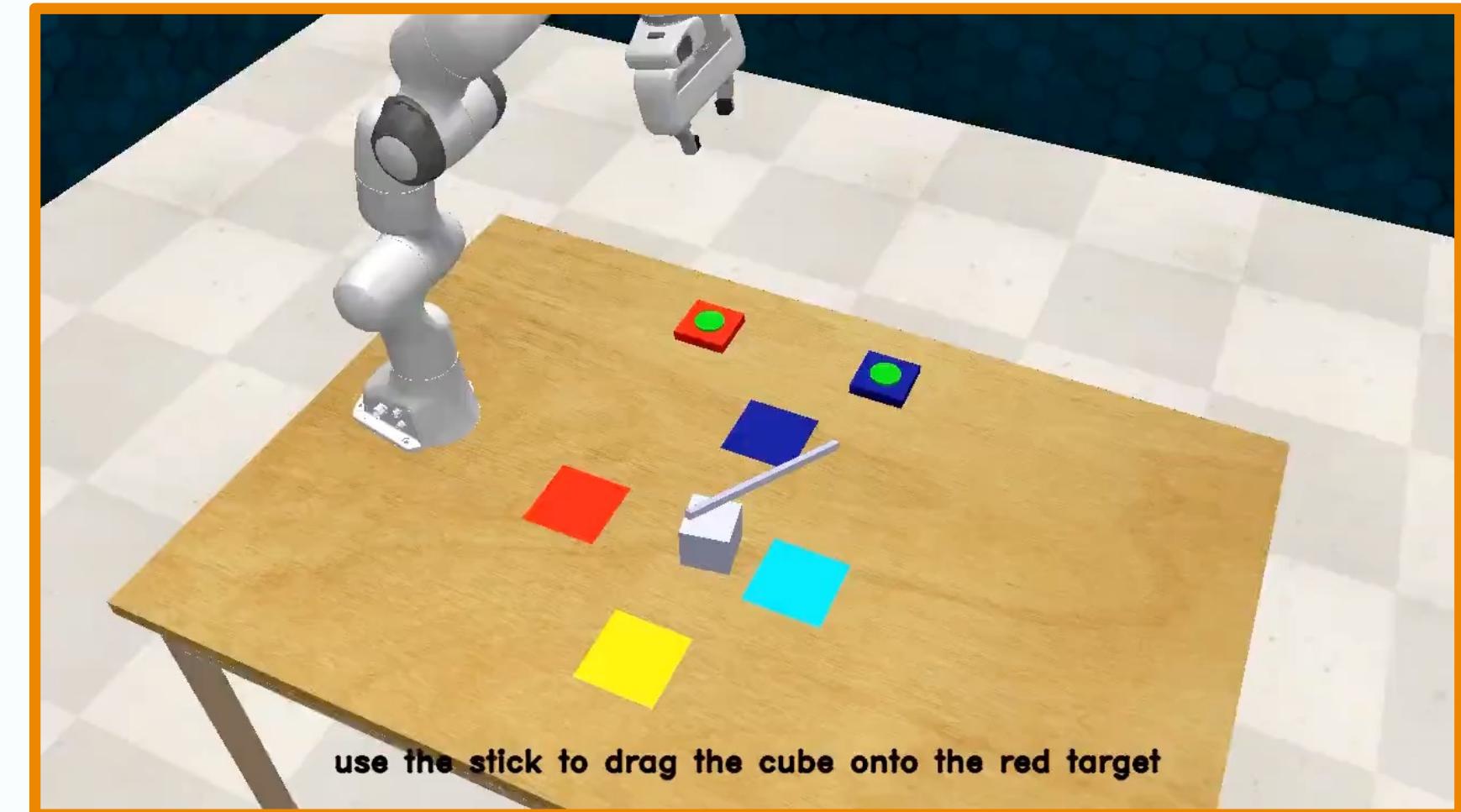
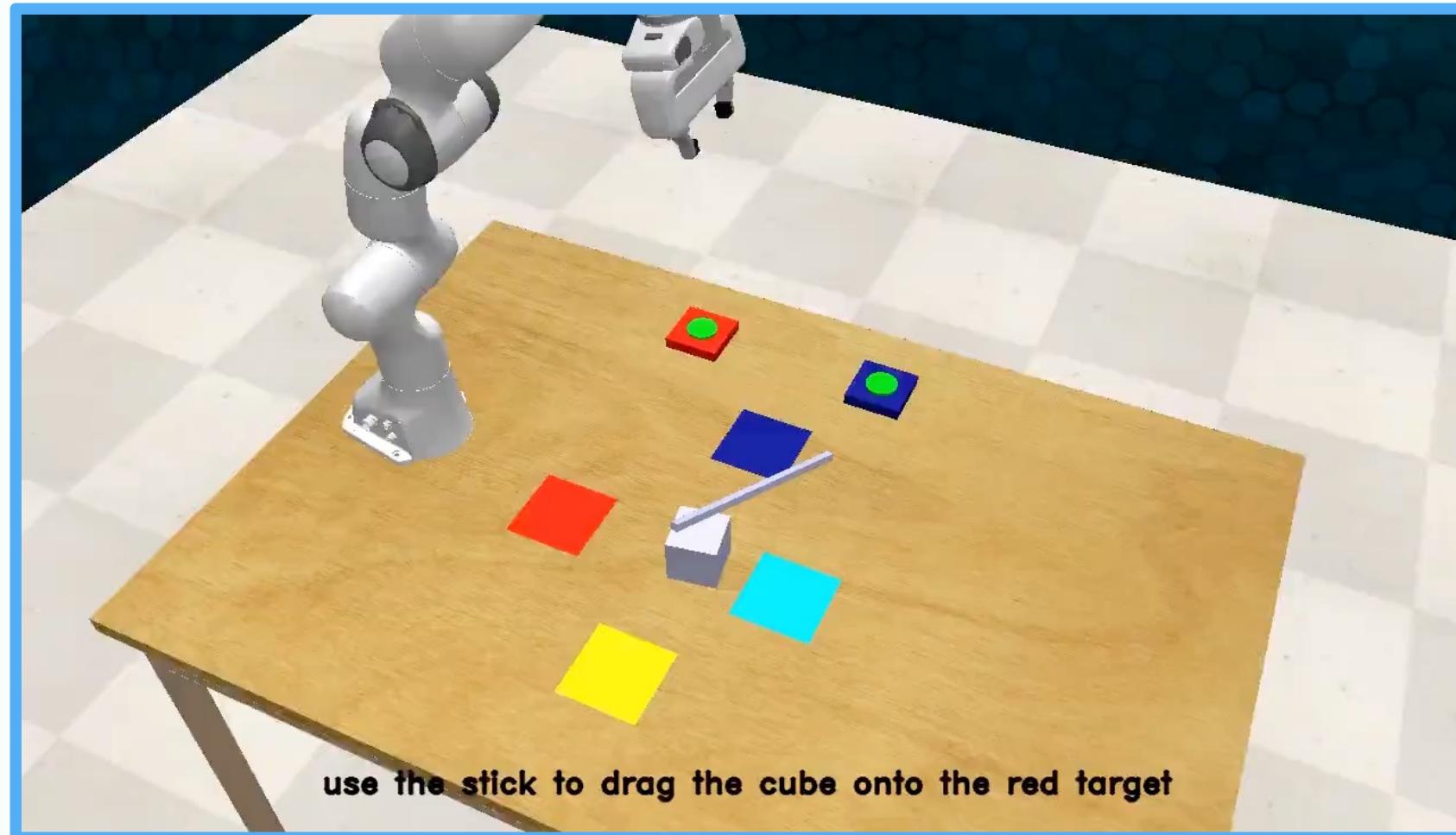
Slide a smaller block to the target



# Generalize to Unseen Scenarios

- PerAct
- GNFactor

Drag the stick with distractors



# Key Factors

Ablation	Success Rate (%)
GNFactor	<b>36.8</b>

# Key Factors: Generalizable NeRF

Ablation	Success Rate (%)
GNFactor	<b>36.8</b>
w/o. GNF objective	24.2
w/o. RGB objective	27.2

# Key Factors: Distilling Foundation Models

Ablation	Success Rate (%)
GNFactor	<b>36.8</b>
w/o. GNF objective	24.2
w/o. RGB objective	27.2
w/o. Diffusion	30.0
Diffusion → DINO	30.4
Diffusion → CLIP	32.0

# Key Factors: Engineering

Ablation	Success Rate (%)
GNFactor	<b>36.8</b>
w/o. GNF objective	24.2
w/o. RGB objective	27.2
w/o. Diffusion	30.0
Diffusion → DINO	30.4
Diffusion → CLIP	32.0
w/o. DGS	29.2
w/o. skip connection	27.6

# Key Factors: Parameter Tuning

Ablation	Success Rate (%)
GNFactor	<b>36.8</b>
w/o. GNF objective	24.2
w/o. RGB objective	27.2
w/o. Diffusion	30.0
Diffusion → DINO	30.4
Diffusion → CLIP	32.0
w/o. DGS	29.2
w/o. skip connection	27.6
$k = 19 \rightarrow 9$	33.2
$\lambda_{\text{feat}} = 0.01 \rightarrow 1.0$	35.2
$\lambda_{\text{recon}} = 0.01 \rightarrow 1.0$	35.2

# Visual Representations for Generalizable Robotic Manipulation

1

## Geometric Prior

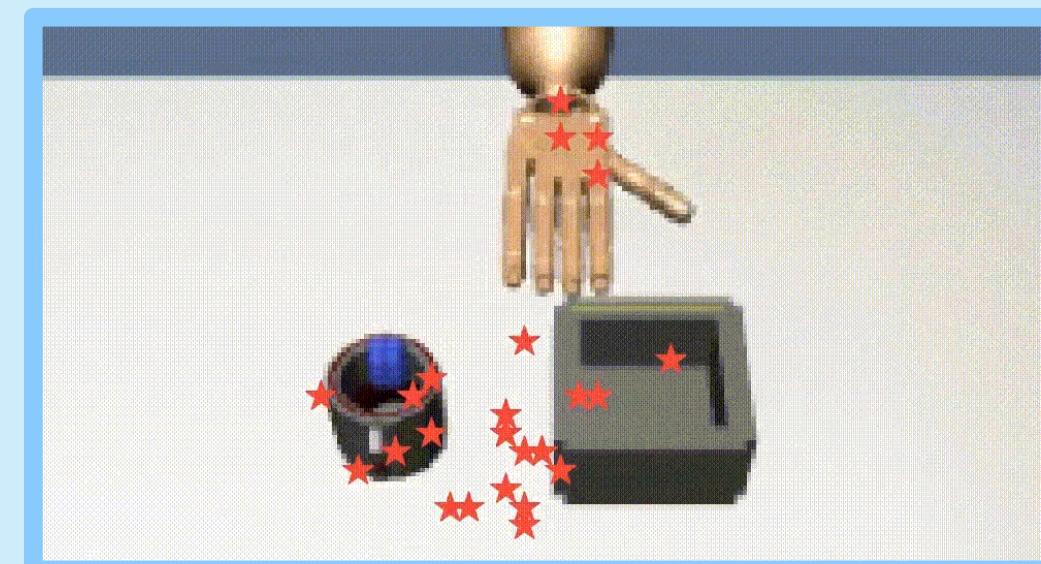


[2] Ze et al., “Visual Reinforcement Learning with Self-Supervised 3D Representations”, RA-L 2023 & IROS 2023.

[3] Ze et al., “GNFactor: Multi-Task Real Robot Learning with Generalizable Neural Feature Fields”, CoRL 2023 Oral.

2

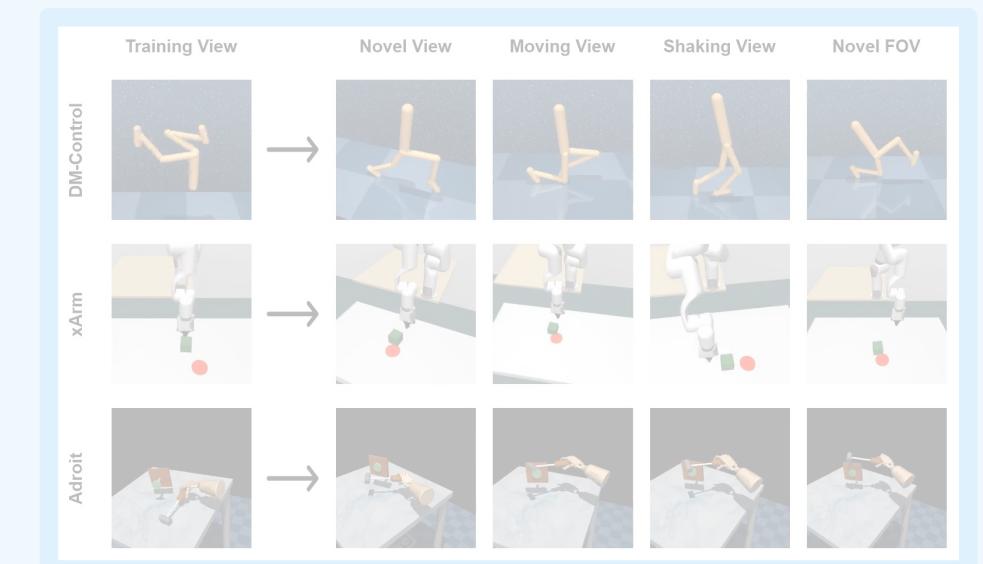
## Human Prior



[4] Ze et al., “H-InDex: Visual Reinforcement Learning with Hand-Informed Representations for Dexterous Manipulation”, NeurIPS 2023.

3

## Dynamics Prior



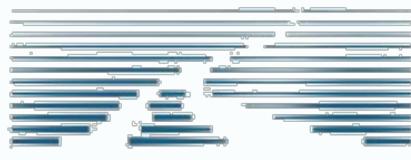
[5] Yang\*, Ze\* et al., “MoVie: Visual Model-Based Policy Adaptation for View Generalization”, NeurIPS 2023.

# H-InDex: Visual Reinforcement Learning with Hand-Informed Representations for Dexterous Manipulation

Yanjie Ze<sup>12</sup> Yuyao Liu<sup>3\*</sup> Ruizhe Shi<sup>3\*</sup> Jiaxin Qin<sup>4</sup>  
Zhecheng Yuan<sup>31</sup> Jiashun Wang<sup>5</sup> Huazhe Xu<sup>316</sup>

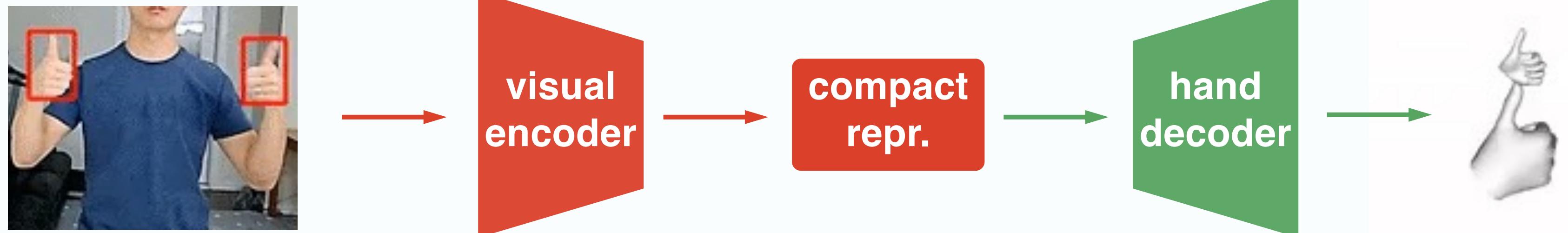
<sup>1</sup>Shanghai Qi Zhi Institute <sup>2</sup>Shanghai Jiao Tong University <sup>3</sup>Tsinghua University, IIIS  
<sup>4</sup>Renmin University of China <sup>5</sup>Carnegie Mellon University <sup>6</sup>Shanghai AI Lab

NeurIPS 2023

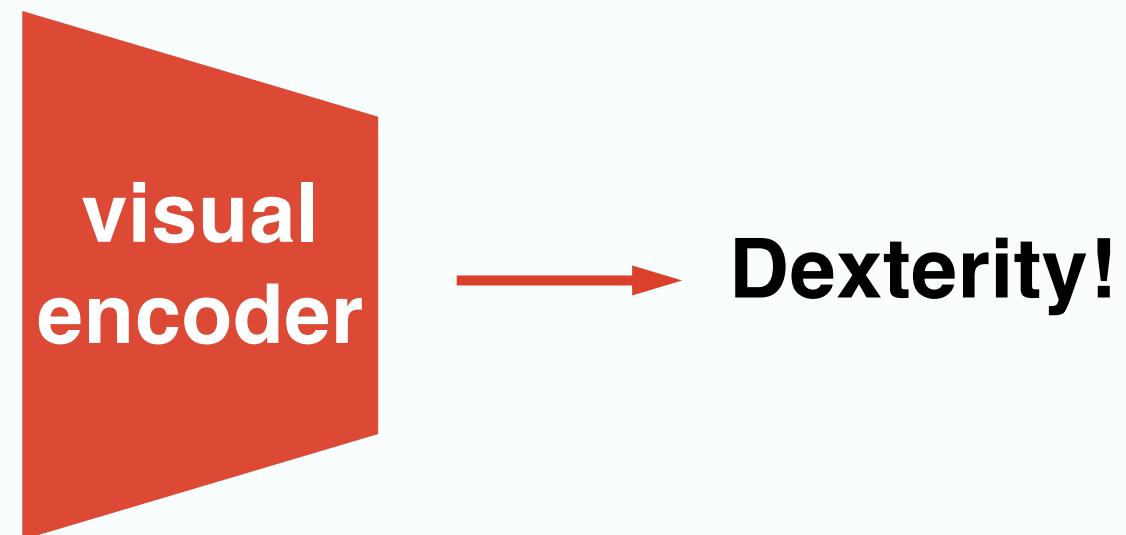


***Can robotic hands leverage representations  
learned from human hands?***

We **directly** adopt the pre-trained visual representation from 3D human hand pose estimation.

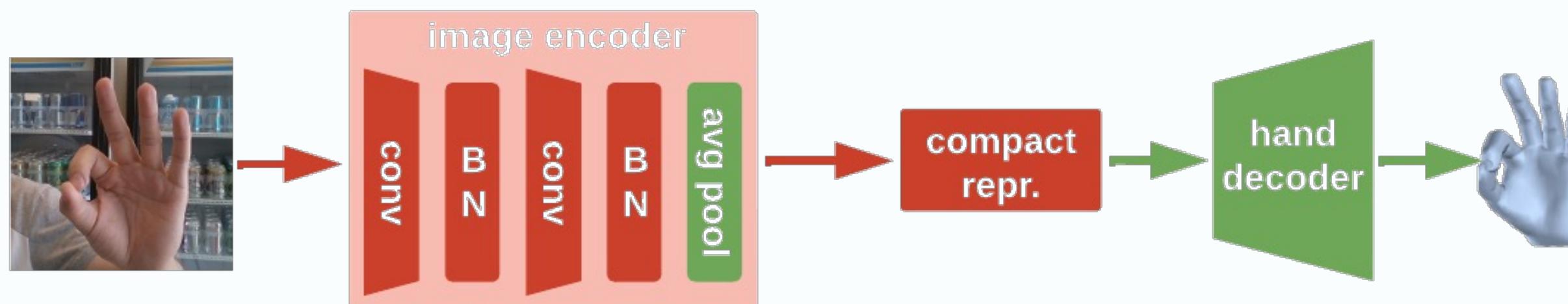


We **directly** adopt the pre-trained visual representation  
from 3D human hand pose estimation.



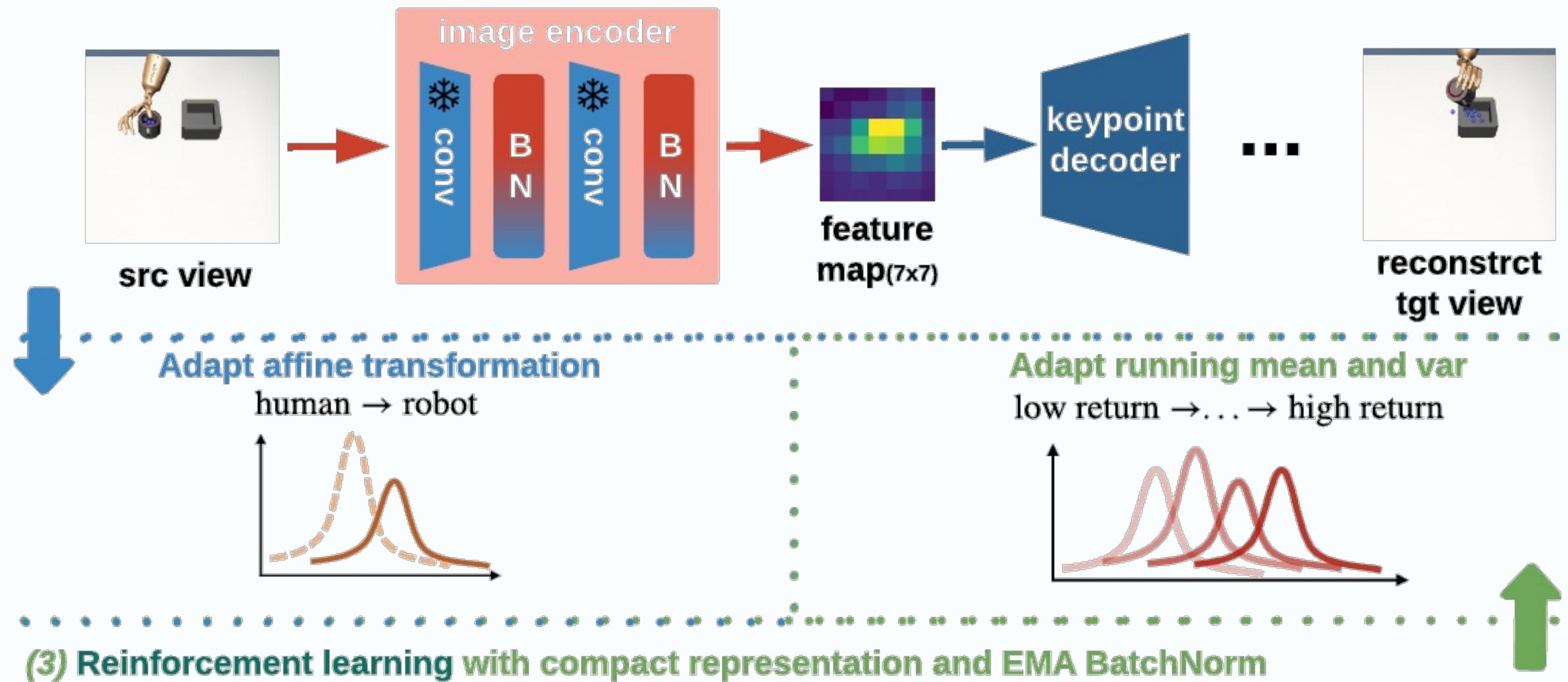
# Overview of H-InDex

## (1) Pre-train representation with 3D human hand pose estimation

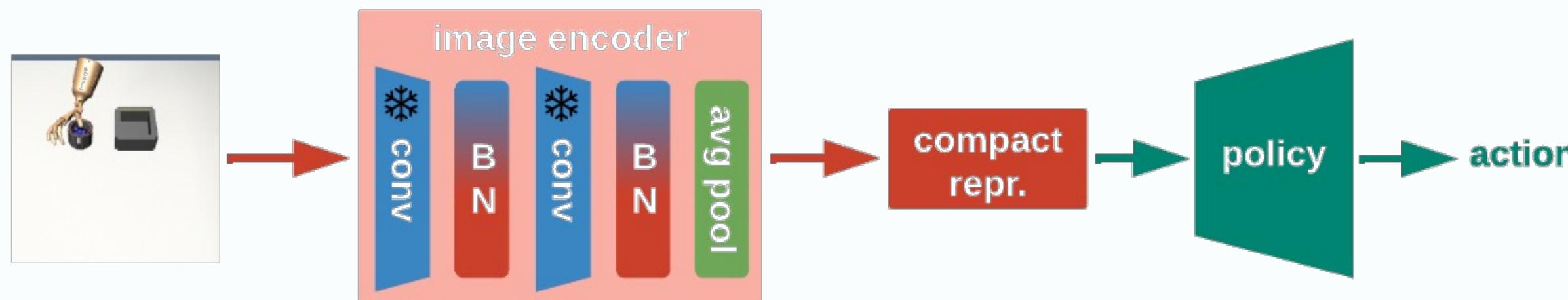


# Overview of H-InDex

## (2) Offline adapt representation with self-supervised keypoint detection

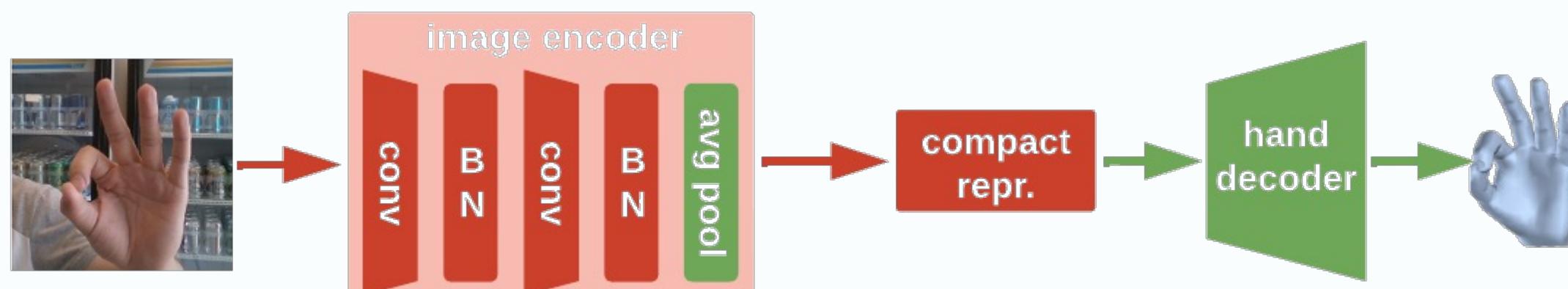


## (3) Reinforcement learning with compact representation and EMA BatchNorm

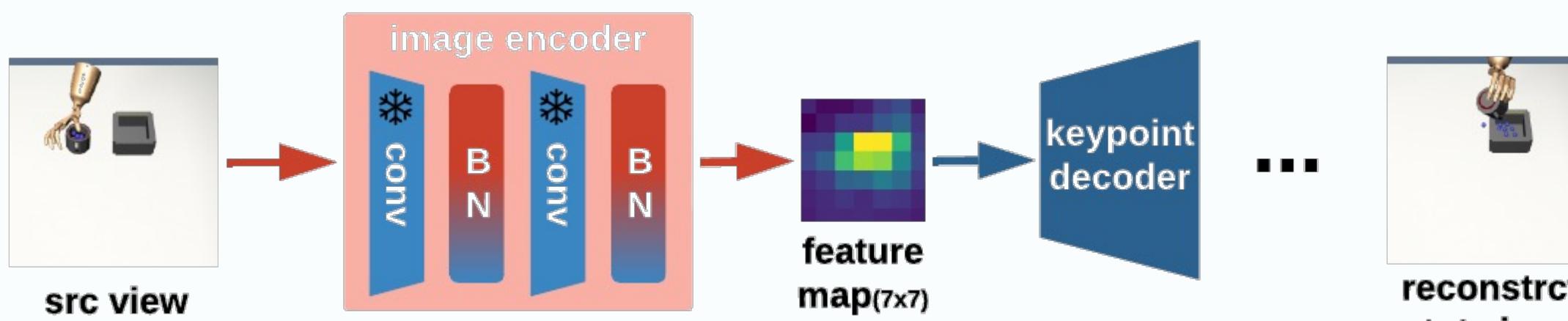


# Overview of H-InDex

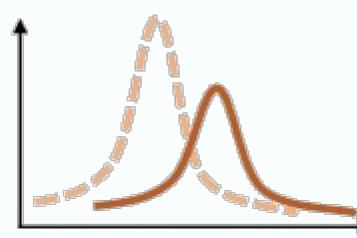
## (1) Pre-train representation with 3D human hand pose estimation



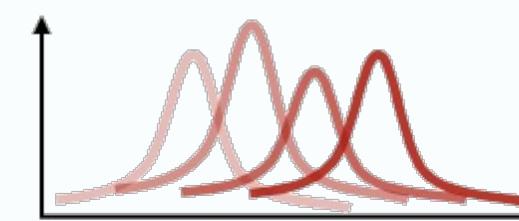
## (2) Offline adapt representation with self-supervised keypoint detection



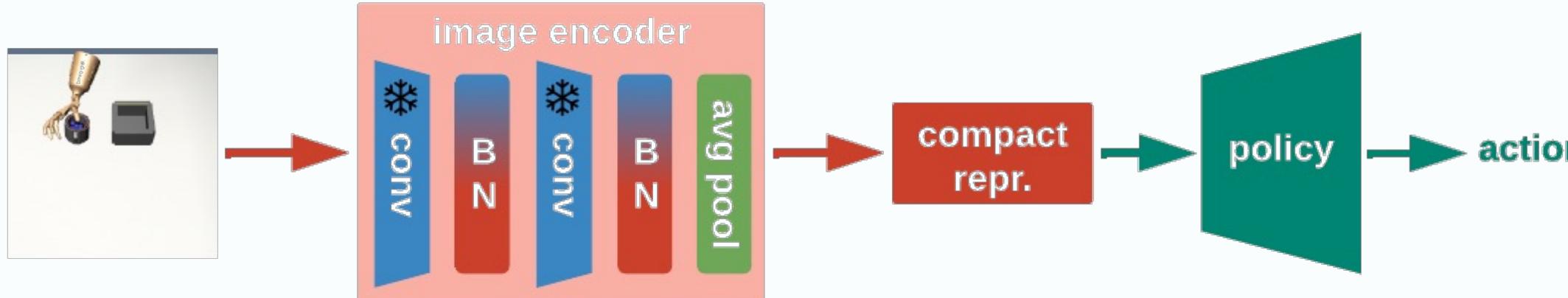
### Adapt affine transformation human → robot



### Adapt running mean and var low return → ... → high return

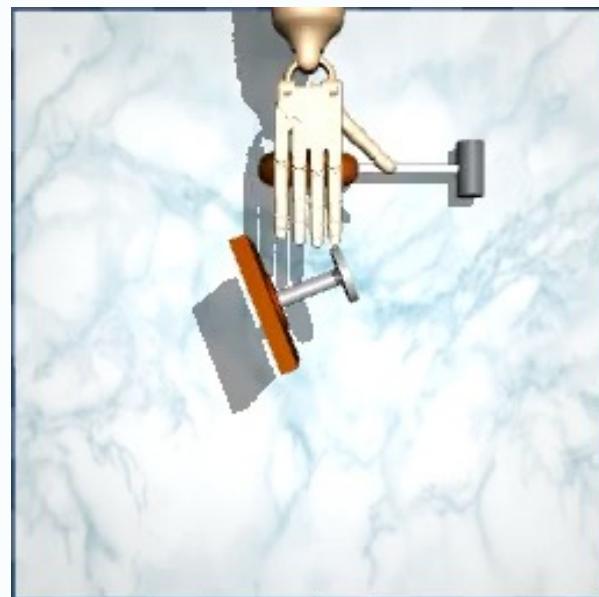


## (3) Reinforcement learning with compact representation and EMA BatchNorm

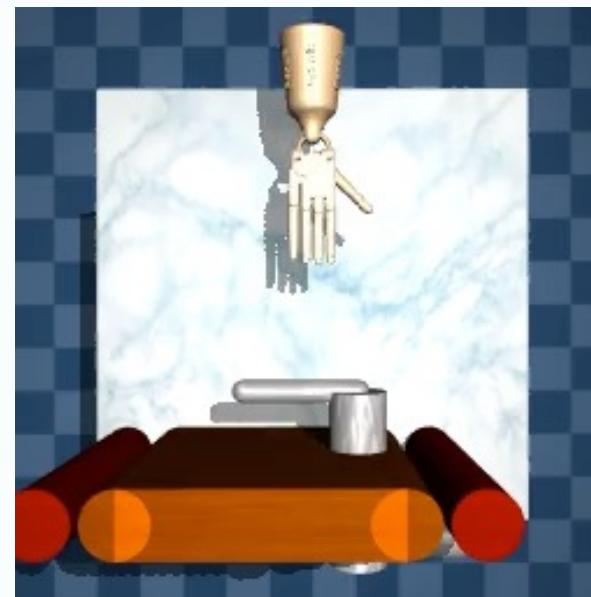


# Strong Results in **12** Dexterous Manipulation Tasks

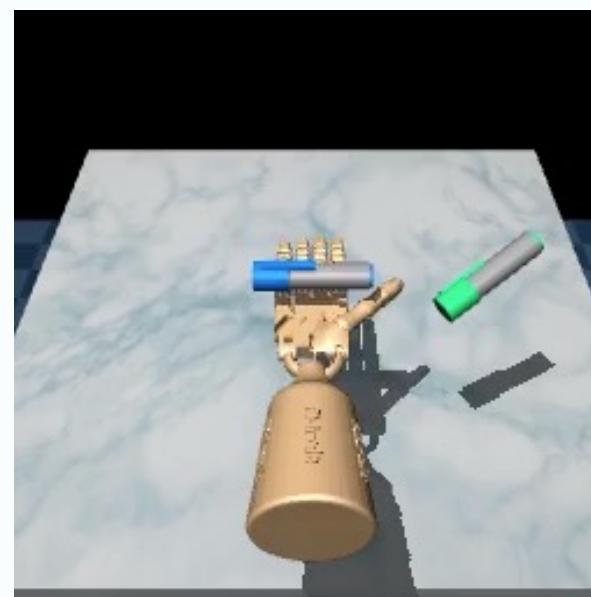
**Hammer**



**Door**



**Pen**



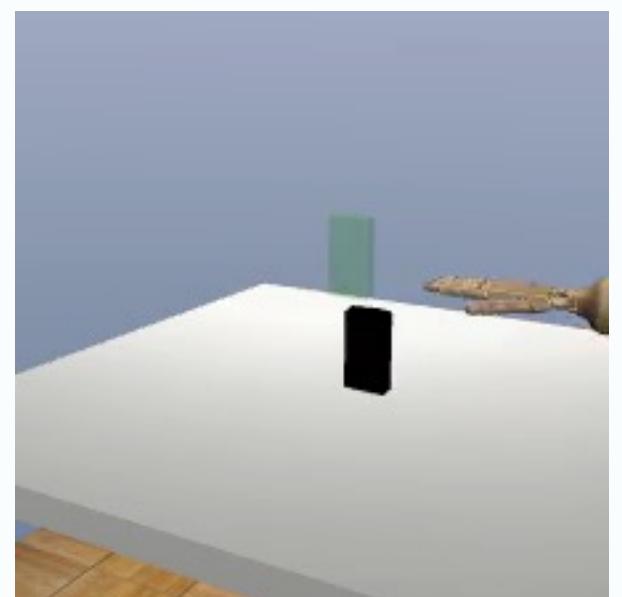
**Pour**



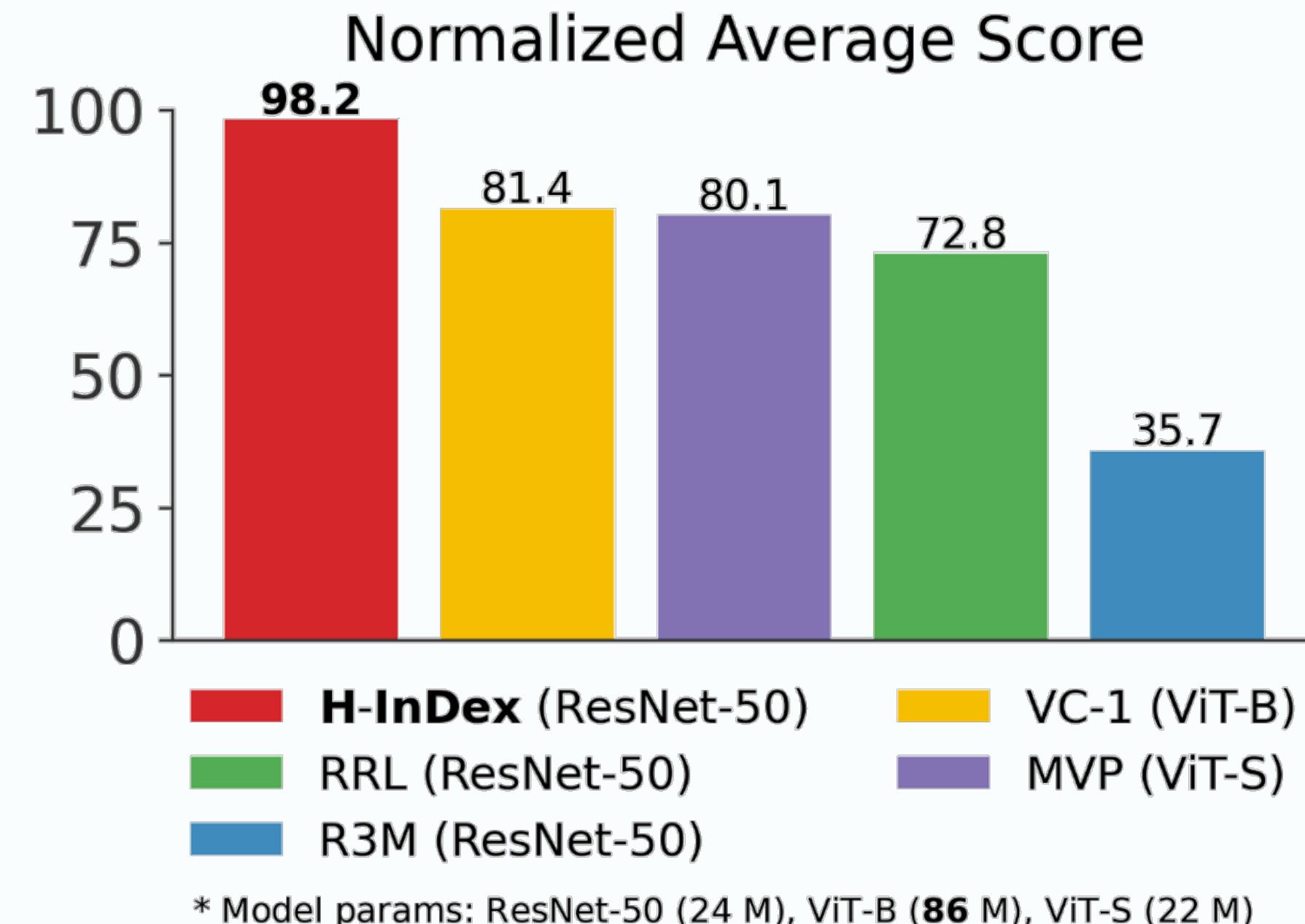
**Place Inside**



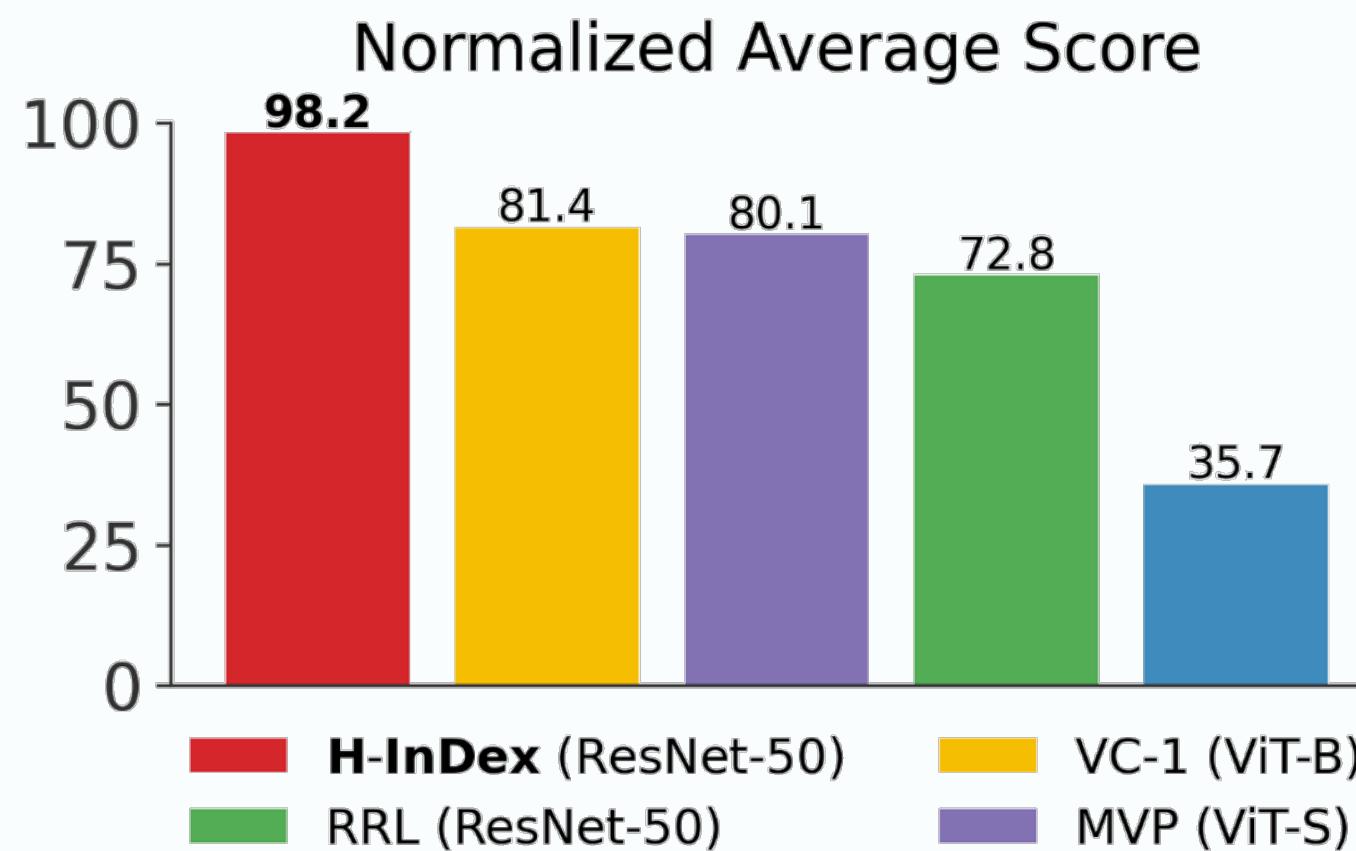
**Relocate**



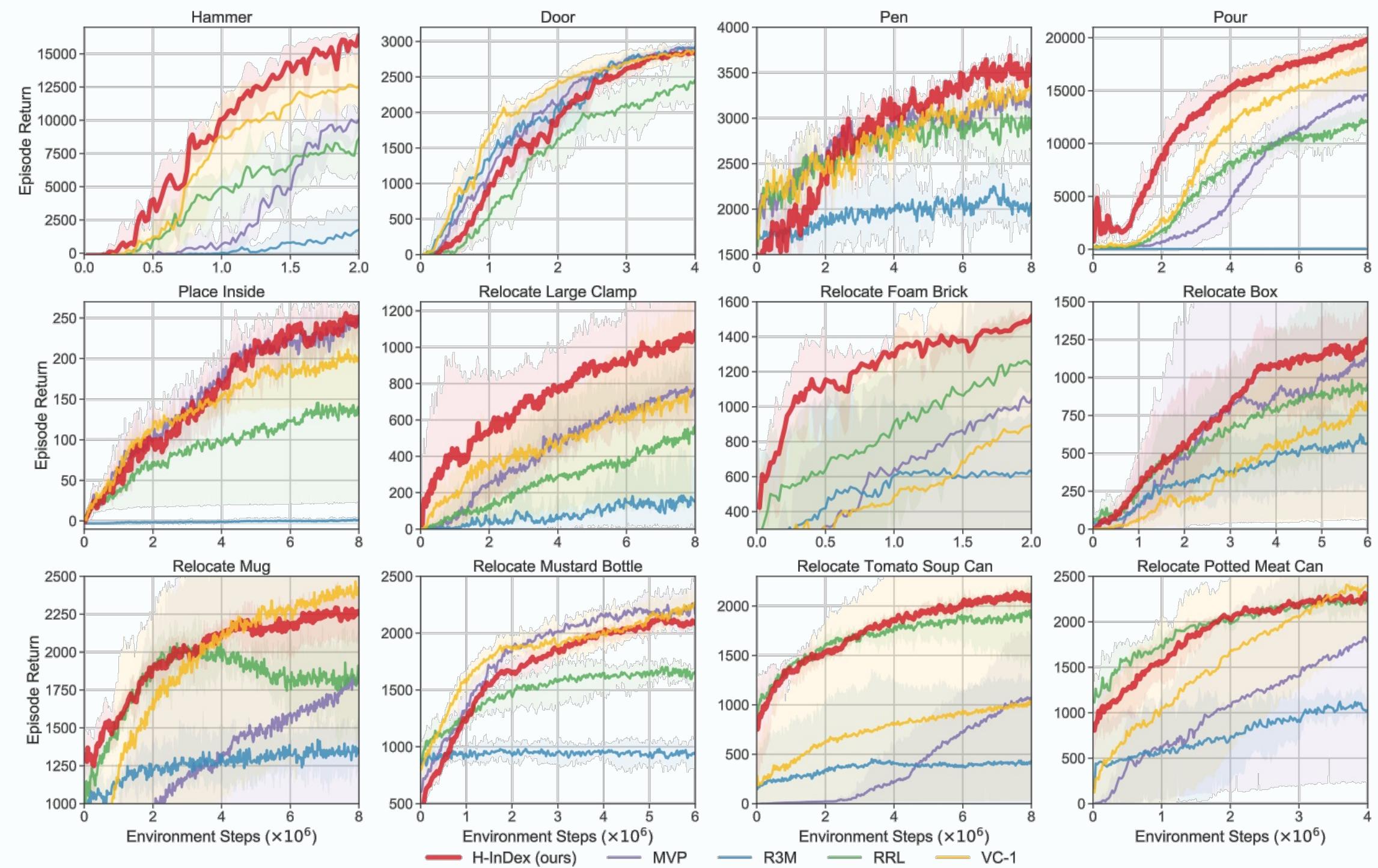
# Strong Results in **12** Dexterous Manipulation Tasks



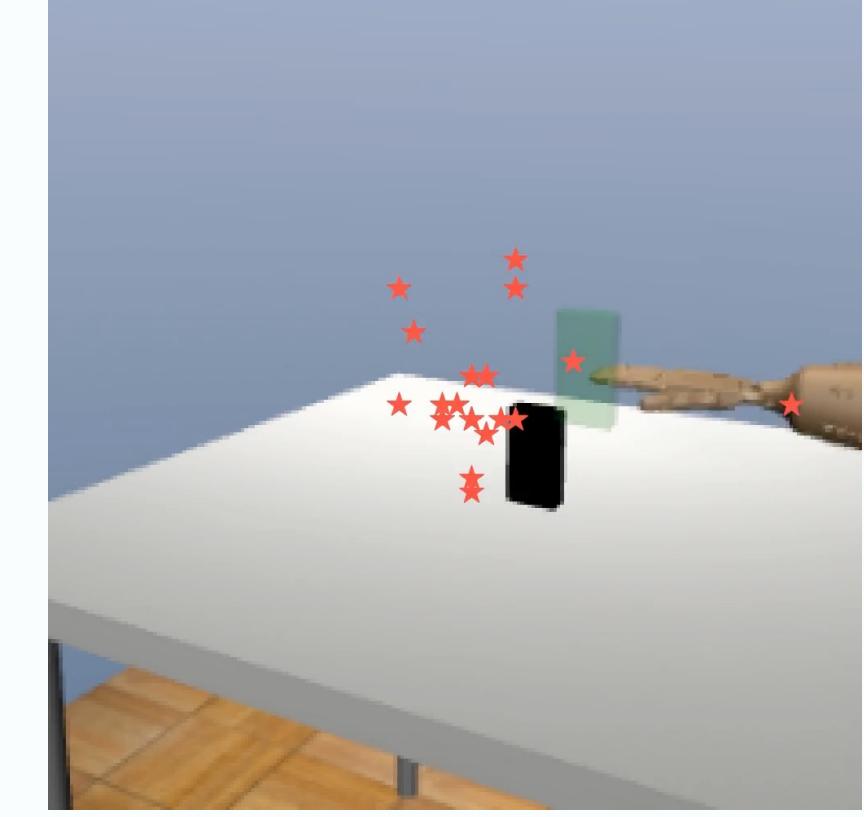
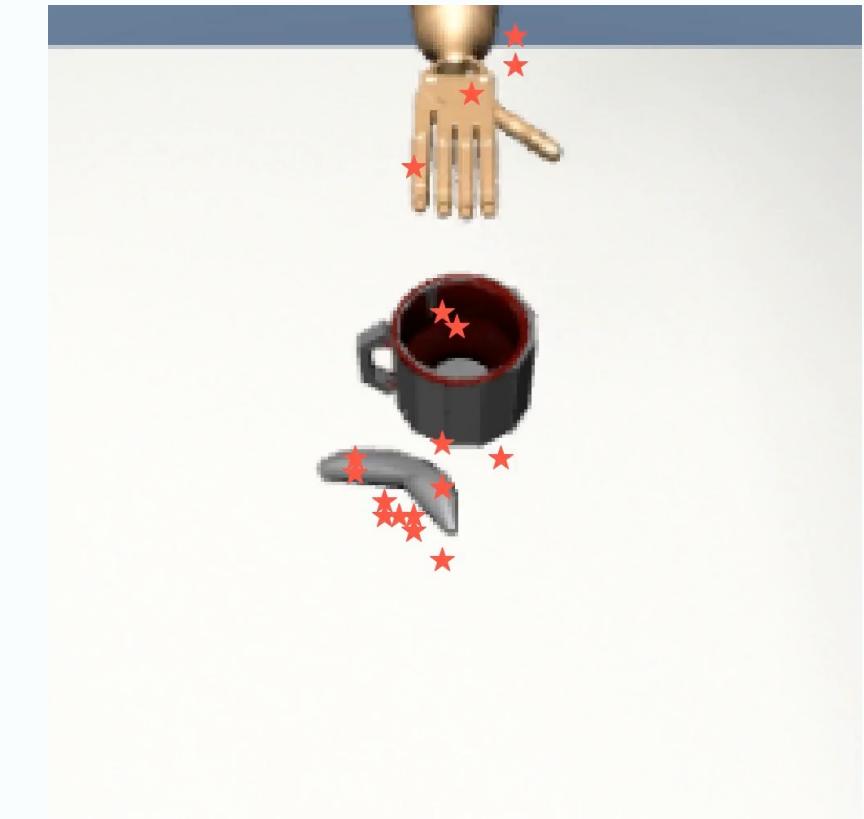
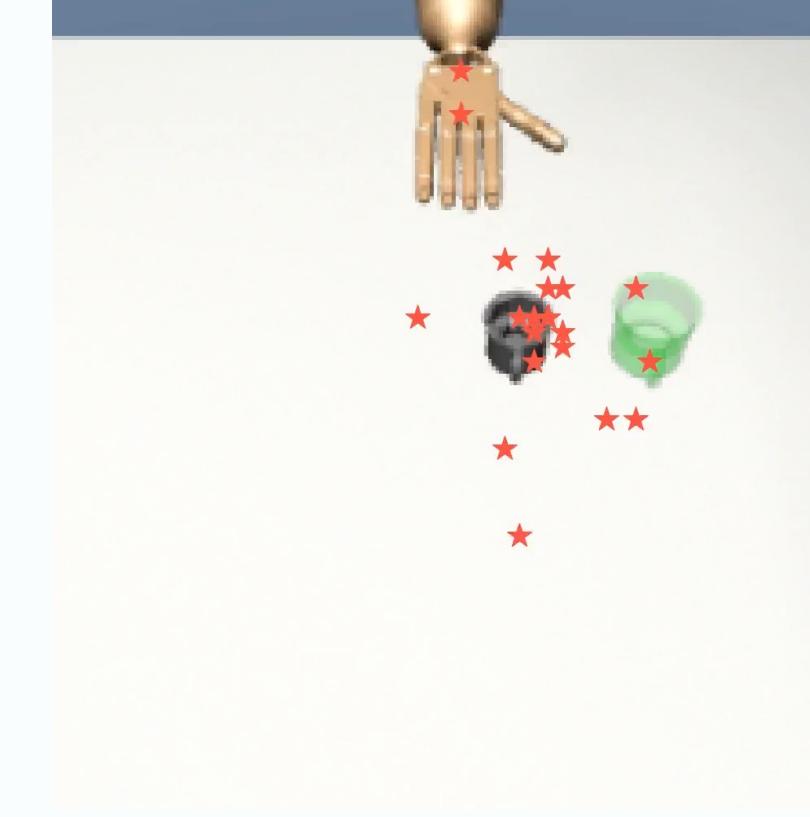
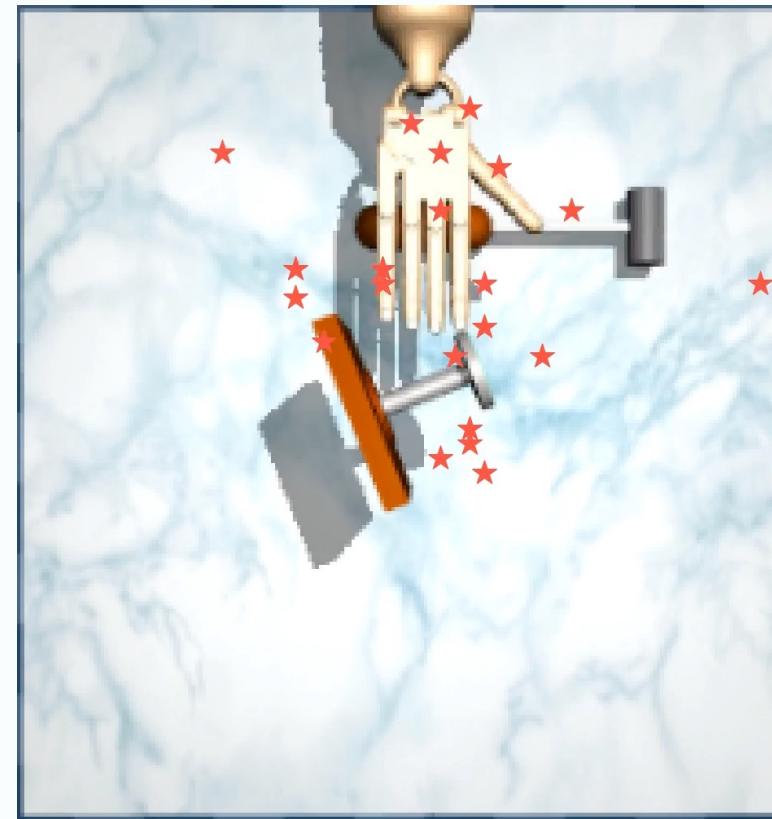
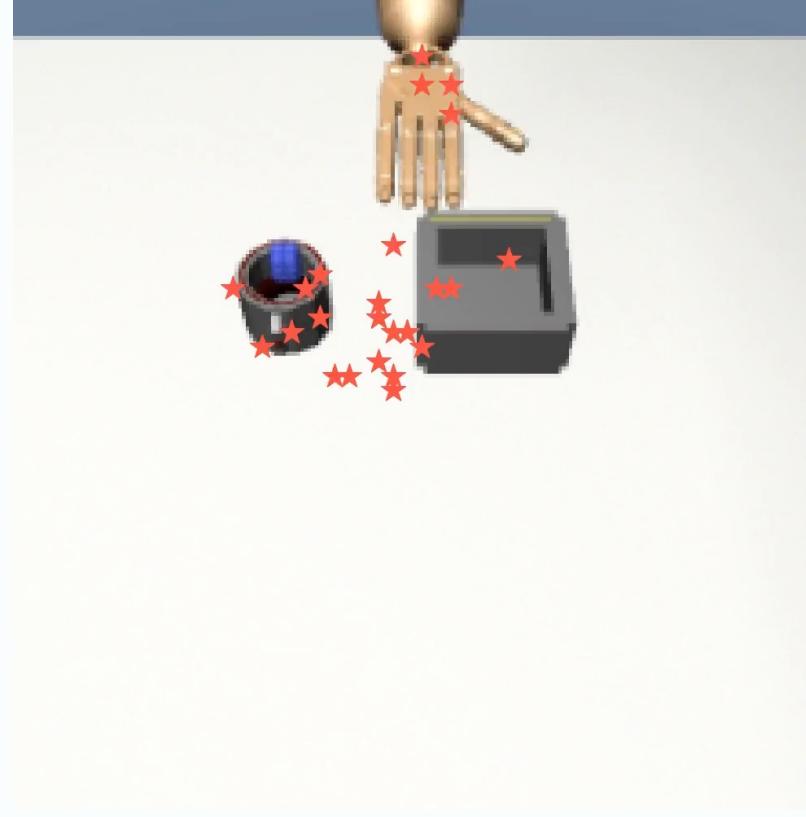
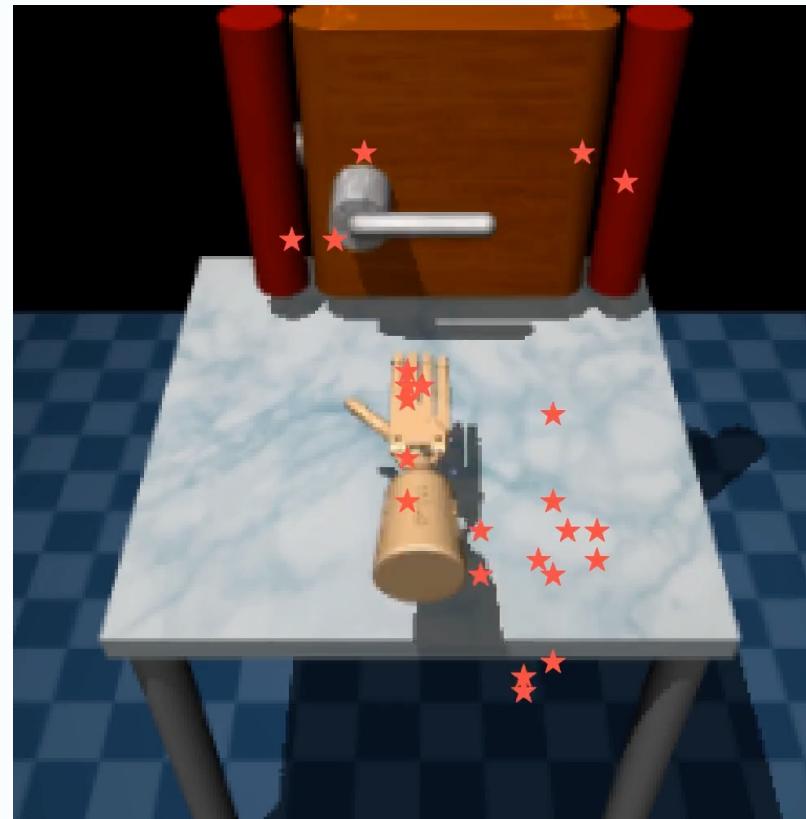
# Strong Results in 12 Dexterous Manipulation Tasks



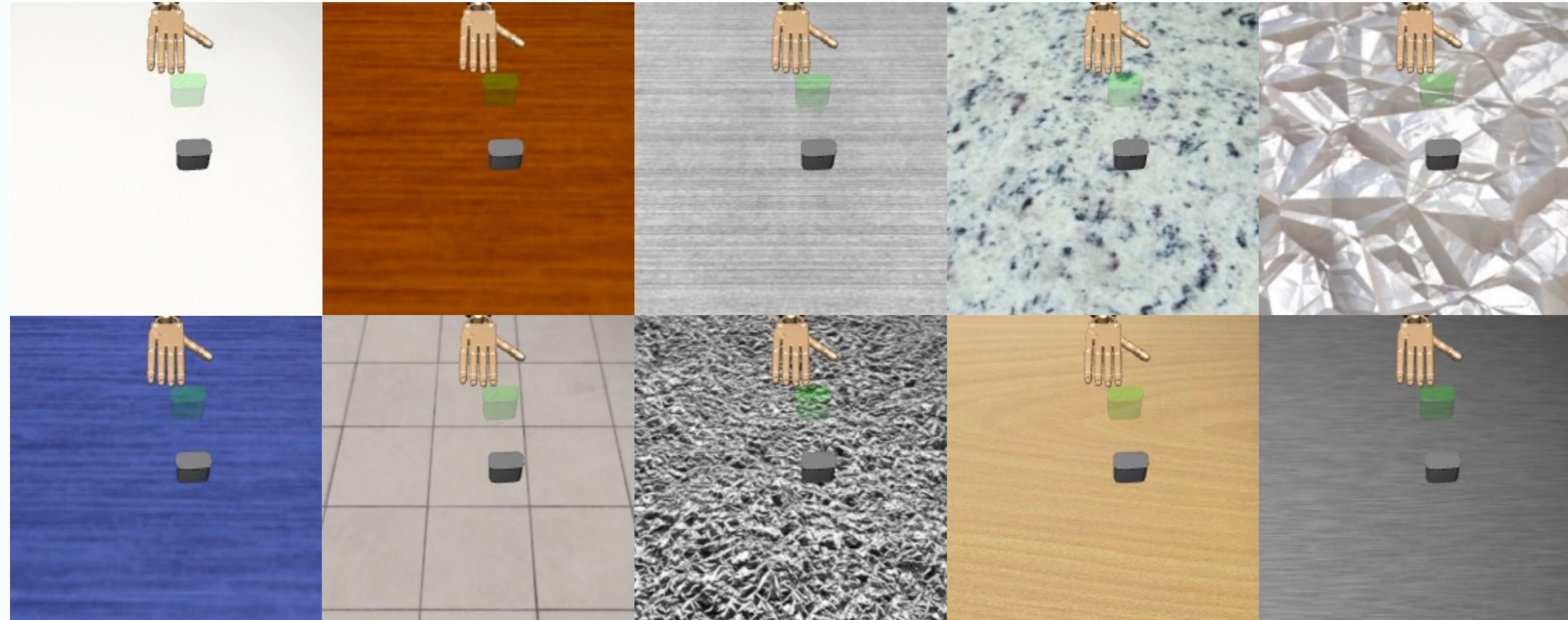
\* Model params: ResNet-50 (24 M), ViT-B (**86 M**), ViT-S (22 M)



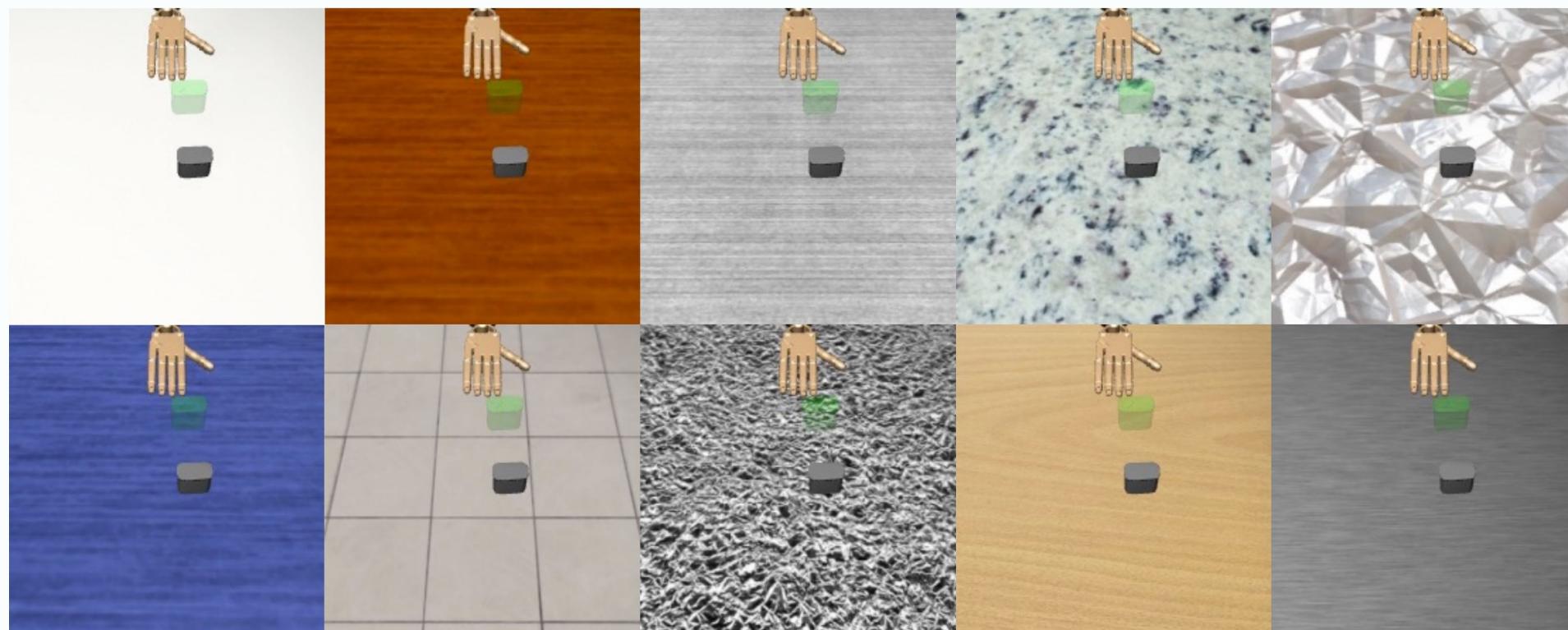
# Self-Supervised Keypoint Detection



# Generalize to Appearance



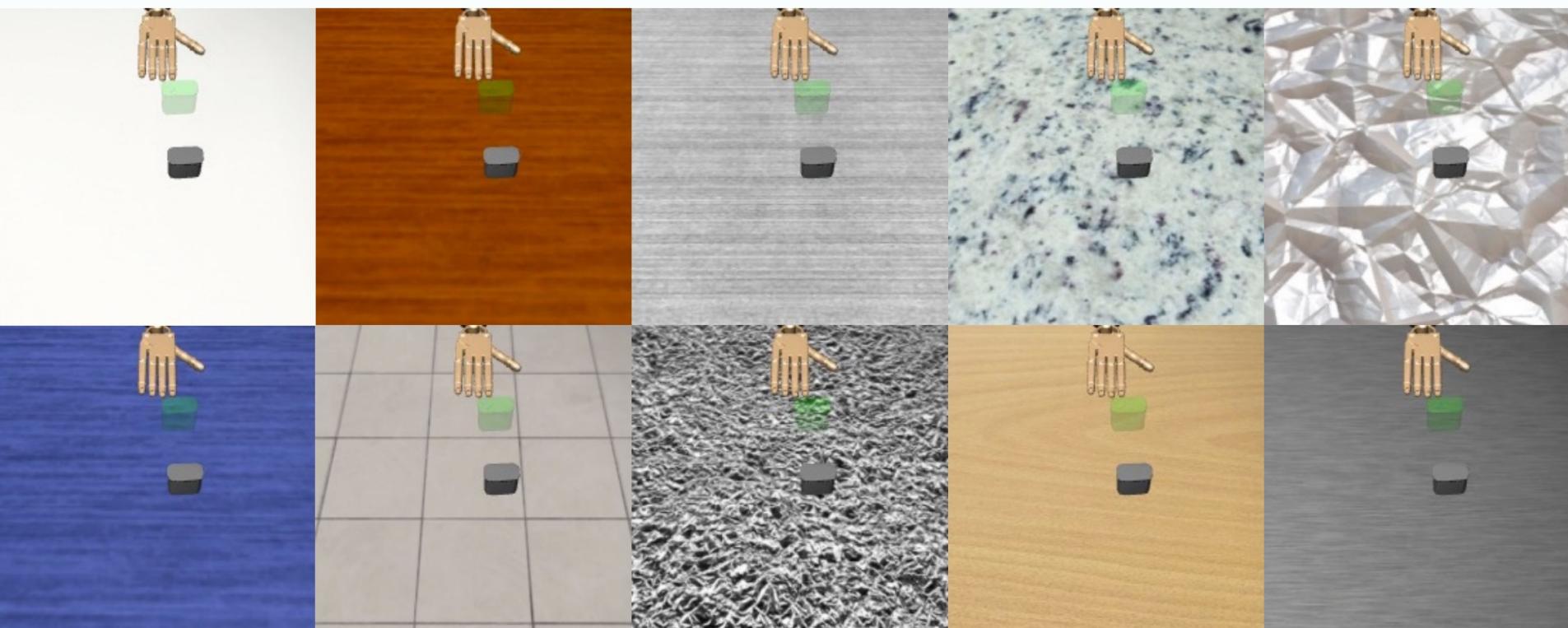
# Generalize to Appearance



Scene ID / Method	VC-1 [17]	H-InDex
Origin	<b>2391.74<math>\pm</math>602.83</b>	2240.37 $\pm$ 85.45

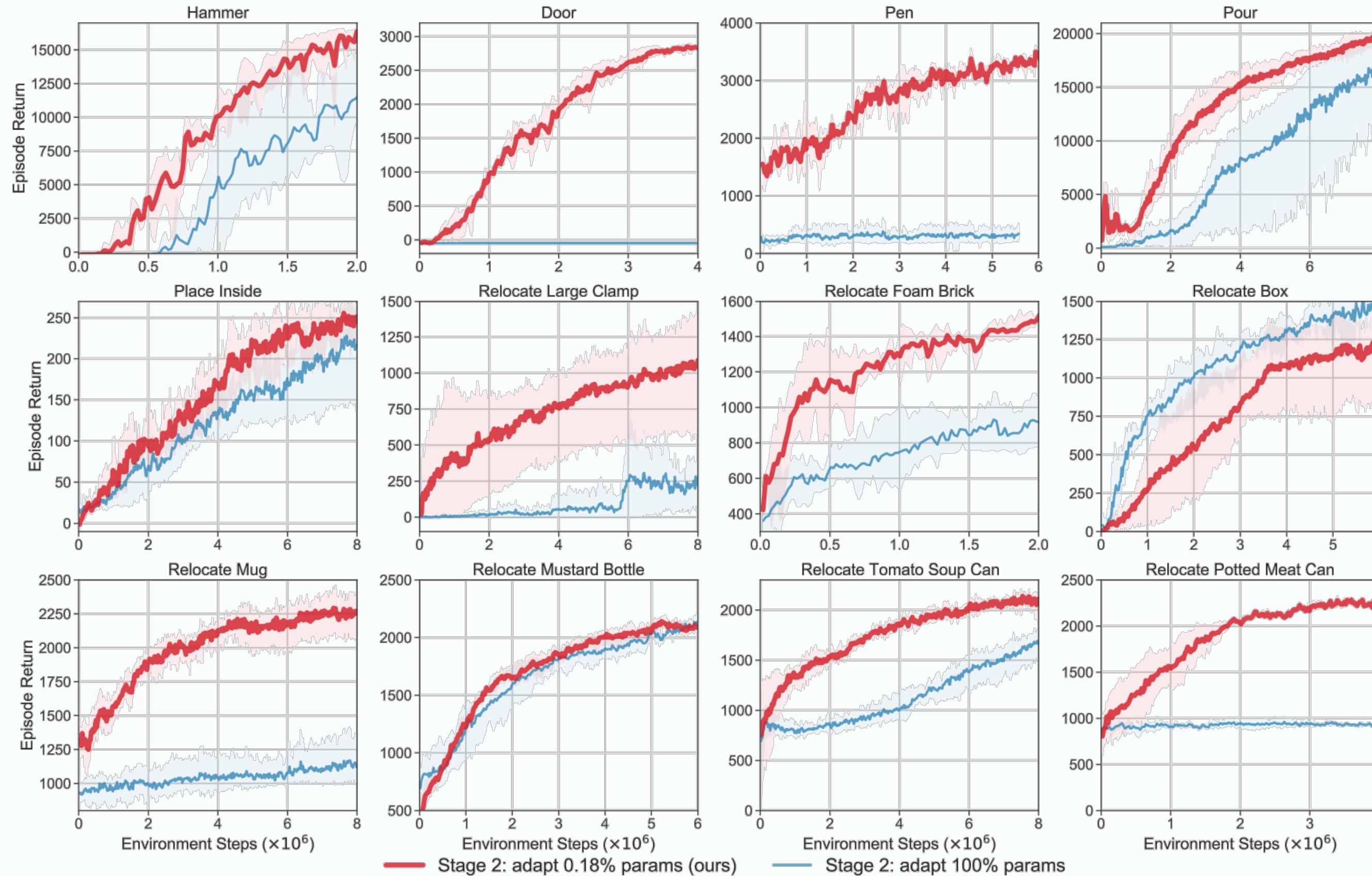
# Generalize to Appearance

**H-InDex** helps visual generalization.



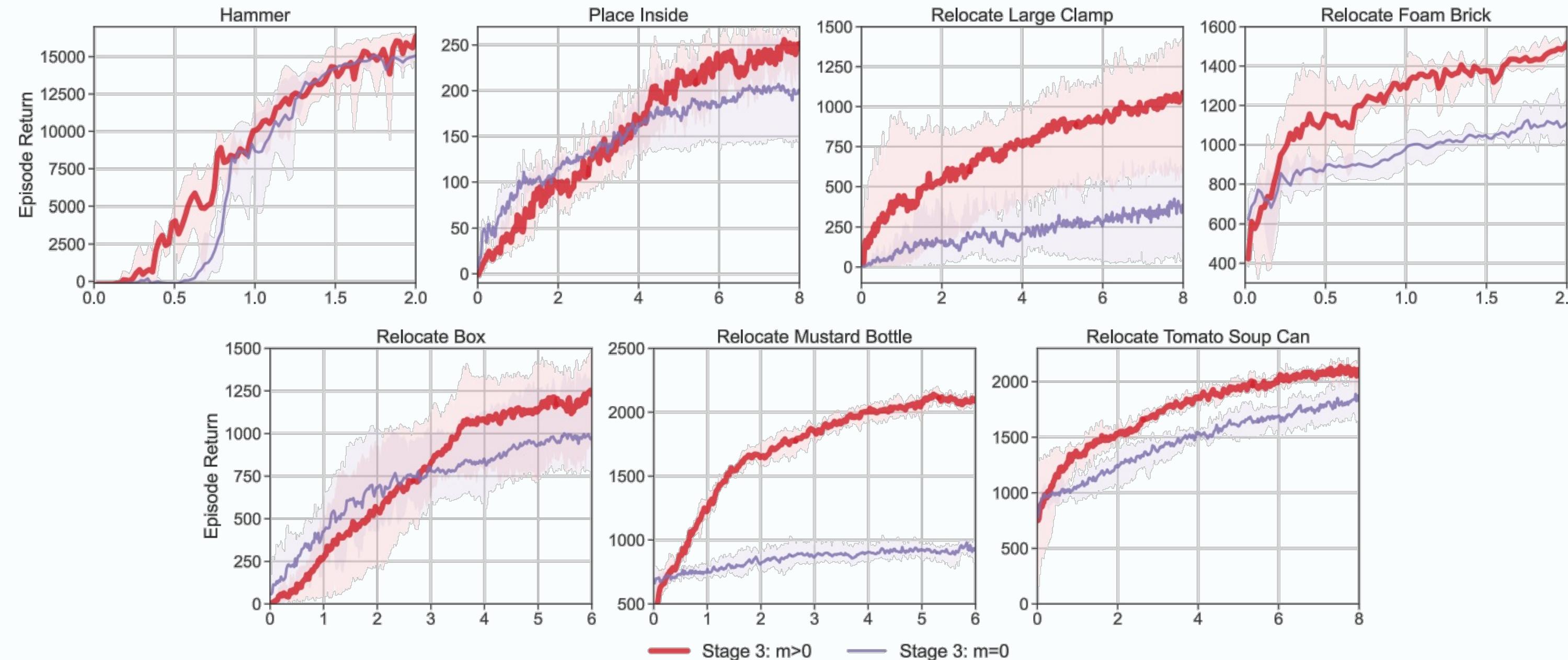
Scene ID / Method	VC-1 [17]	H-InDex
Origin	<b>2391.74±602.83</b>	2240.37±85.45
1	896.28±1006.55	<b>915.95±922.65</b>
2	603.26±920.48	<b>771.28±793.51</b>
3	451.36±839.45	<b>578.42±764.09</b>
4	360.21±772.64	<b>472.66±715.03</b>
5	300.02±718.05	<b>393.32±676.41</b>
6	256.80±673.16	<b>340.07±639.68</b>
7	224.21±635.56	<b>298.20±608.54</b>
8	226.59±610.58	<b>265.82±581.03</b>
9	214.30±581.76	<b>239.60±556.80</b>
Average	392.56	475.04

# Full Finetuning in Stage 2 is Harmful



# Updating BatchNorm Slightly during RL is Critical

$$\mu^{(i)} \leftarrow (1 - m) \cdot \mu^{(i)} + m \cdot \mathbb{E}[x^{(i)}],$$
$$(\sigma^{(i)})^2 \leftarrow (1 - m) \cdot (\sigma^{(i)})^2 + m \cdot \text{Var}[x^{(i)}]$$



# Visual Representations for Generalizable Robotic Manipulation

1

## Geometric Prior

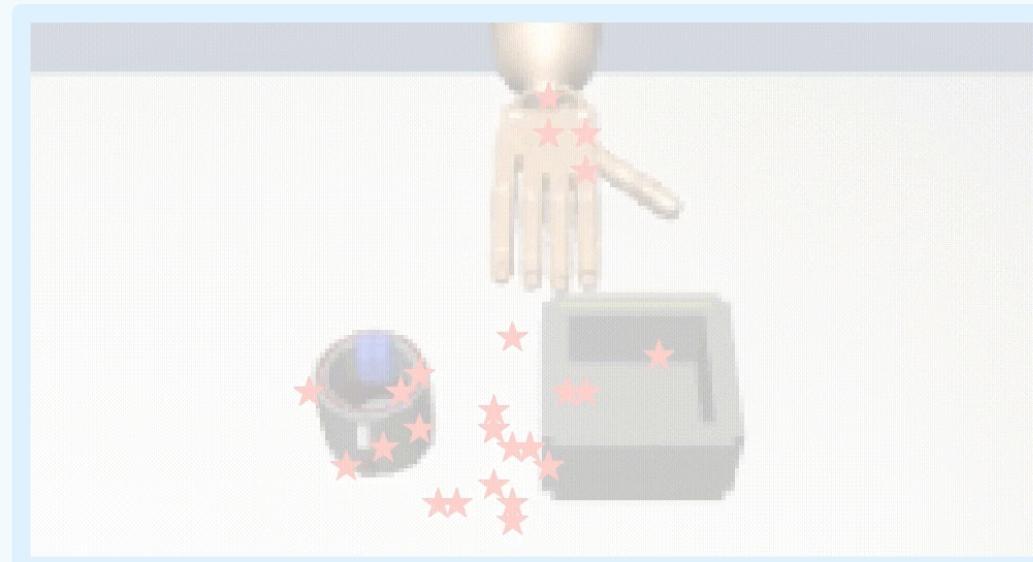


[2] Ze et al., “Visual Reinforcement Learning with Self-Supervised 3D Representations”, RA-L 2023 & IROS 2023.

[3] Ze et al., “GNFactor: Multi-Task Real Robot Learning with Generalizable Neural Feature Fields”, CoRL 2023 Oral.

2

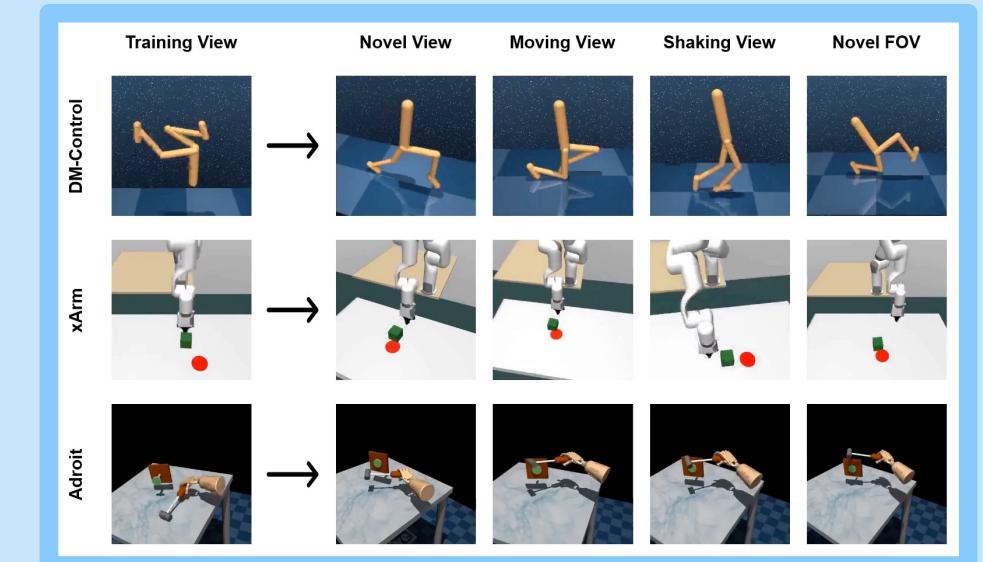
## Human Prior



[4] Ze et al., “H-InDex: Visual Reinforcement Learning with Hand-Informed Representations for Dexterous Manipulation”, NeurIPS 2023.

3

## Dynamics Prior



[5] Yang\*, Ze\* et al., “MoVie: Visual Model-Based Policy Adaptation for View Generalization”, NeurIPS 2023.

### Training View

DM-Control



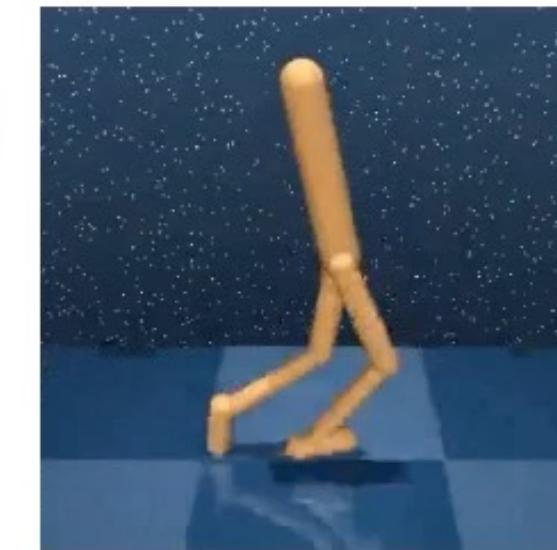
### Novel View



### Moving View



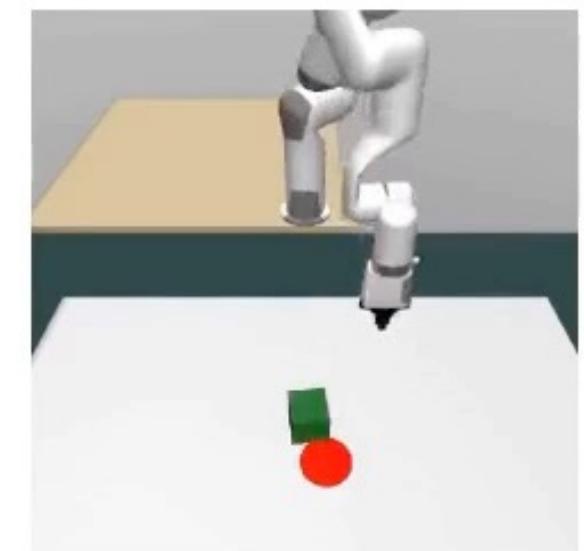
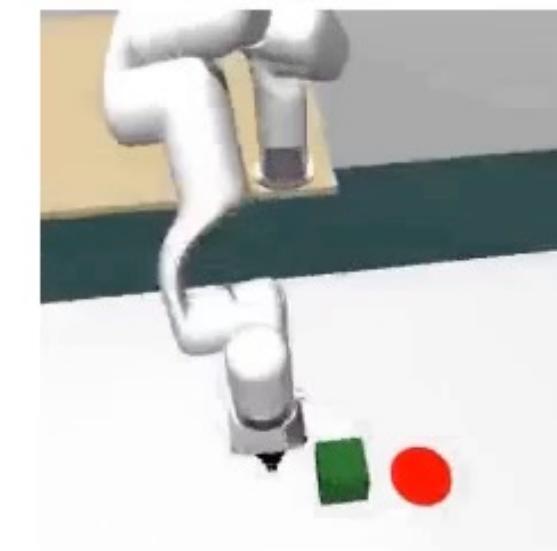
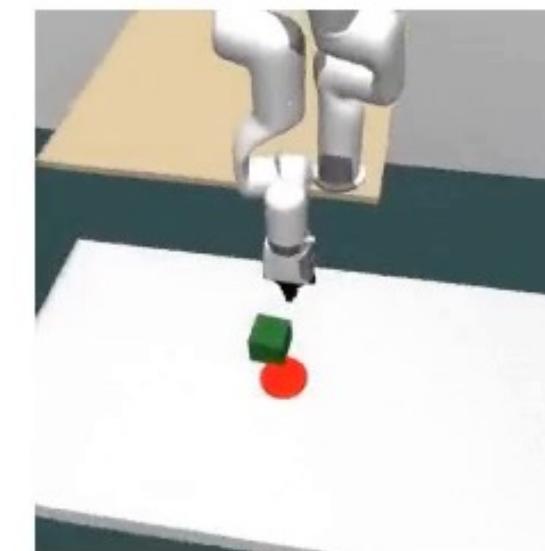
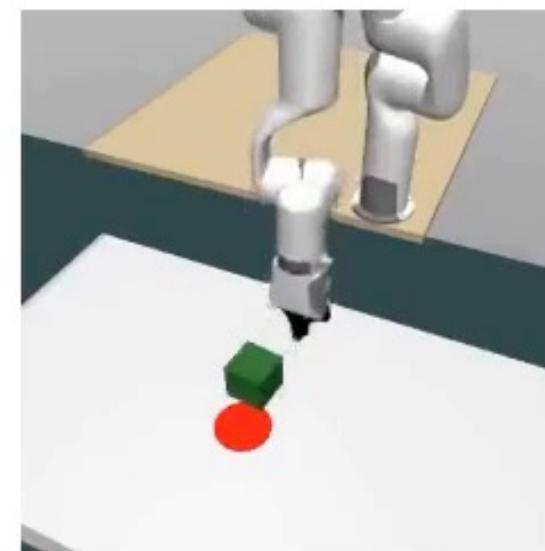
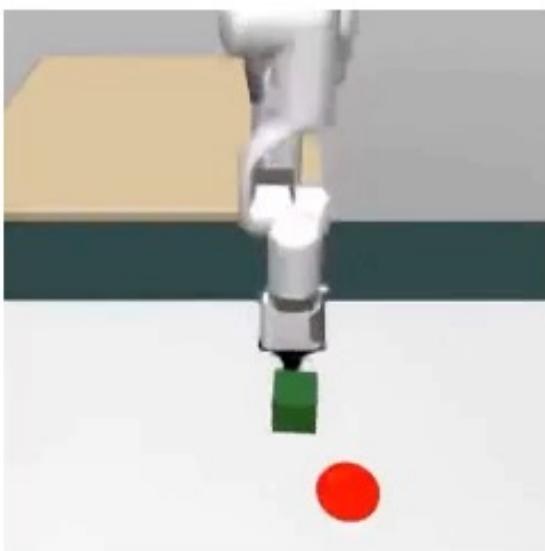
### Shaking View



### Novel FOV



xArm



Adroit



**Only trained on one fixed camera view,  
our agents generalize to diverse unseen views.**

# MoVie: Visual Model-Based Policy Adaptation for View Generalization

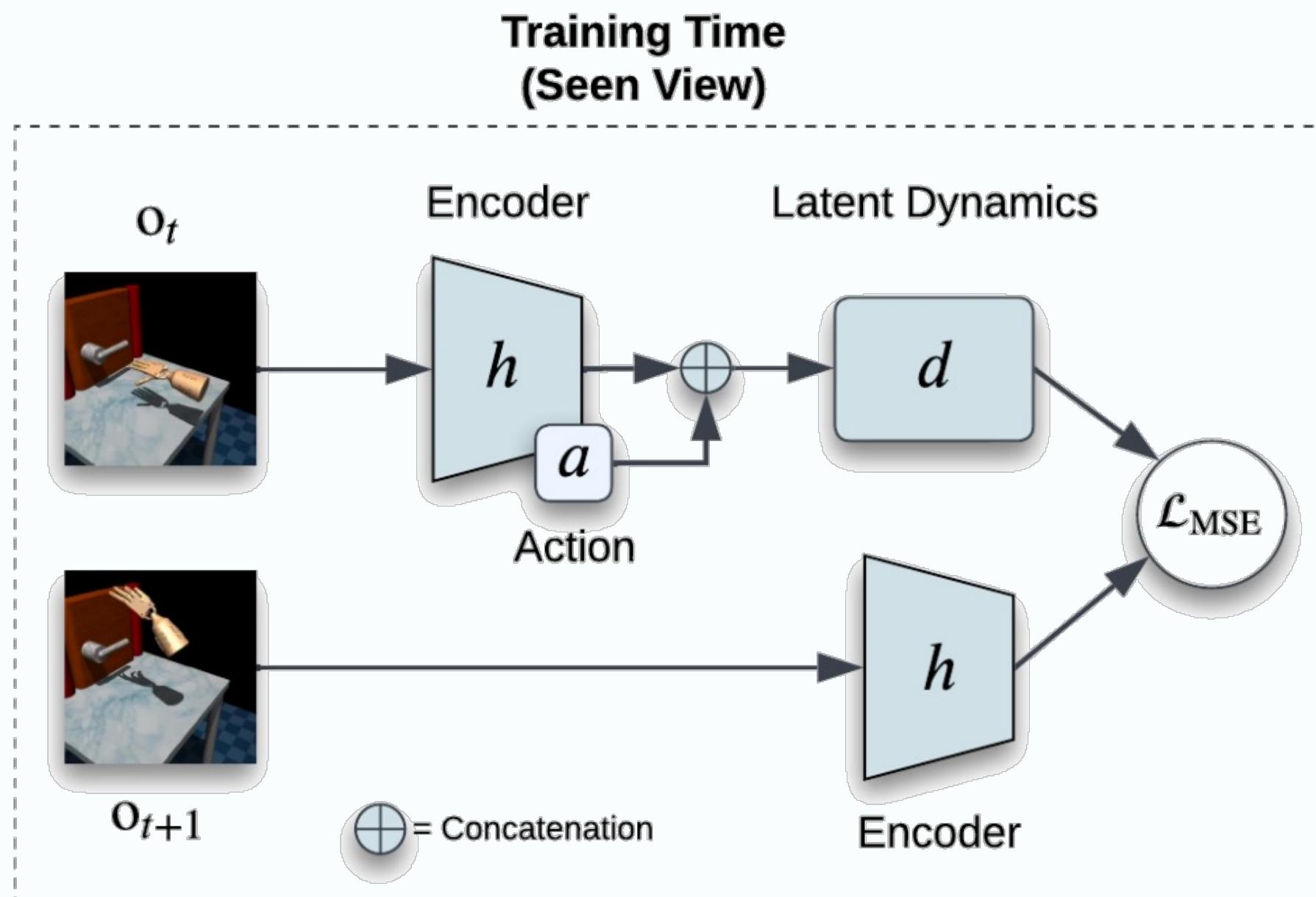
Sizhe Yang<sup>12\*</sup> Yanjie Ze<sup>13\*</sup> Huazhe Xu<sup>415</sup>

<sup>1</sup>Shanghai Qi Zhi Institute    <sup>2</sup>UESTC    <sup>3</sup>Shanghai Jiao Tong University  
<sup>4</sup>Tsinghua University, IIIS    <sup>5</sup>Shanghai AI Lab

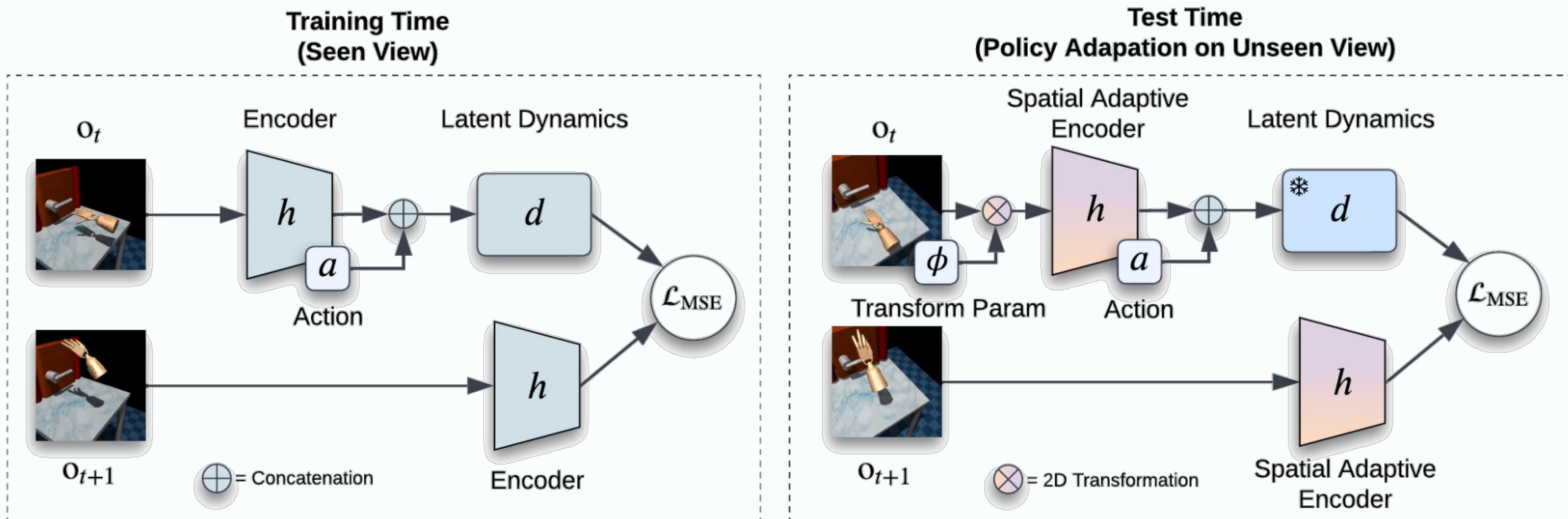
NeurIPS 2023



# Test-Time Adaptation by Dynamics

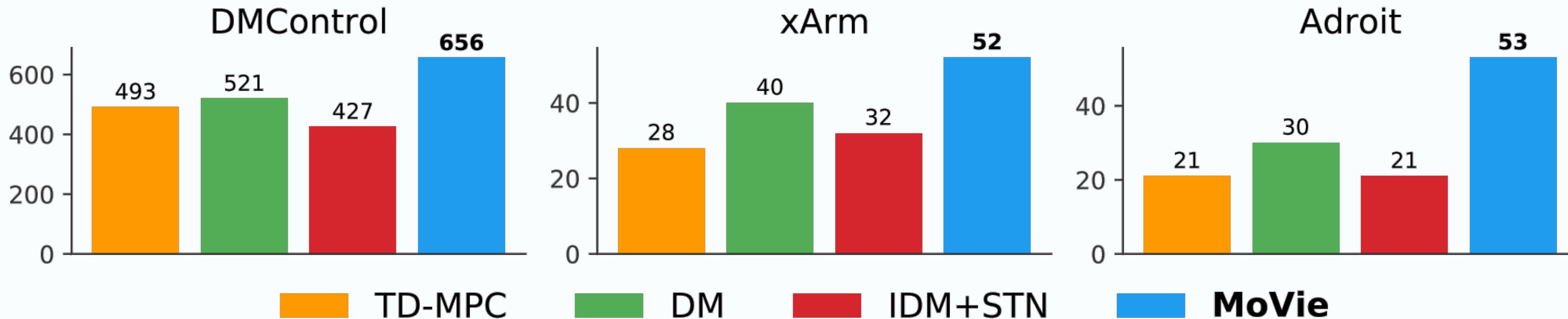


# Test-Time Adaptation by Dynamics



# Experiments

18 tasks across 3 domains

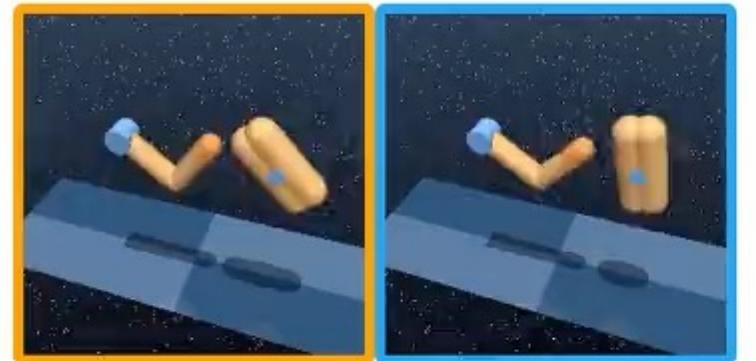


Finger, spin  
Pendulum, swingup  
Cheetah, run  
Walker, walk

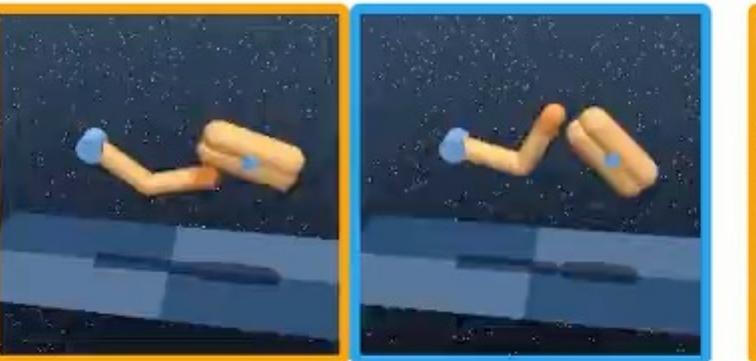
Training View



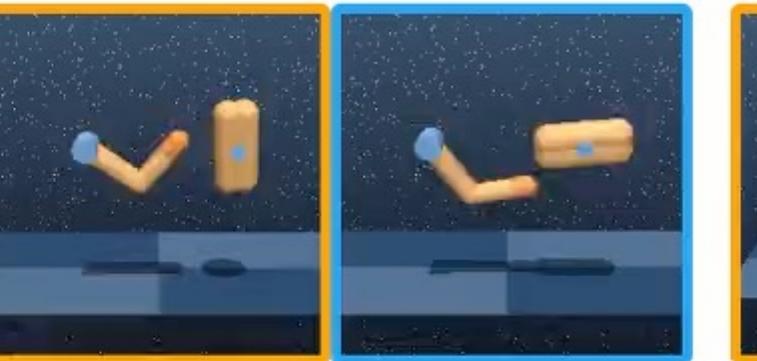
Novel View



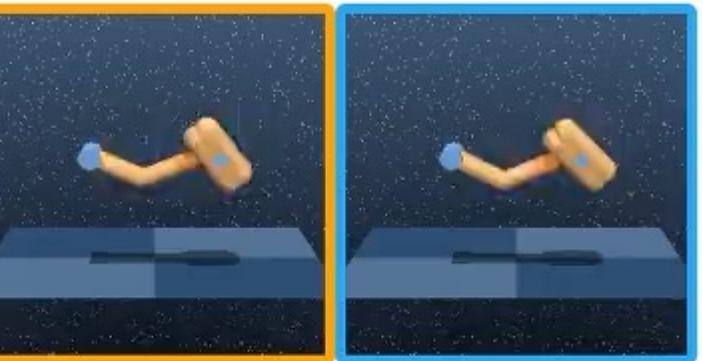
Moving View

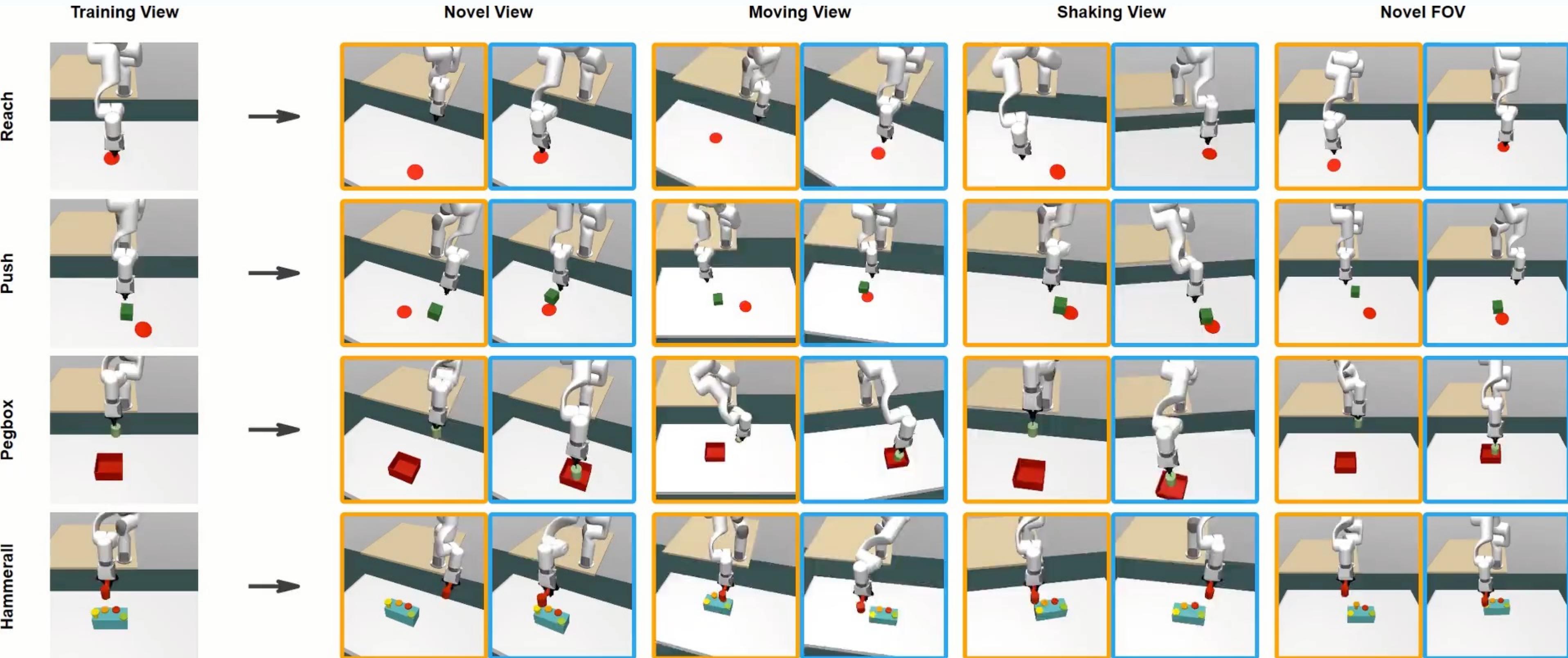


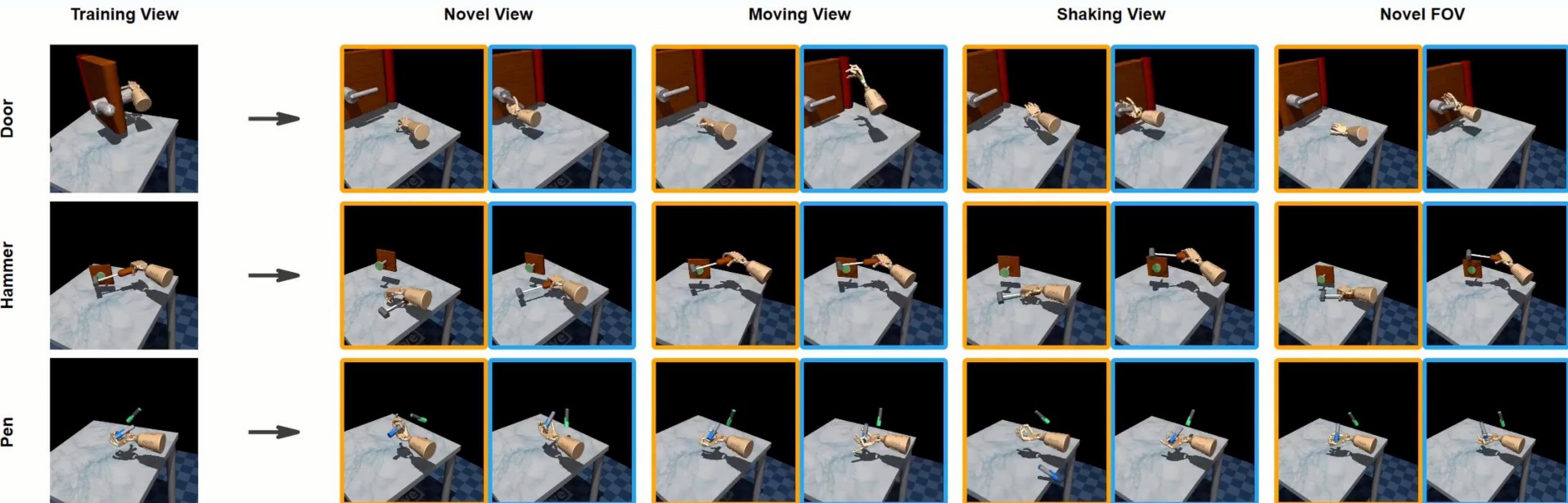
Shaking View



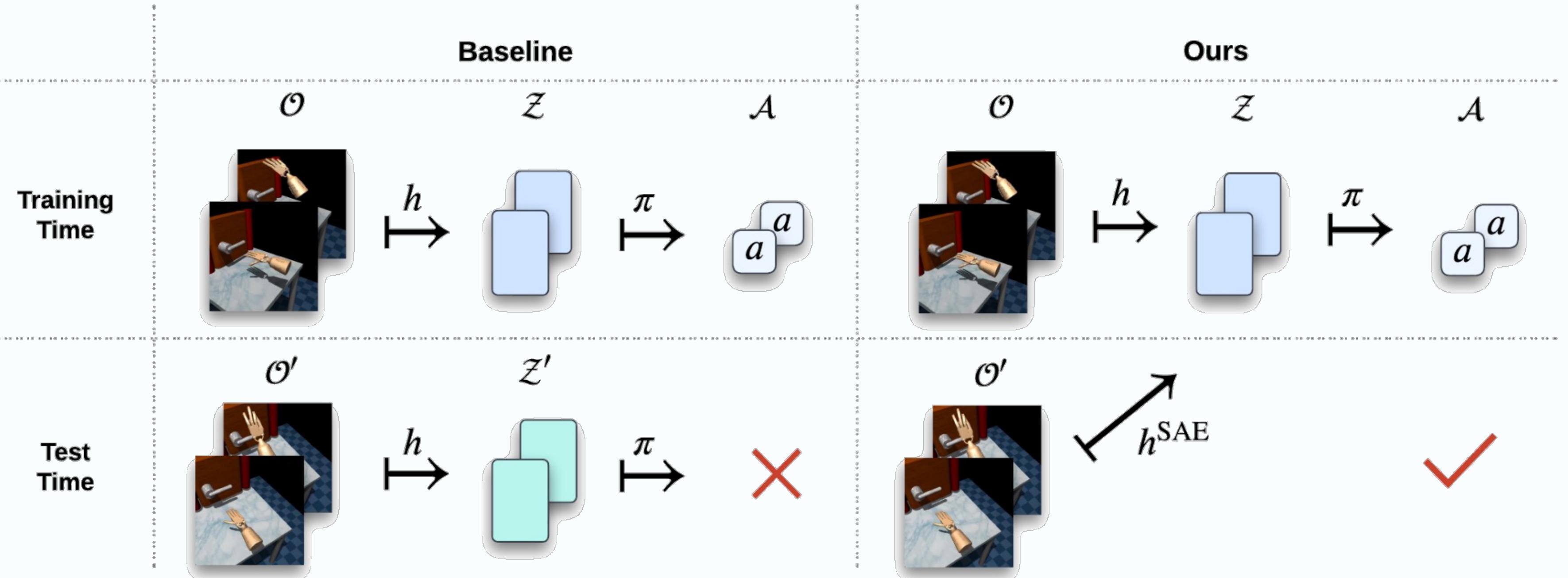
Novel FOV







# Explanation: Dynamics as Supervision



# Explanation: Dynamics as Supervision

Fintuning dynamics → only use test-time information

Freezing dynamics → dynamics consistency between training and test

Cheetah-run	Novel View	Moving View	Shaking View	Novel FOV
MoVie	<b>342.39±54.95</b>	<b>365.22±42.66</b>	<b>493.54±56.80</b>	$532.94\pm19.74$
Finetune DM	$273.01\pm11.29$	$331.55\pm22.22$	$476.25\pm30.25$	<b>561.94±37.73</b>

# Visual Representations for Generalizable Robotic Manipulation

1

## Geometric Prior

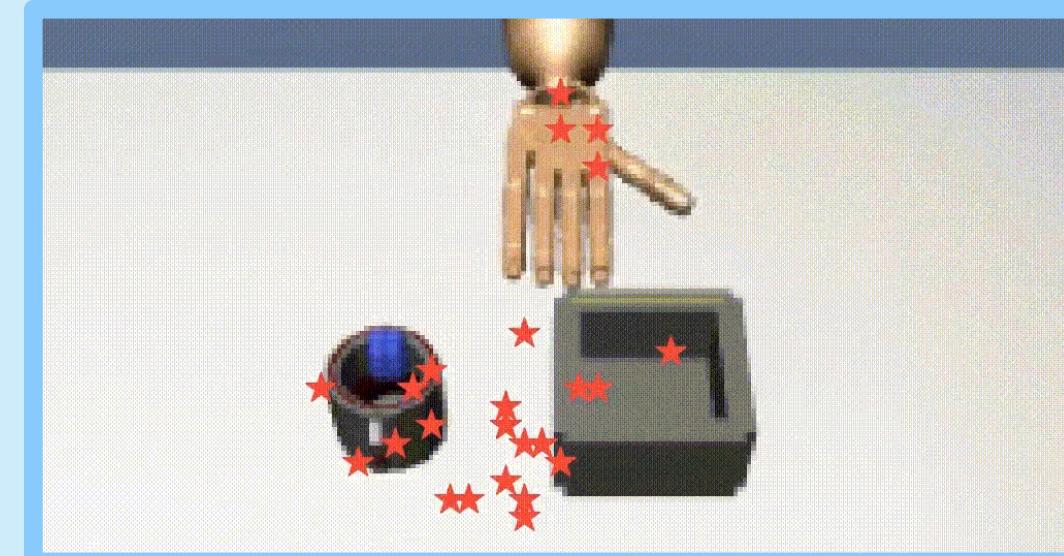


[2] Ze et al., “Visual Reinforcement Learning with Self-Supervised 3D Representations”, **RA-L** 2023 & **IROS** 2023.

[3] Ze et al., “GNFactor: Multi-Task Real Robot Learning with Generalizable Neural Feature Fields”, **CoRL** 2023 Oral.

2

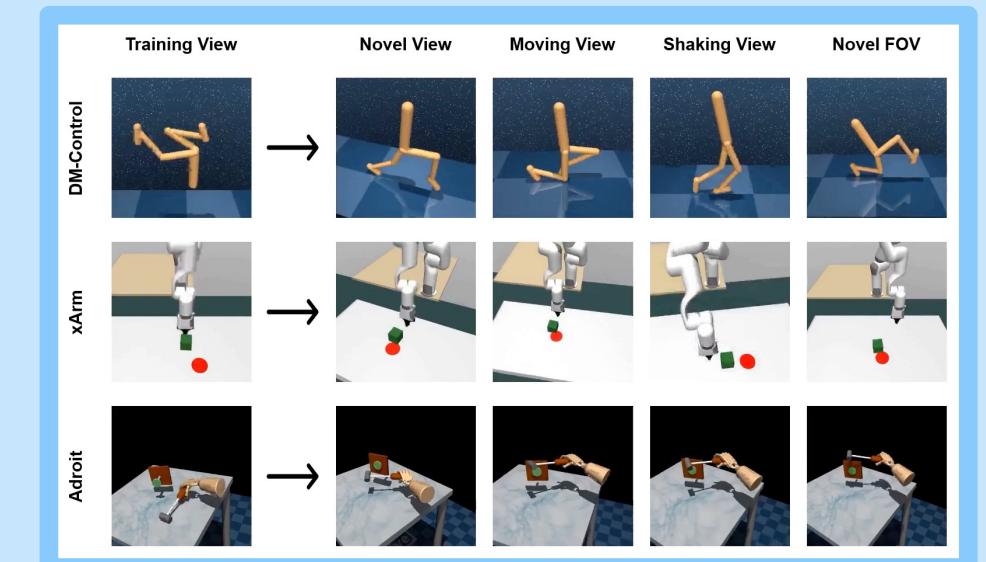
## Human Prior



[4] Ze et al., “H-InDex: Visual Reinforcement Learning with Hand-Informed Representations for Dexterous Manipulation”, **NeurIPS** 2023.

3

## Dynamics Prior



[5] Yang\*, Ze\* et al., “MoVie: Visual Model-Based Policy Adaptation for View Generalization”, **NeurIPS** 2023.

# More info



<https://yanjieze.com>

To Be Continued

