

# **Structure from Action:**

## **Articulated Object Structure Discovery with Active Interactions**

Structural and Compositional Learning on 3D Data workshop @ ICCV



Shuran Song

Columbia University  
Artificial Intelligence & Robotics Lab

# **Structure from Action:**

## **Articulated Object Structure Discovery with Active Interactions**

Structural and Compositional Learning on 3D Data workshop @ ICCV

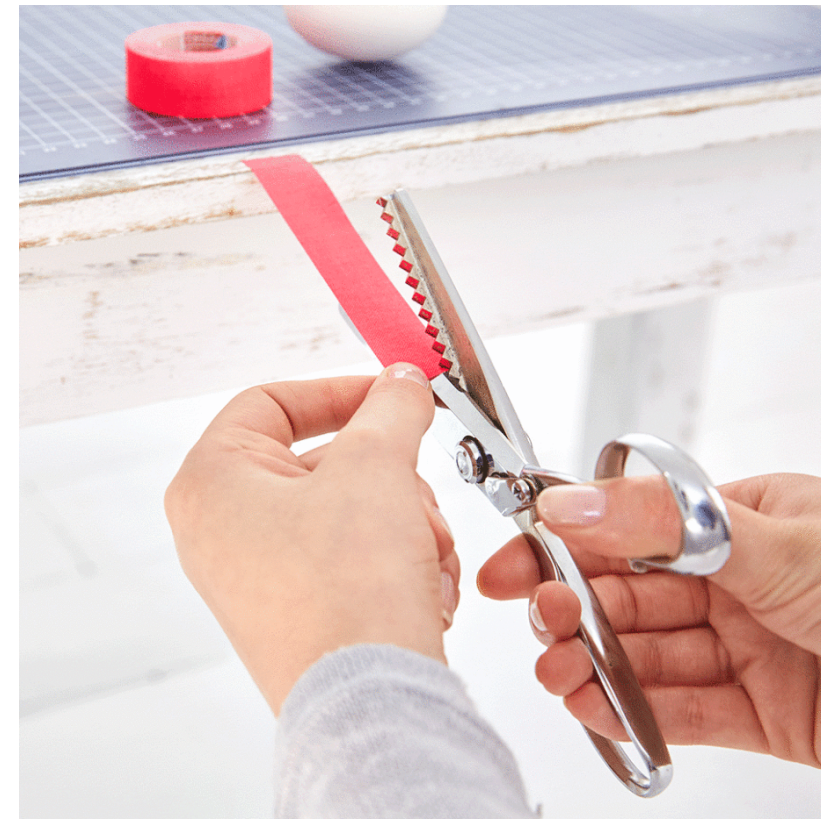


Shuran Song

Columbia University  
Artificial Intelligence & Robotics Lab



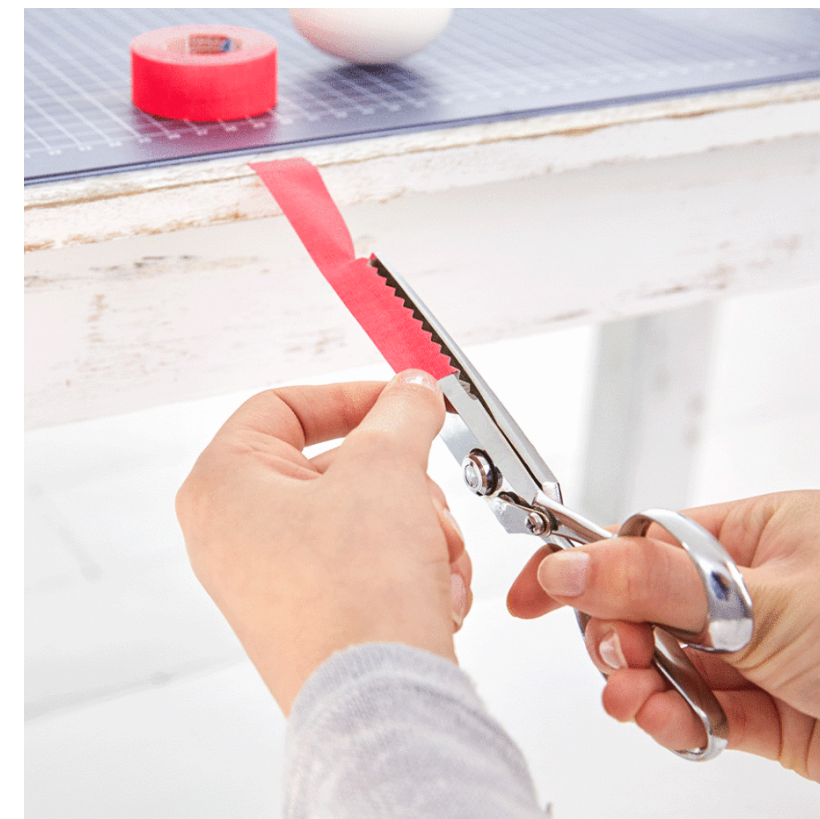
# Structures in Articulated Object



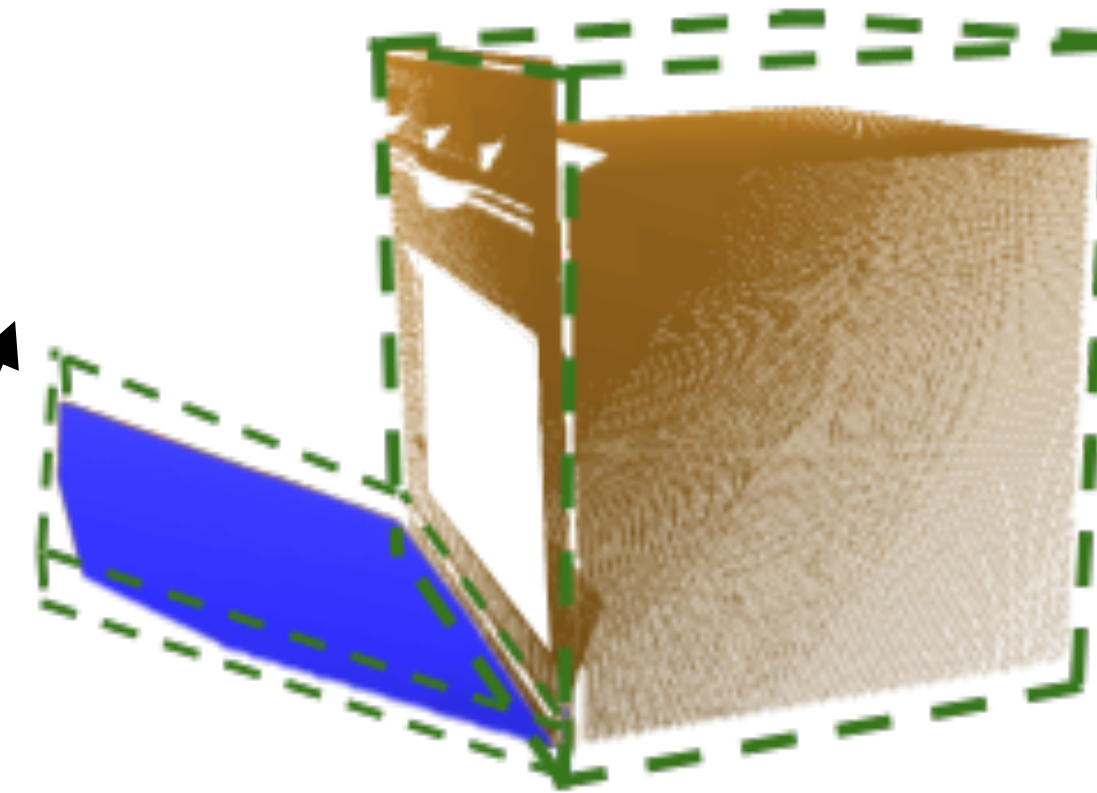
**Articulated Objects**



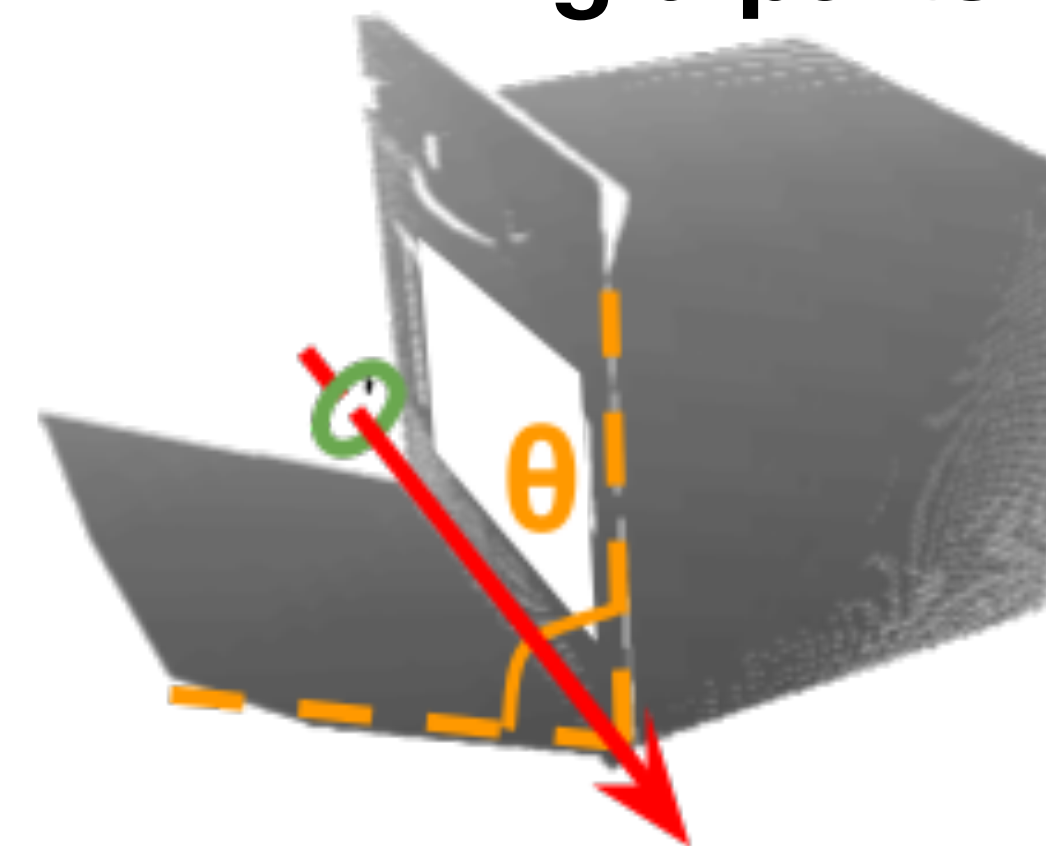
# Structures in Articulated Object



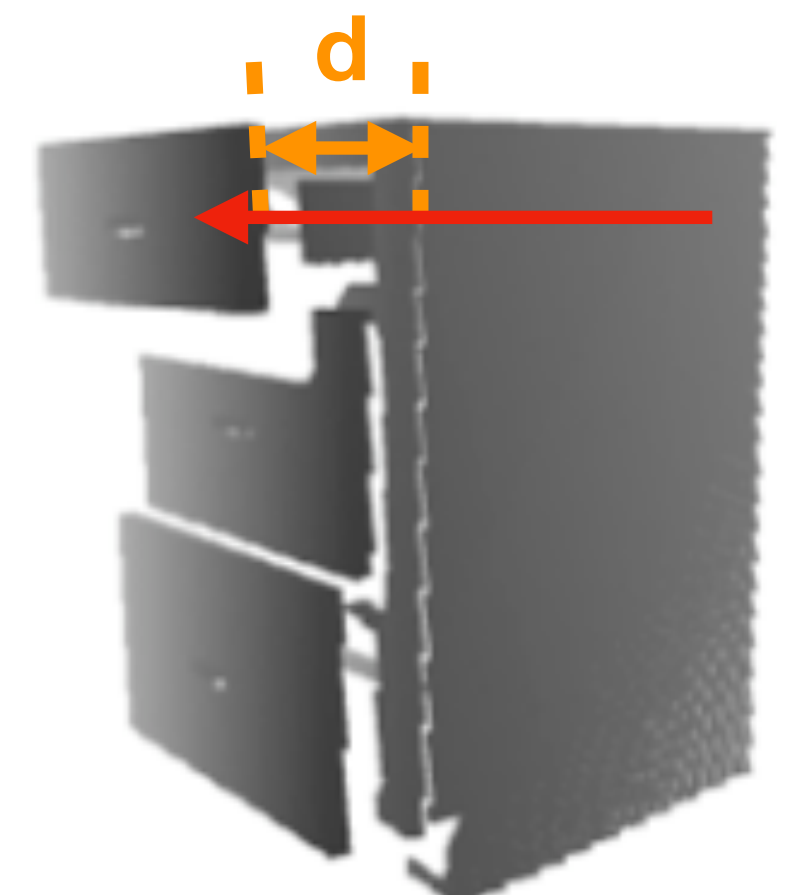
**Articulated Objects**



**Rigid parts**



**Joints (revolute)**

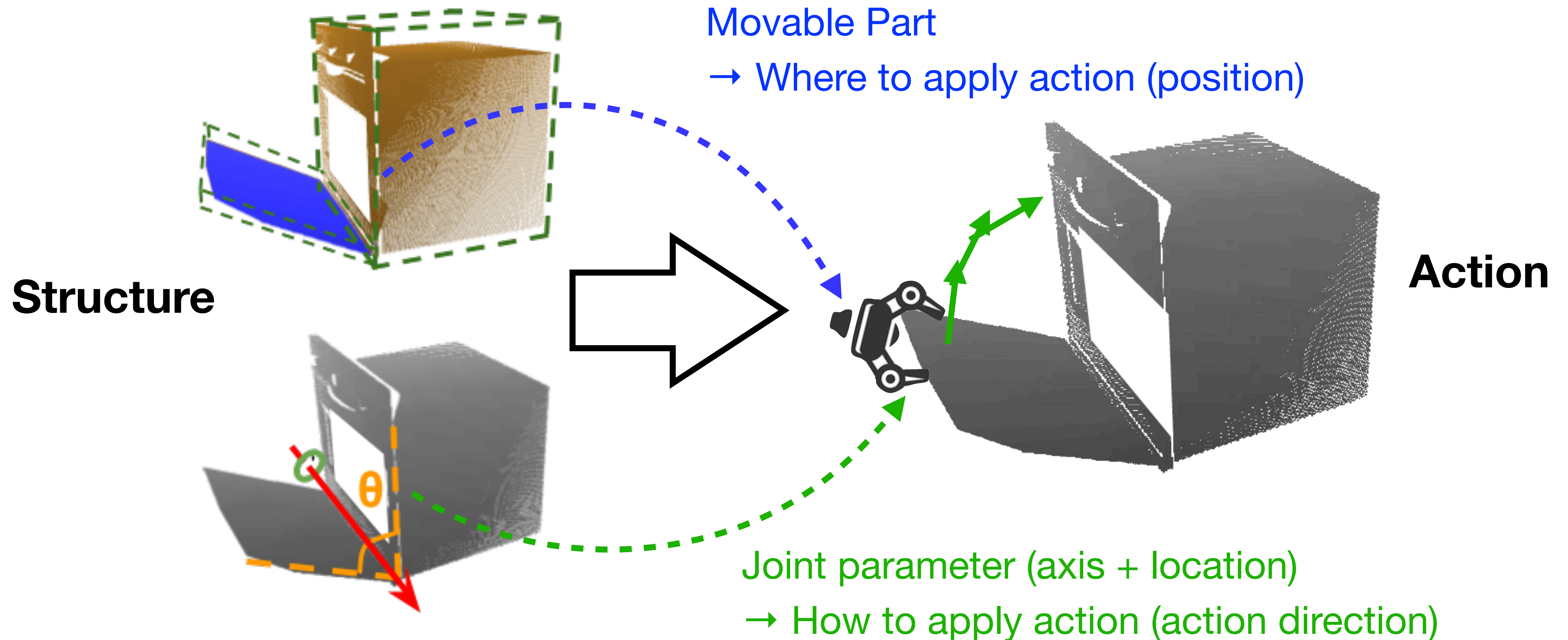


**Joint (Prismatic)**

**Kinematic Structure**

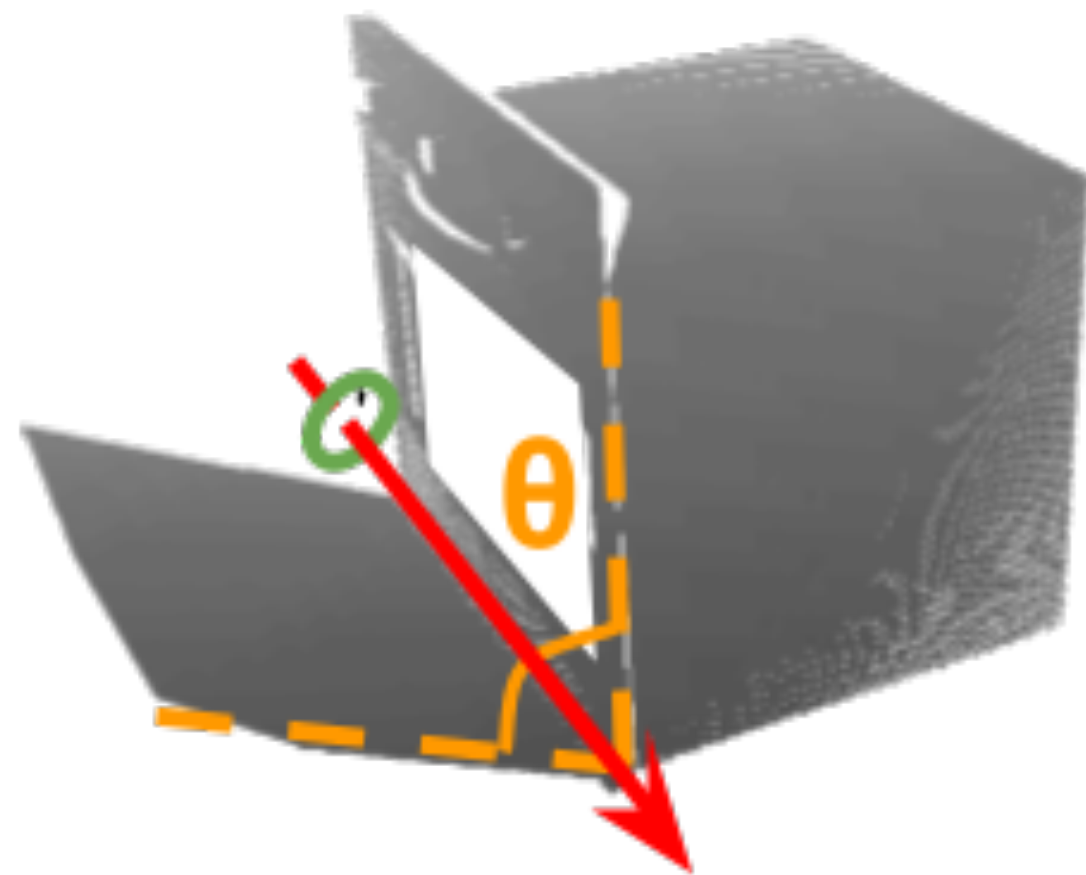
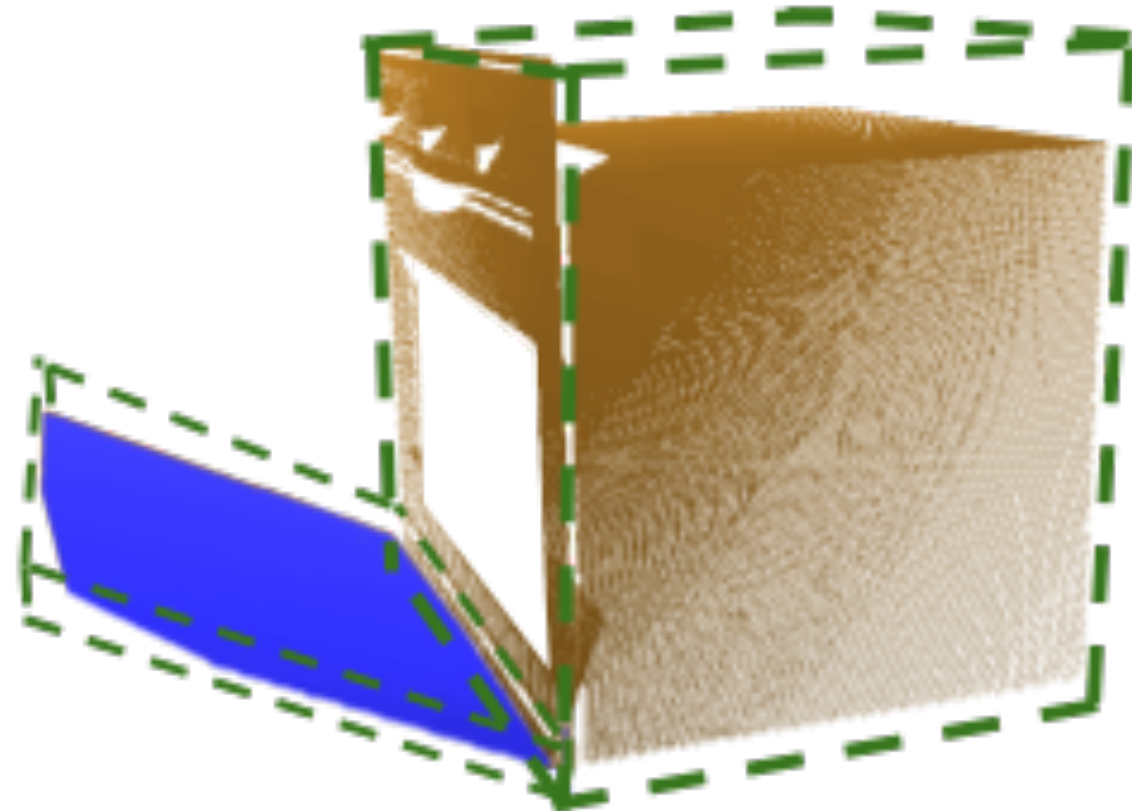


# Structure → Action



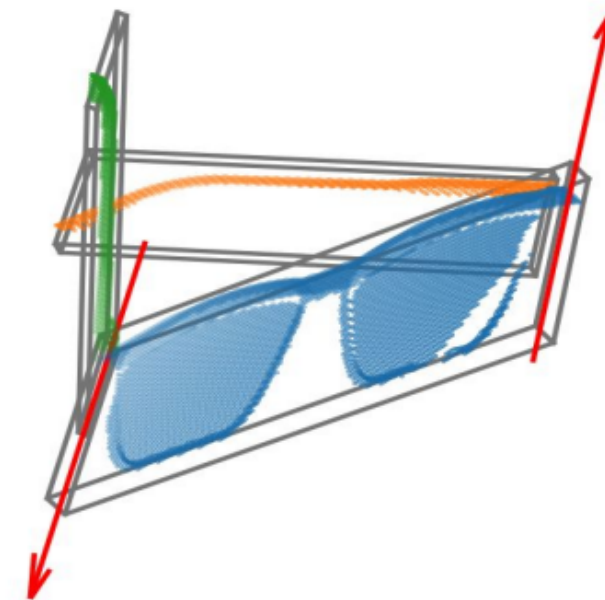
# Structure → Action

**Structure**



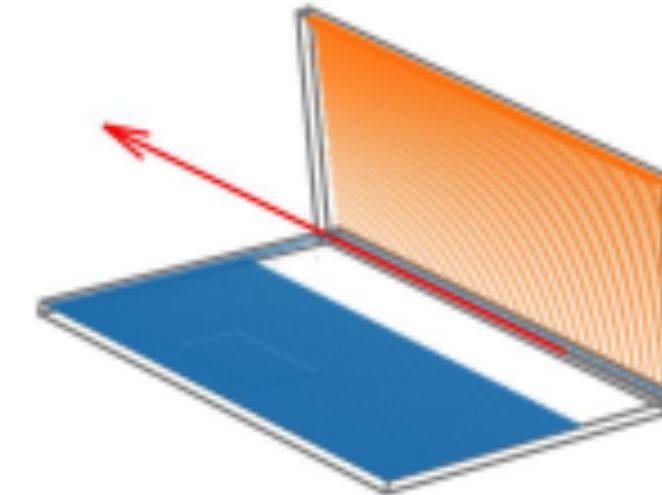
**Inferring structure from  
a single image**

**Model 1**



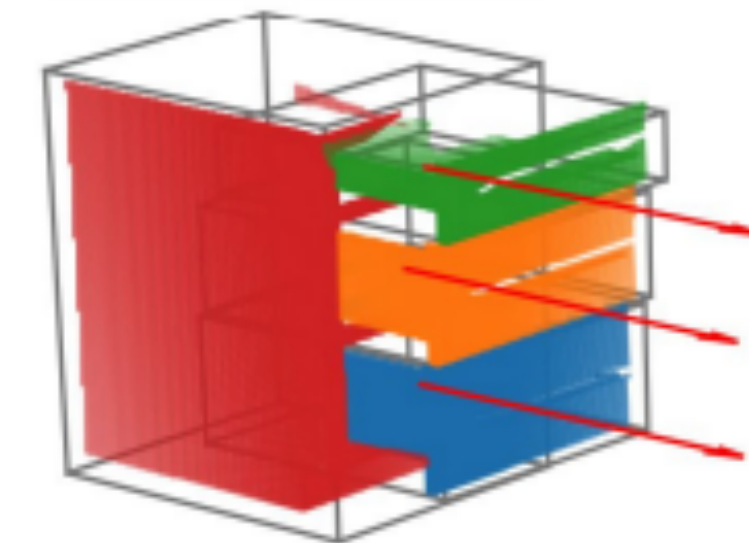
**3 Parts**

**Model 2**



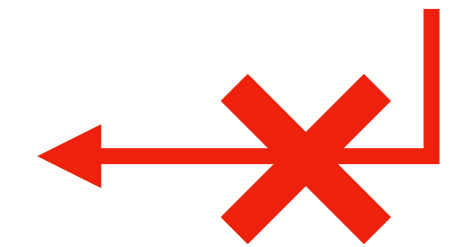
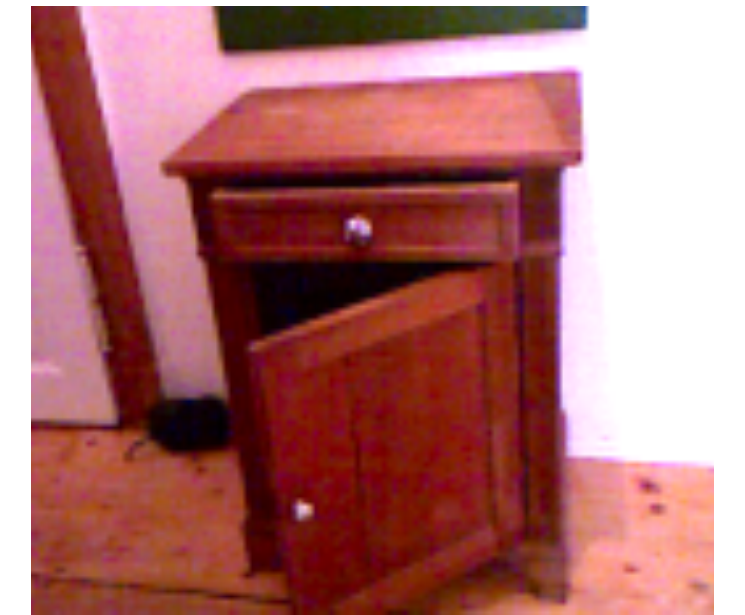
**2 Parts**

**Model 3**



**4 Parts**

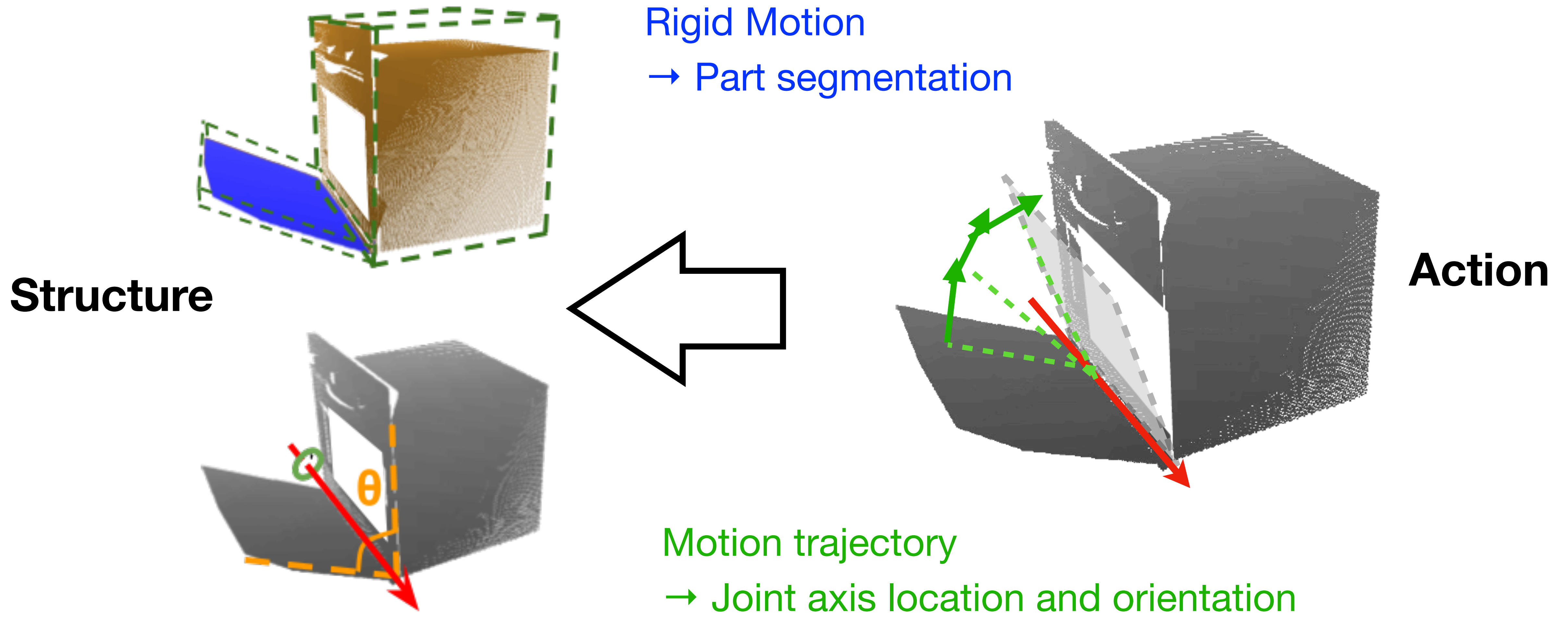
...



\* Category-Level Articulated Object Pose Estimation. Li et al, CVPR 2019



# Action → Structure





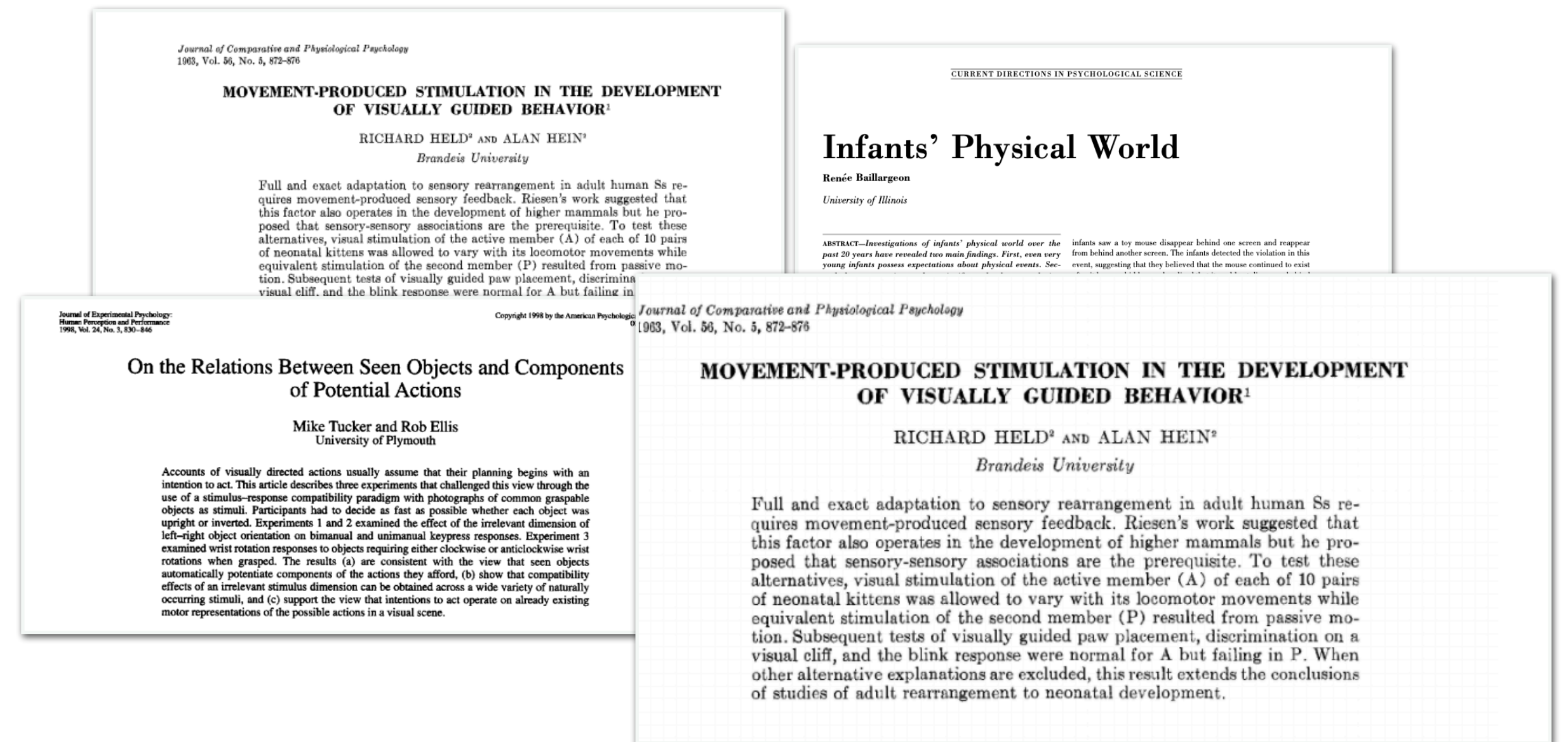
# Action → Structure



Video Credit: Samir Y. Gadre

How do kids learn to understand an articulated objects?

Interacting + observing!

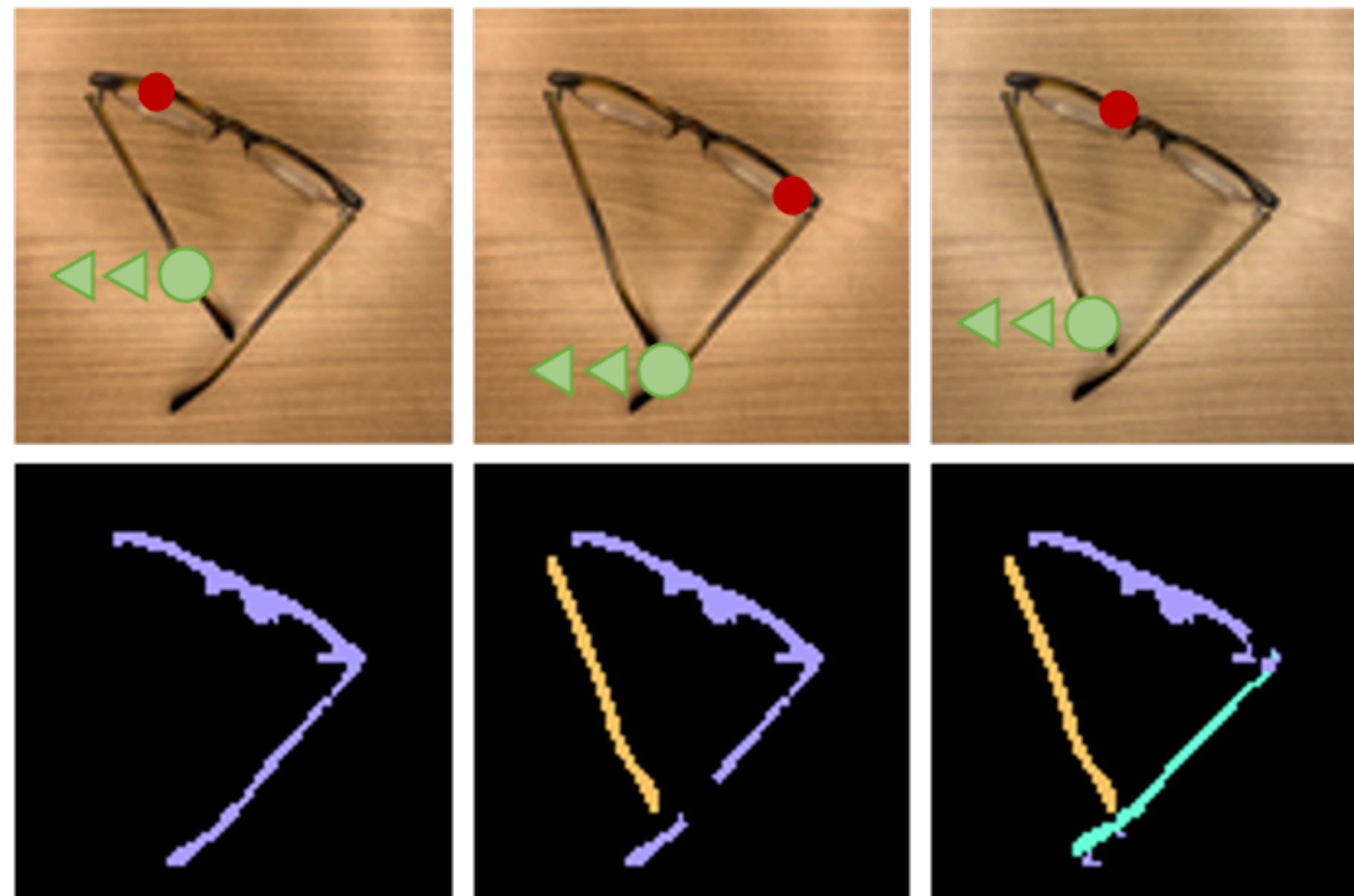


Can we allow our robot to do the same?

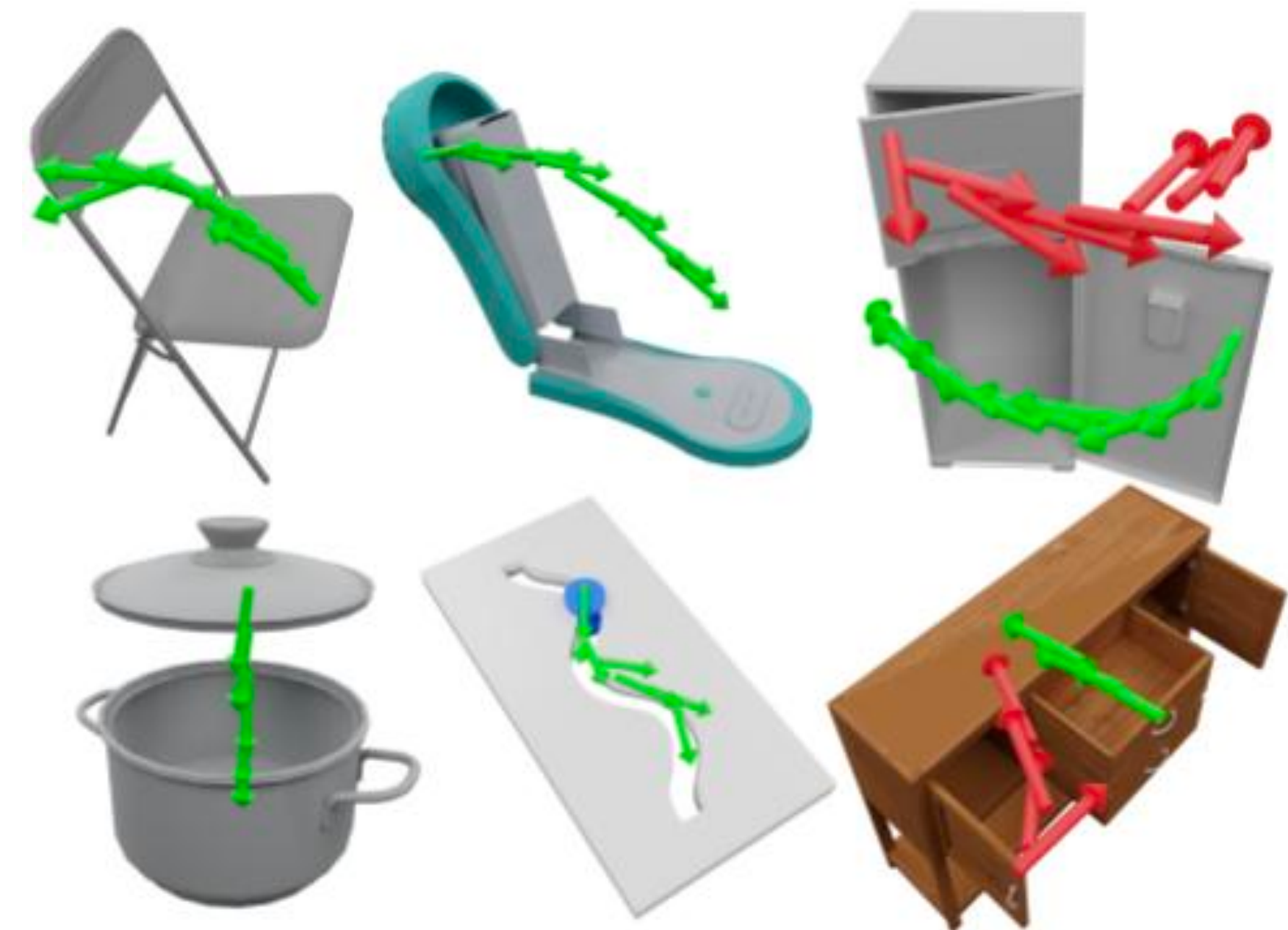


# Structure from Action

**Act the Part**



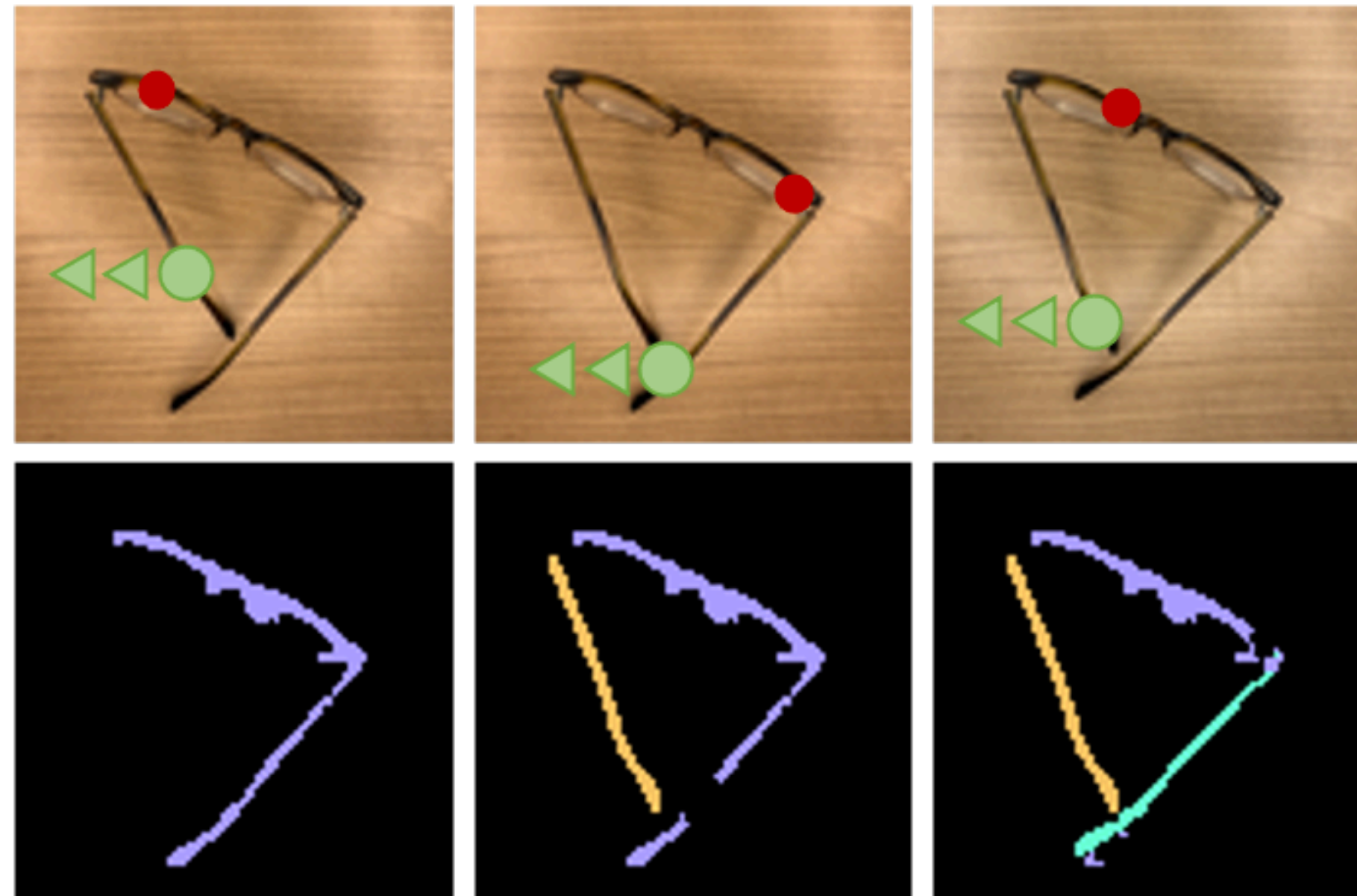
**Universal Manipulation Policy**



Leverage active interaction as way to discover Articulated Object Structure

# Structure from Action

## Act the Part



Act the Part: Learning to Interact to  
Discover Articulated Object Structure  
**ICCV 2021**

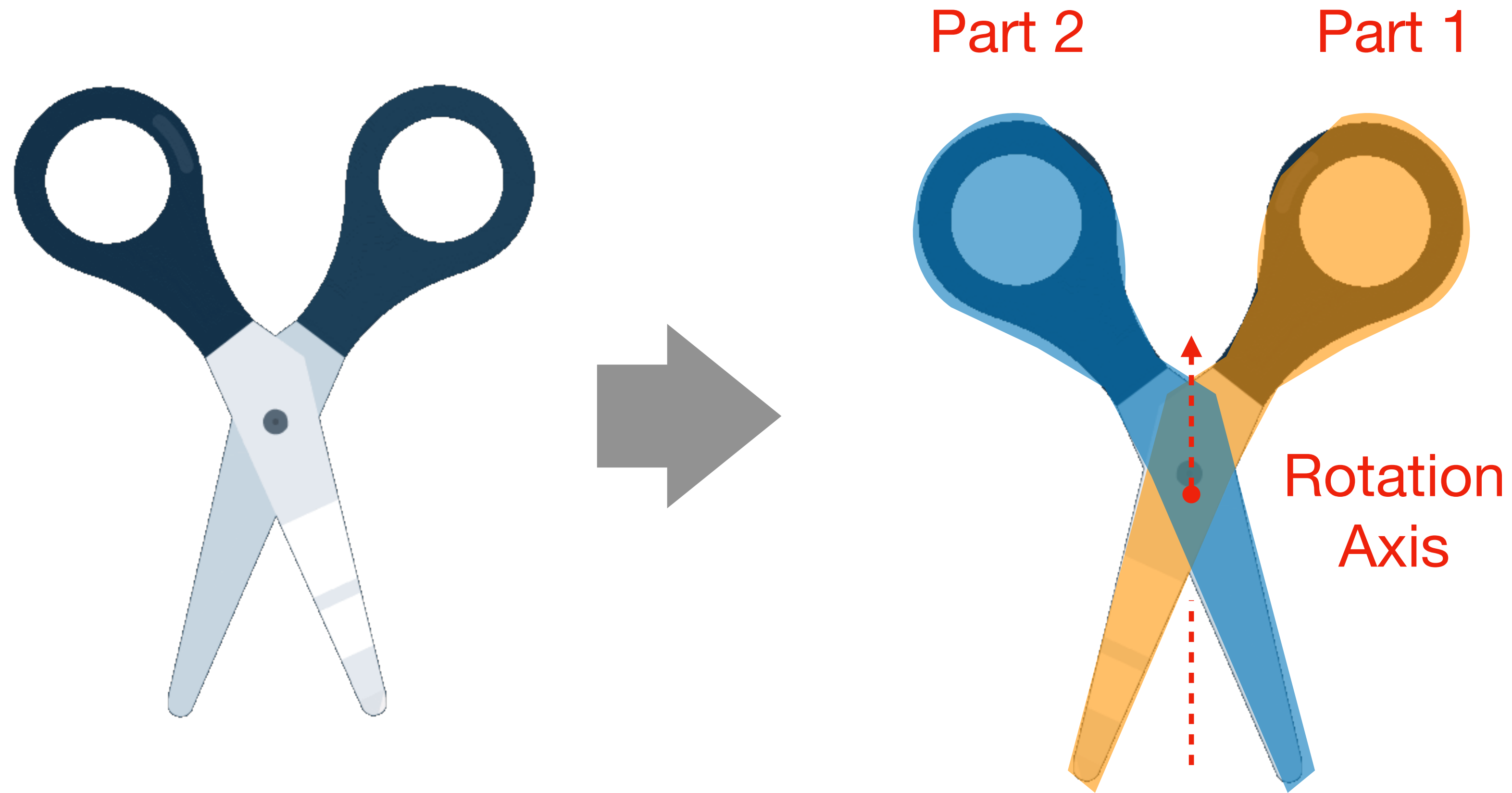
Samir Y. Gadre, Kiana Ehsani, Shuran Song





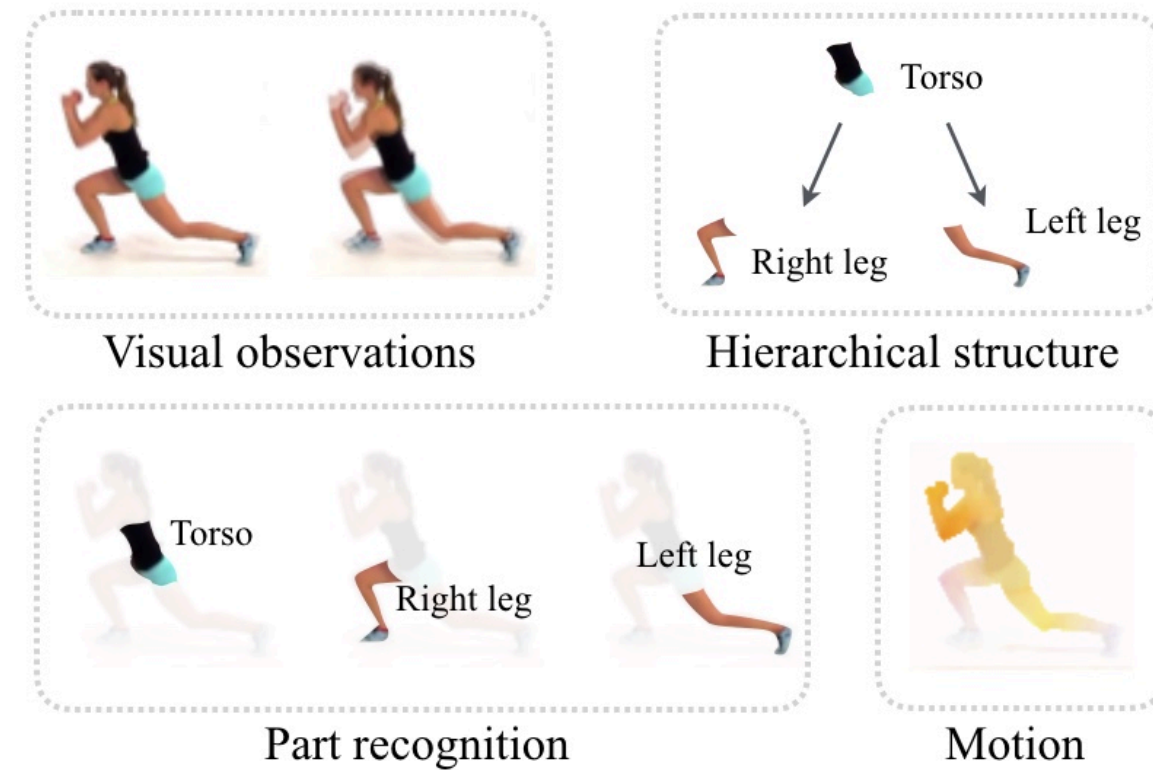
# Goal

**Discover the kinematic structure for articulated objects through interaction**

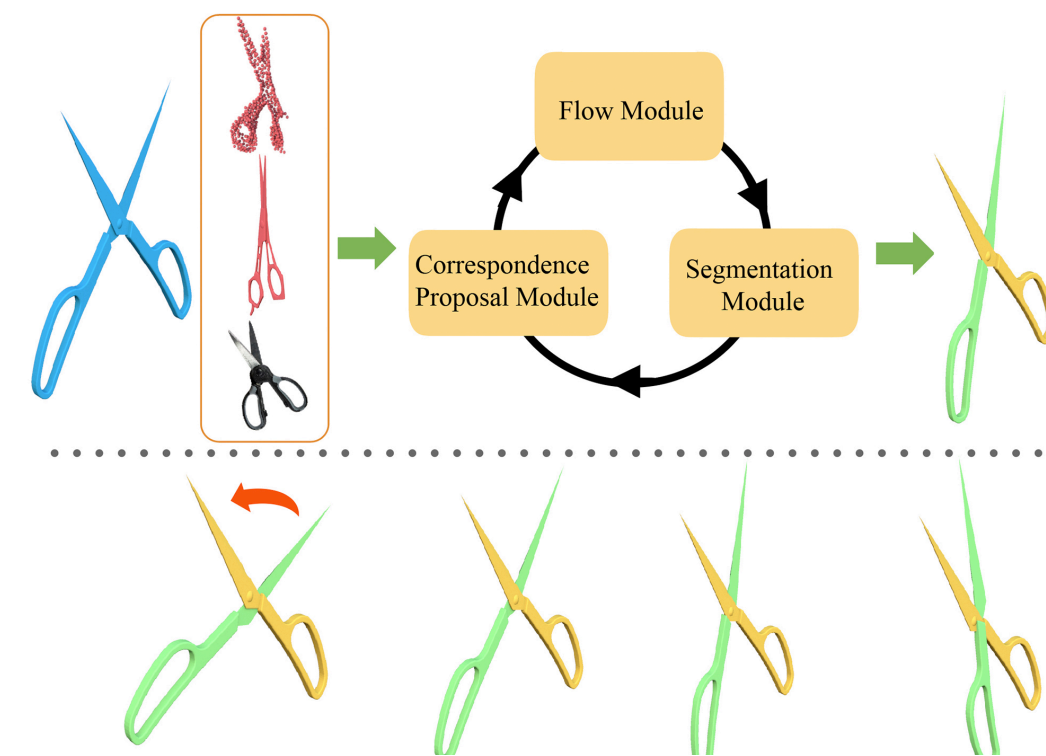


*Object part discovery and segmentation with unknown objects*

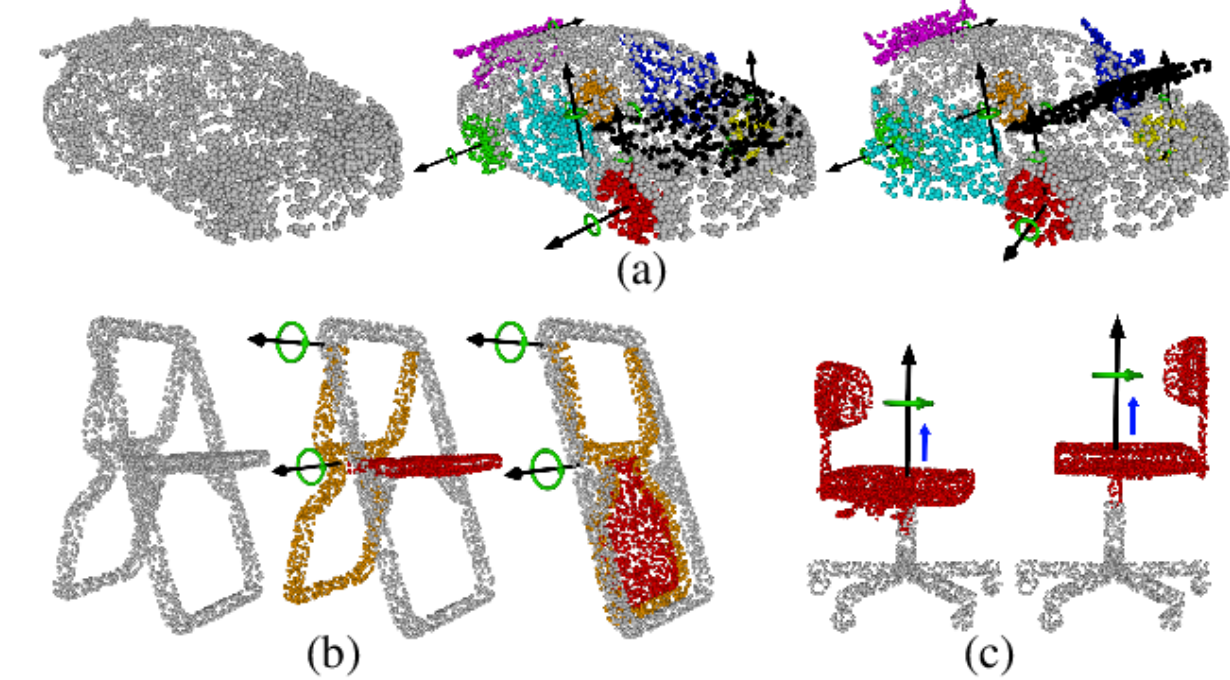
# Unsupervised Part Discovery from Videos



Xu et al.



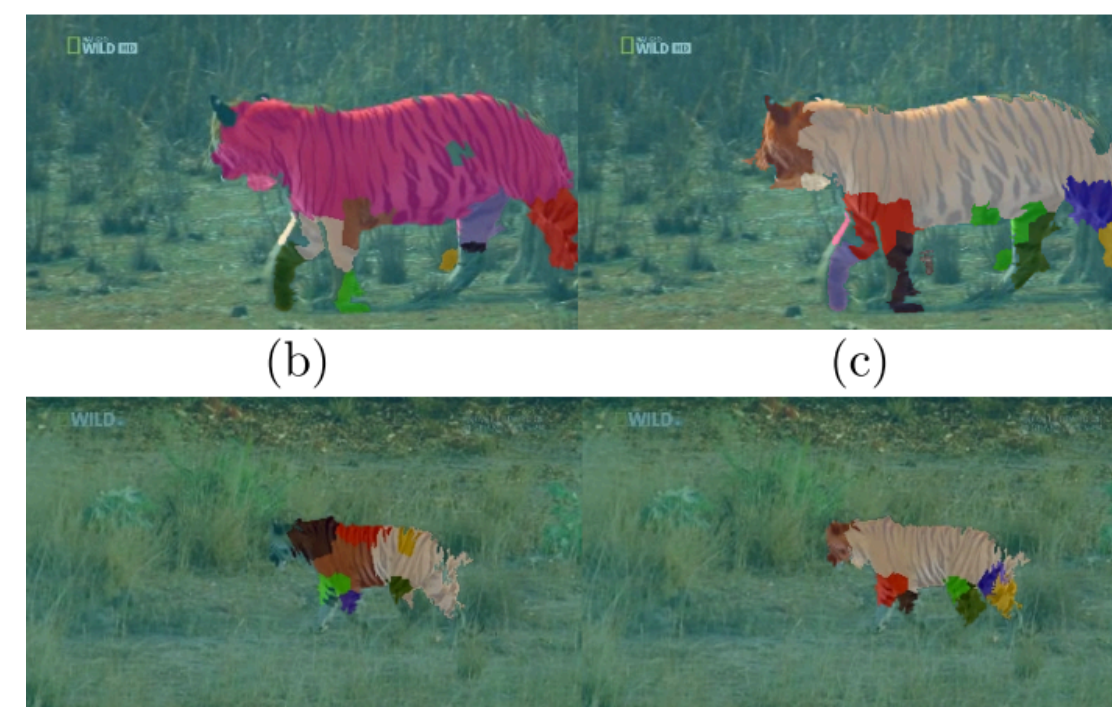
Yi et al.



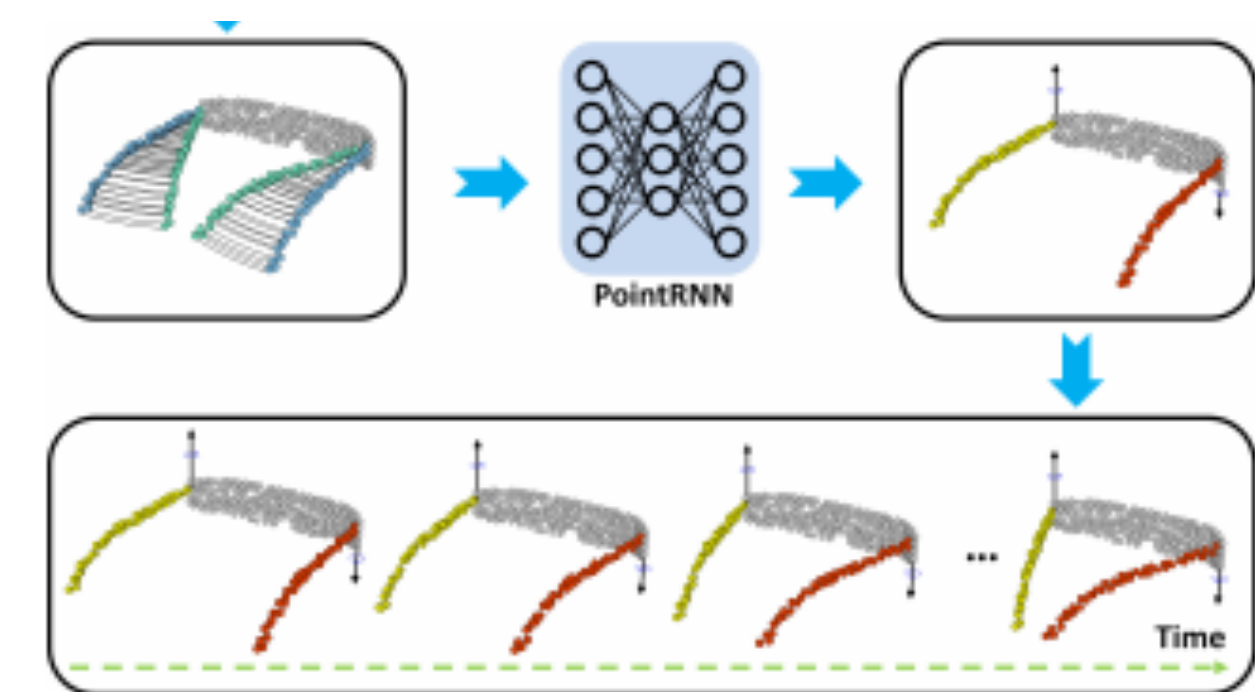
Wang et al.



Tokmakov et al.



Pero et al.



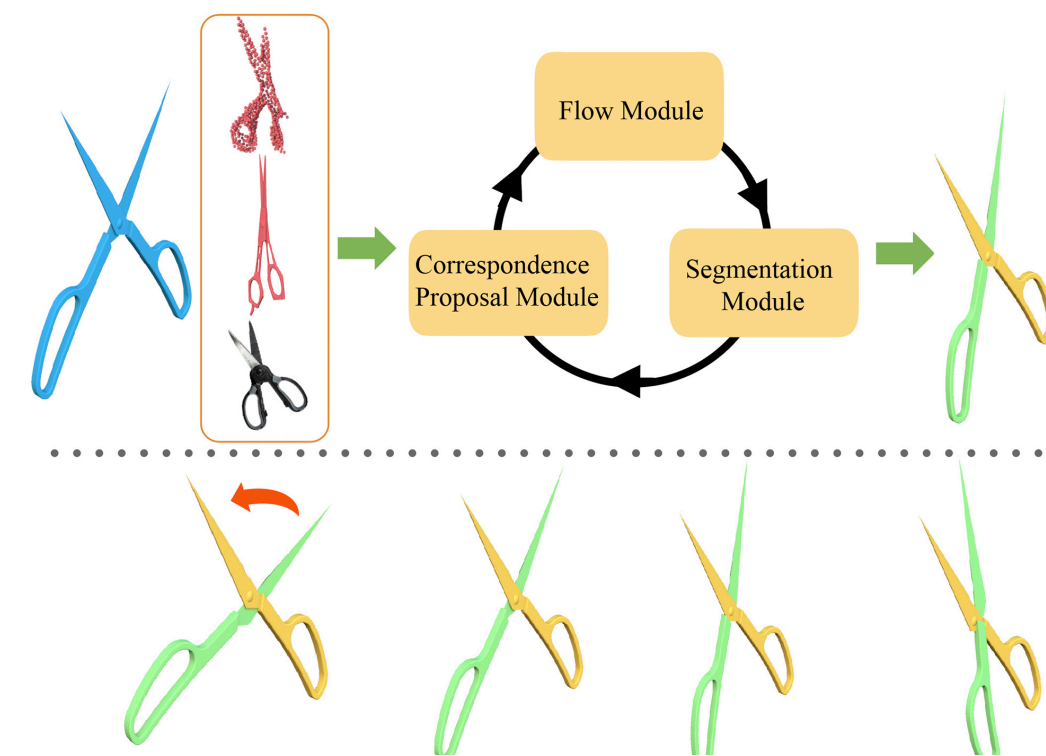
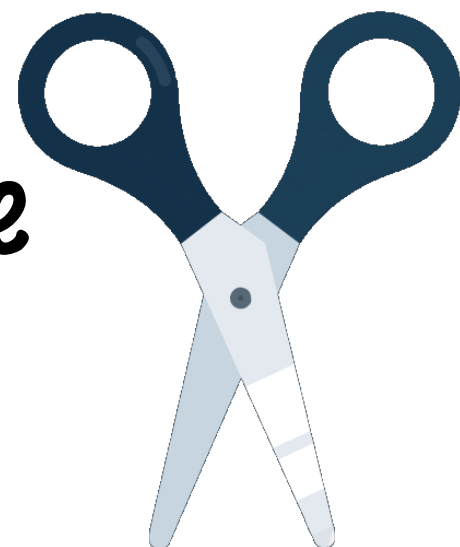
Shi et al.

Prior works: Motion Consistency for Part Segmentation

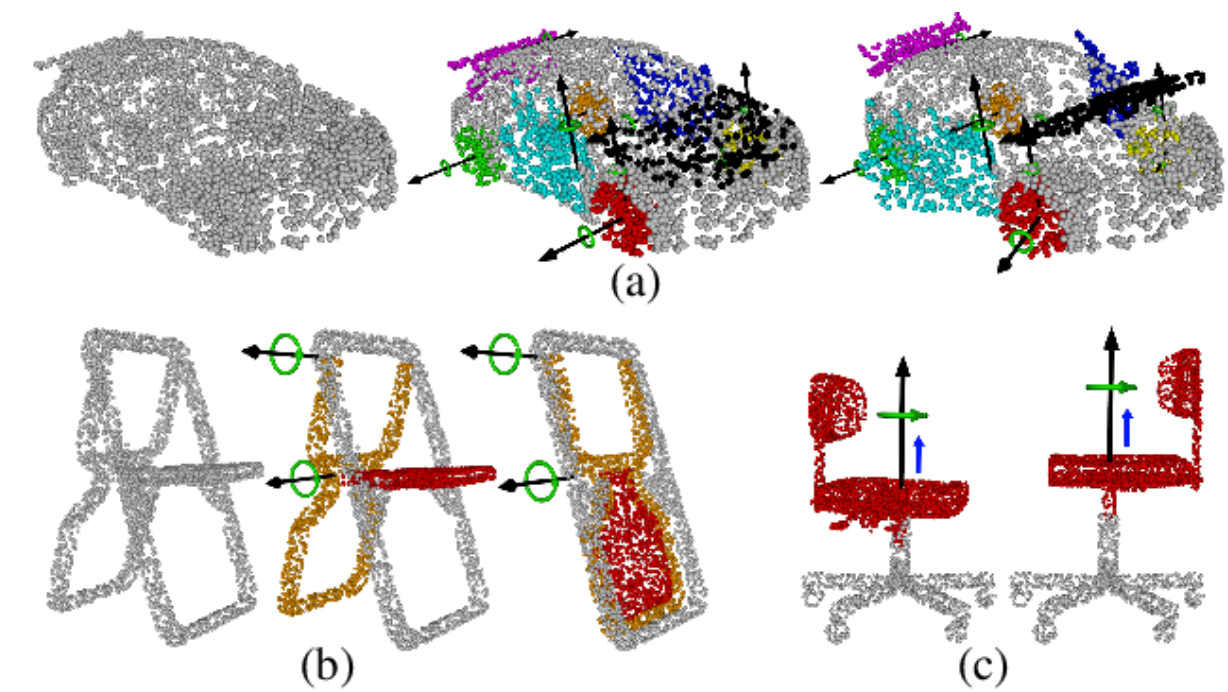


# Object Part Segmentations from Video

Who causes all these motions ?

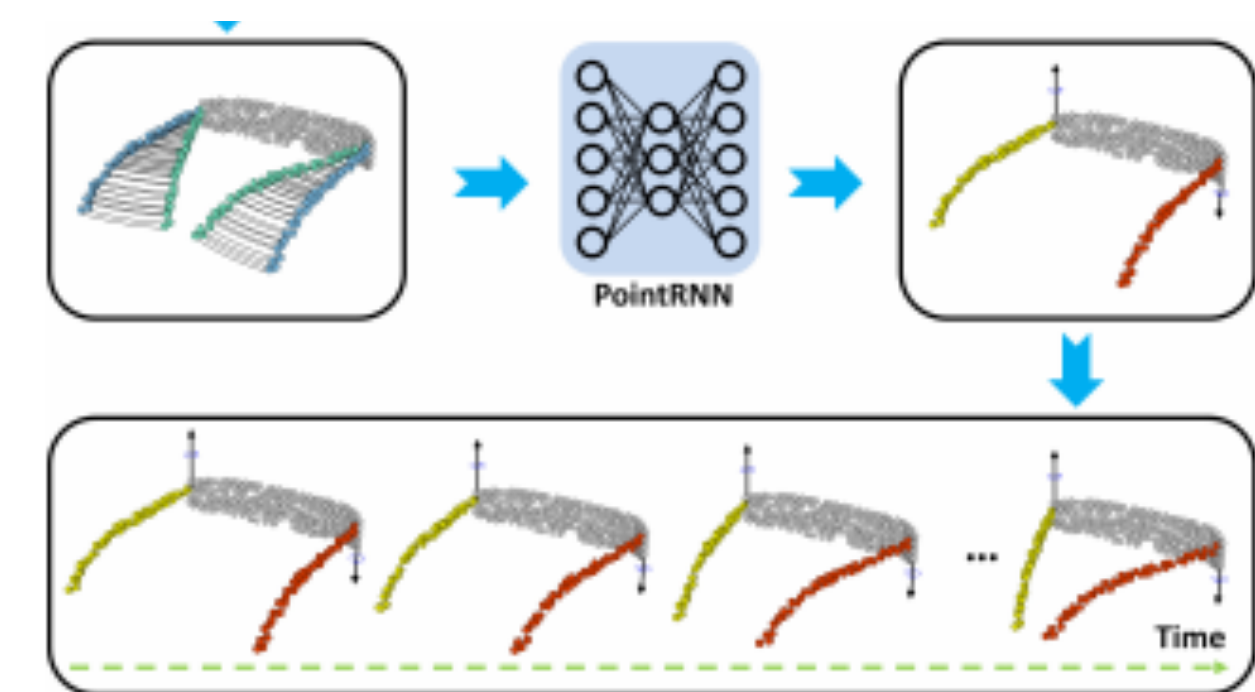


Yi et al.



Wang et al.

Most of Objects don't move by themselves ...



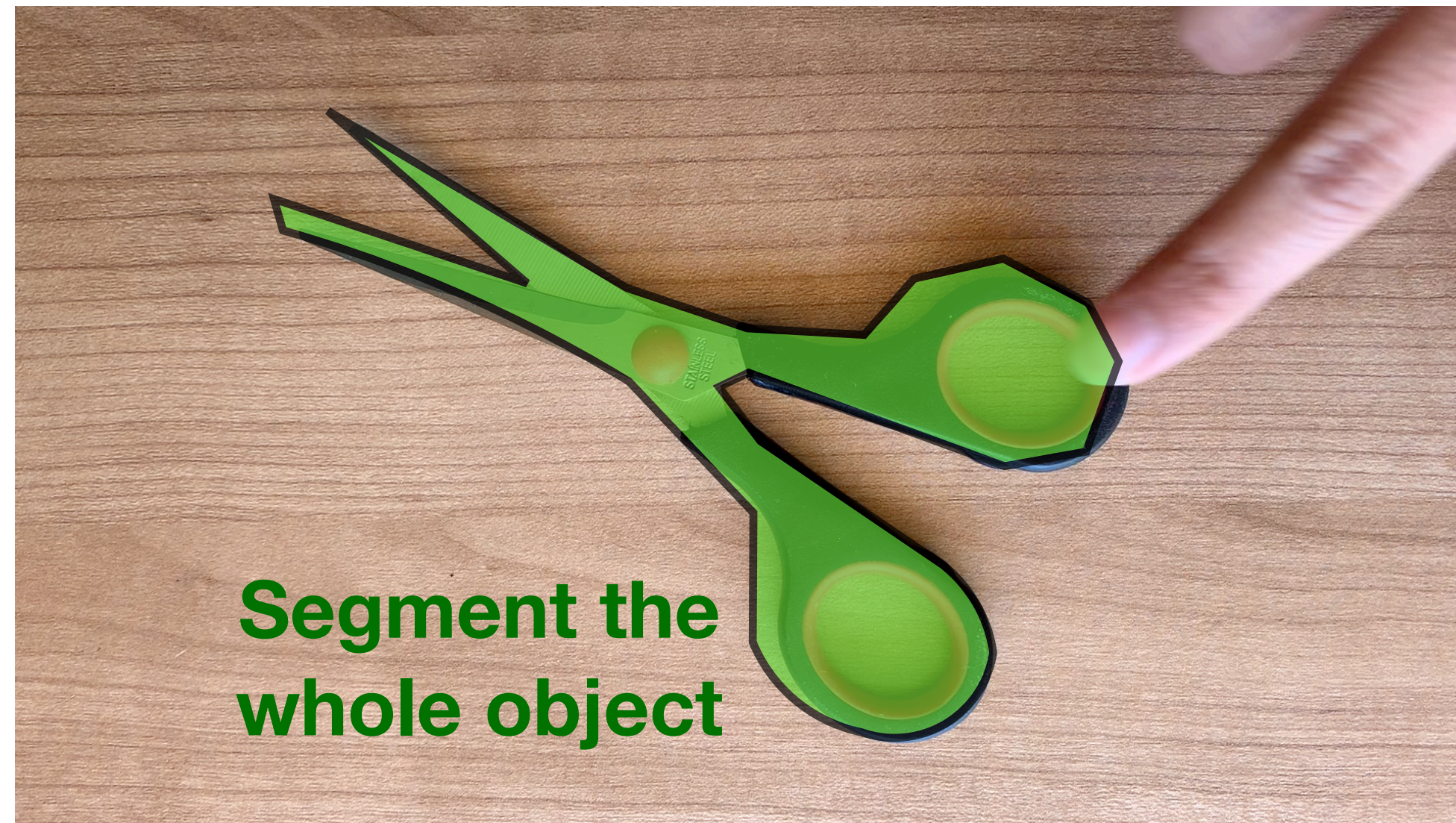
Shi et al.

Prior works: Motion Consistency for Part Segmentation

# **Not all Motions are Informative**



# Not all Motions are Informative



Some interactions don't give insight about articulation.



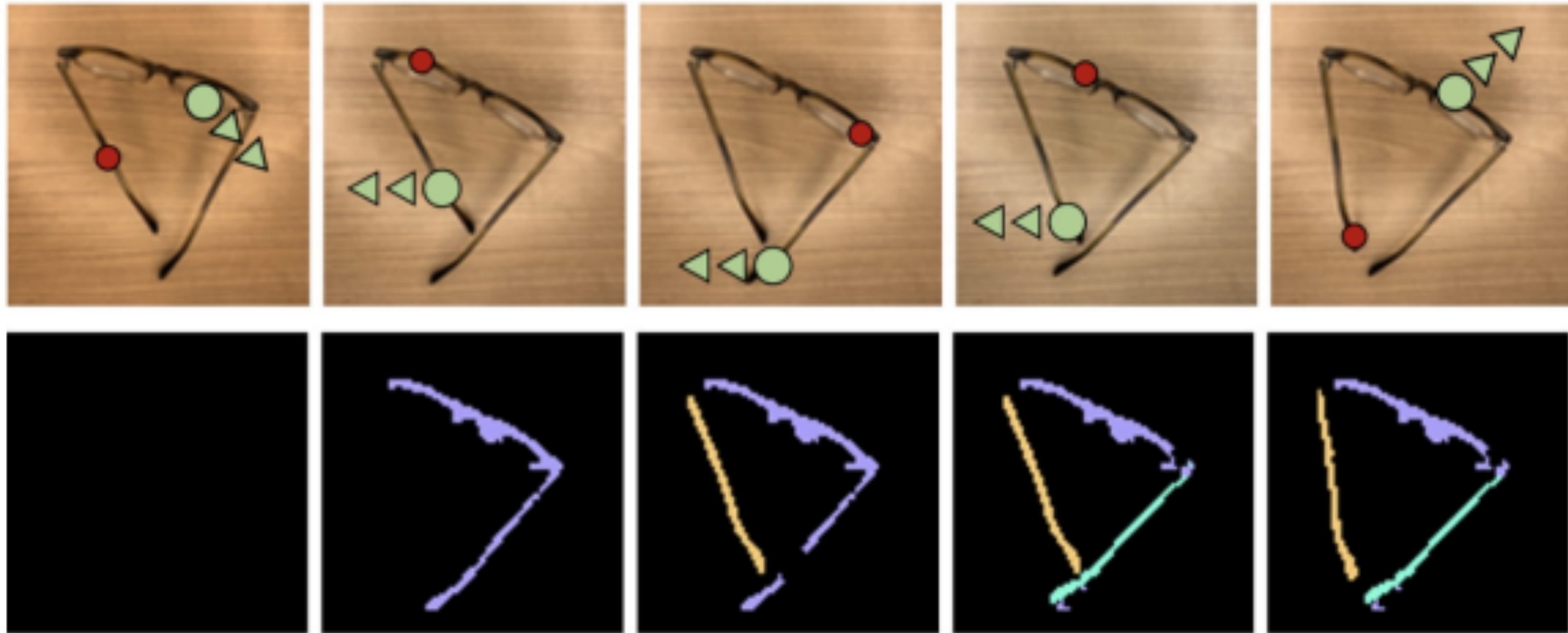
Single step interaction is not enough for multi-link object

*“Informative or not” really depends on “what the system already know (i.e., belief)”*



# Act the Part

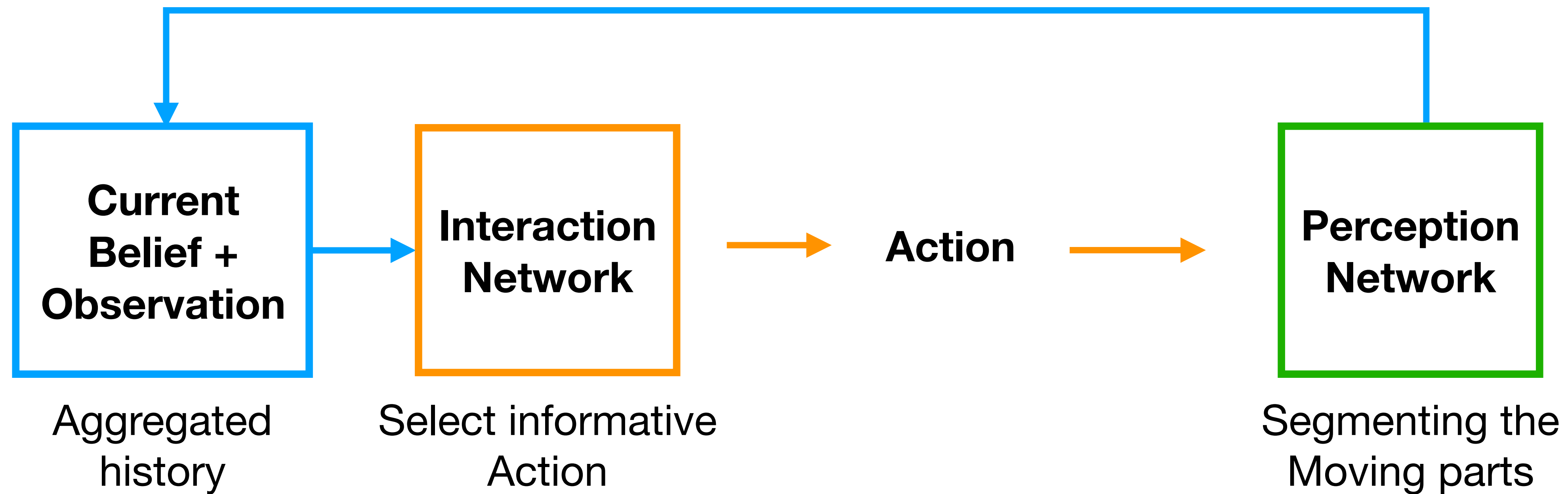
Learning to Interact to Discover Articulated Object Structure



*Generate a sequence of hold and push actions that would results in informative motion.*

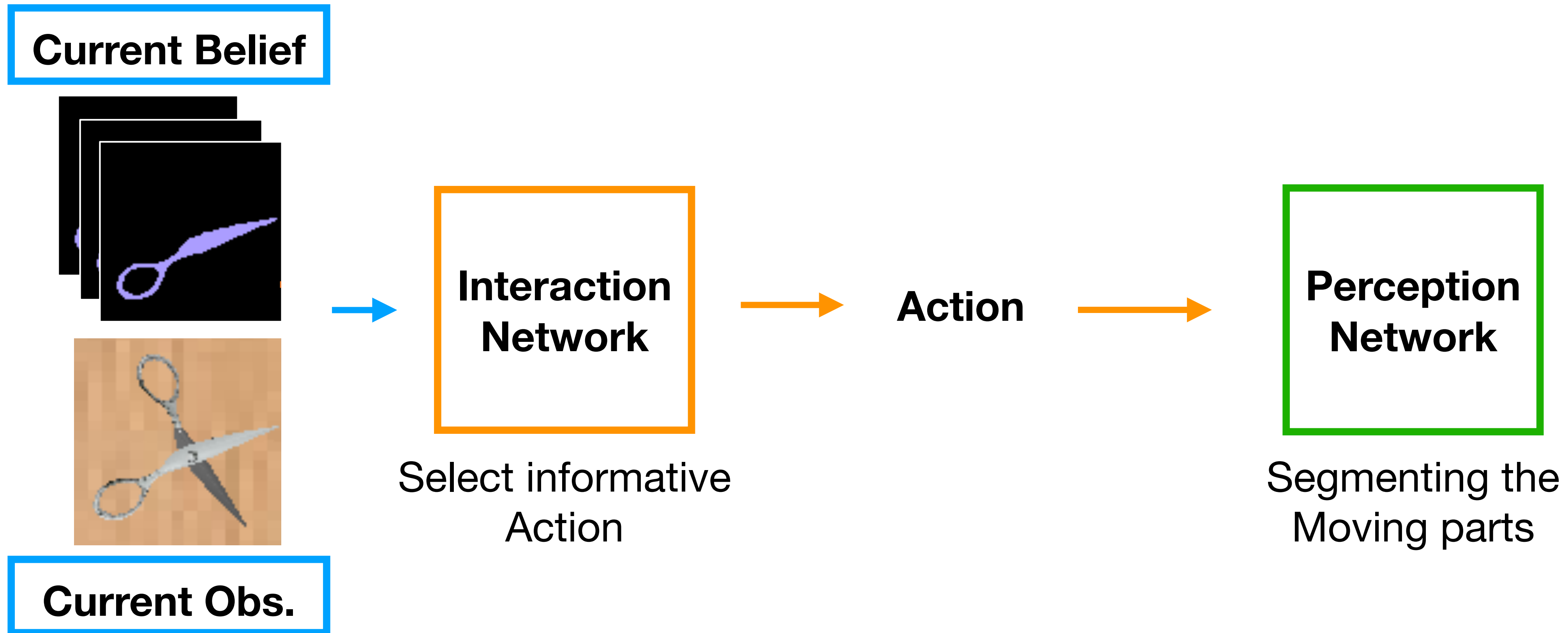
# Act the Part

Key Idea: couple action selection and motion segmentation.



# Act the Part

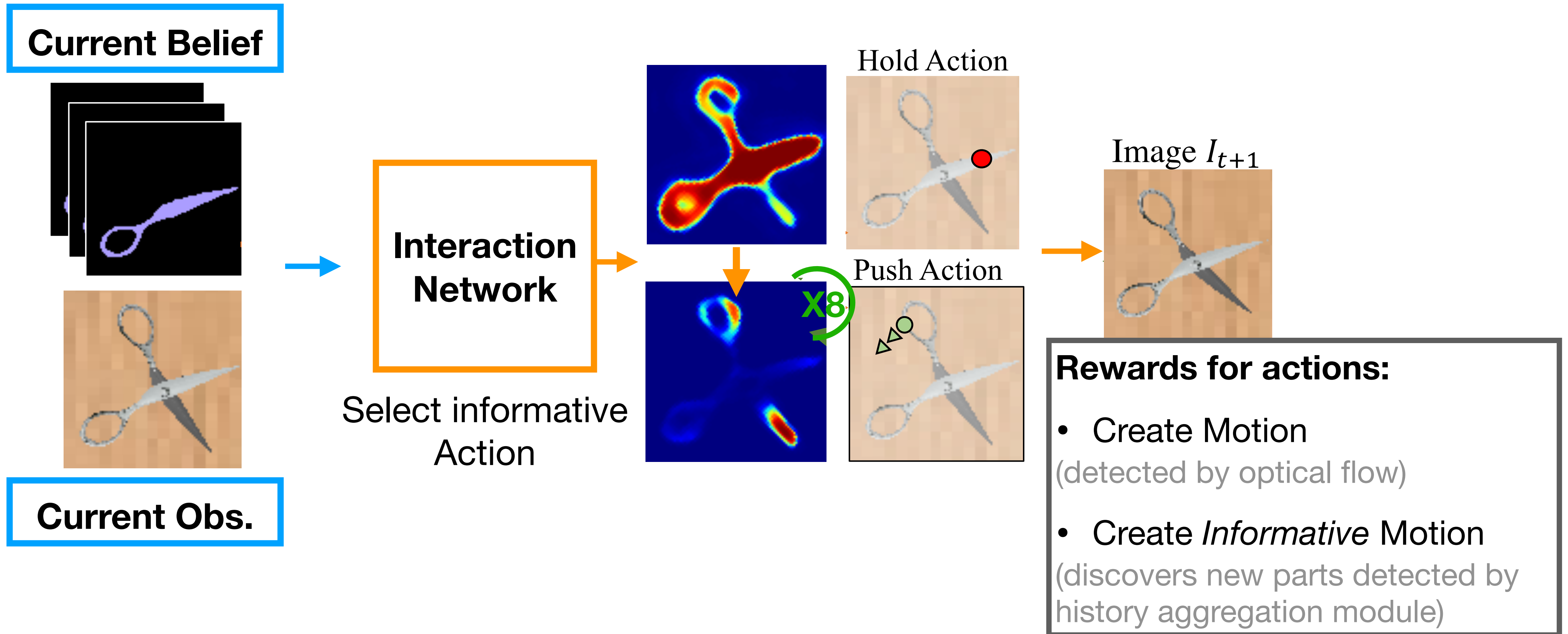
Key Idea: couple action selection and motion segmentation.





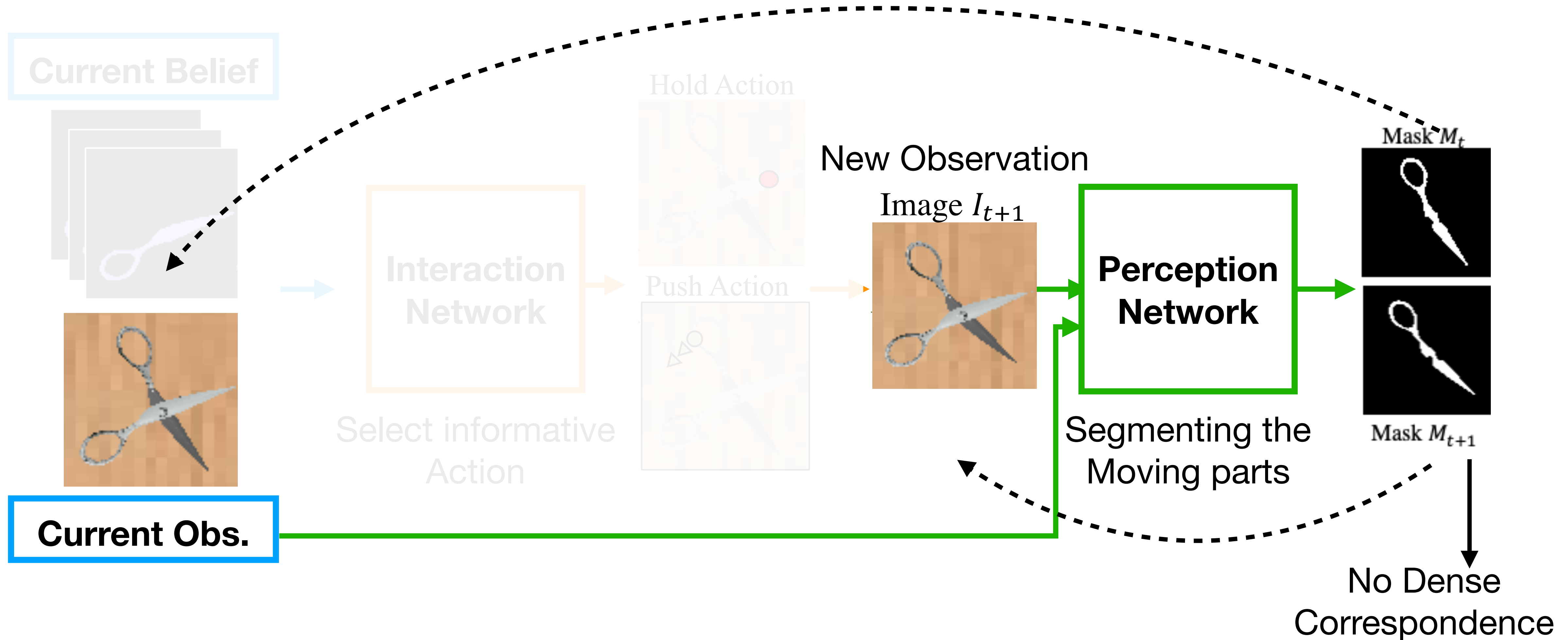
# Act the Part

Key Idea: couple action selection and motion segmentation.



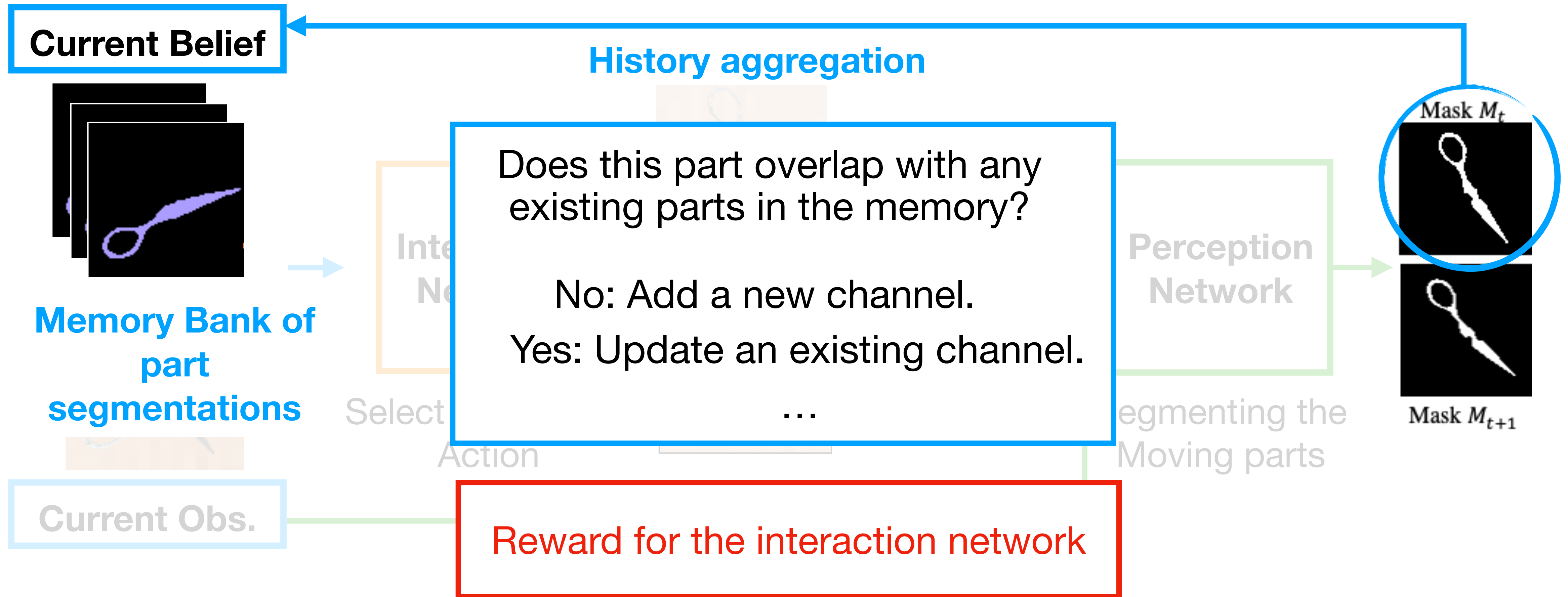
# Act the Part

Key Idea: couple action selection and motion segmentation.



# Act the Part

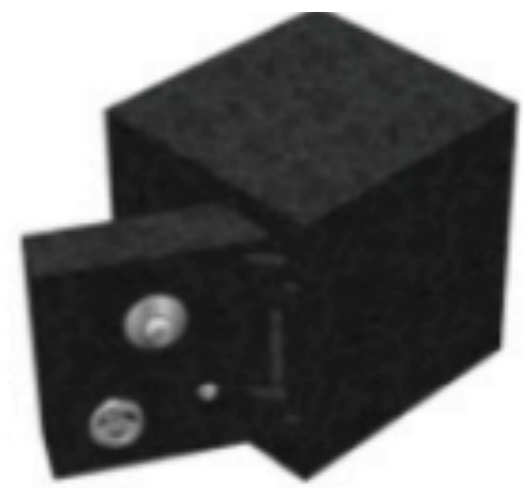
Key Idea: couple action selection and motion segmentation.





# Training Testing Objects

## Training



Safe



Scissors



USB



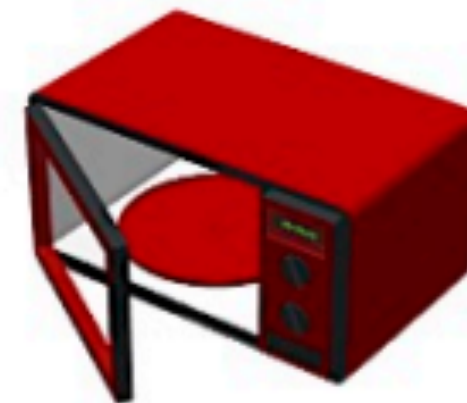
Knife

## Testing Categories

### *Simulation*



Lighter



Microwave



Eyeglasses



Pliers



Multilink

### *Realword*



Earbuds



Eyeglasses



Tea Bag



Keys  
(Two Link)



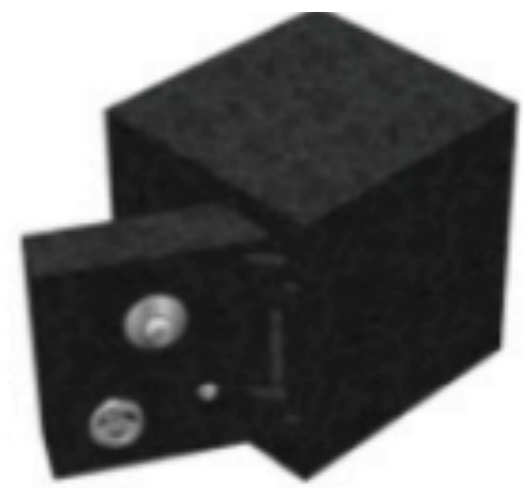
Keys  
(Three Link)

*One Model for All Object Categories*



# Training Testing Objects

## Training



Safe



Scissors



USB



Knife

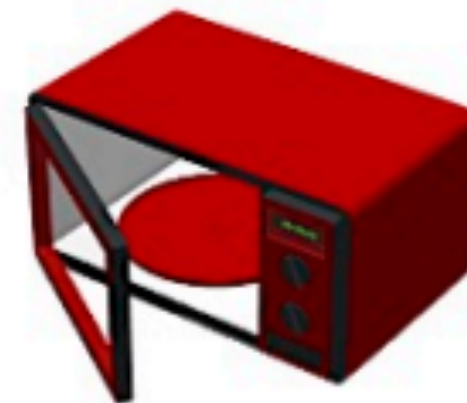
*Two links*

## Testing Categories

### Simulation



Lighter



Microwave



Eyeglasses



Pliers



Multilink

### Realword



Earbuds



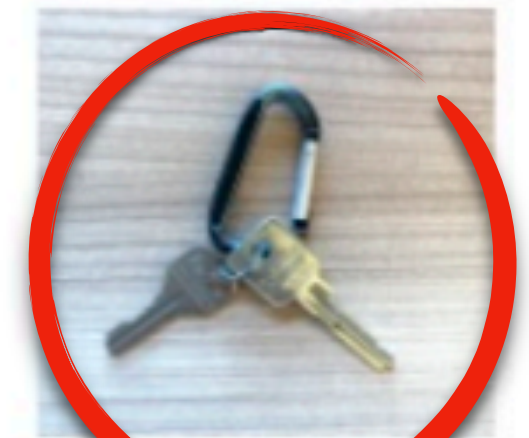
Eyeglasses



Tea Bag



Keys  
(Two Link)

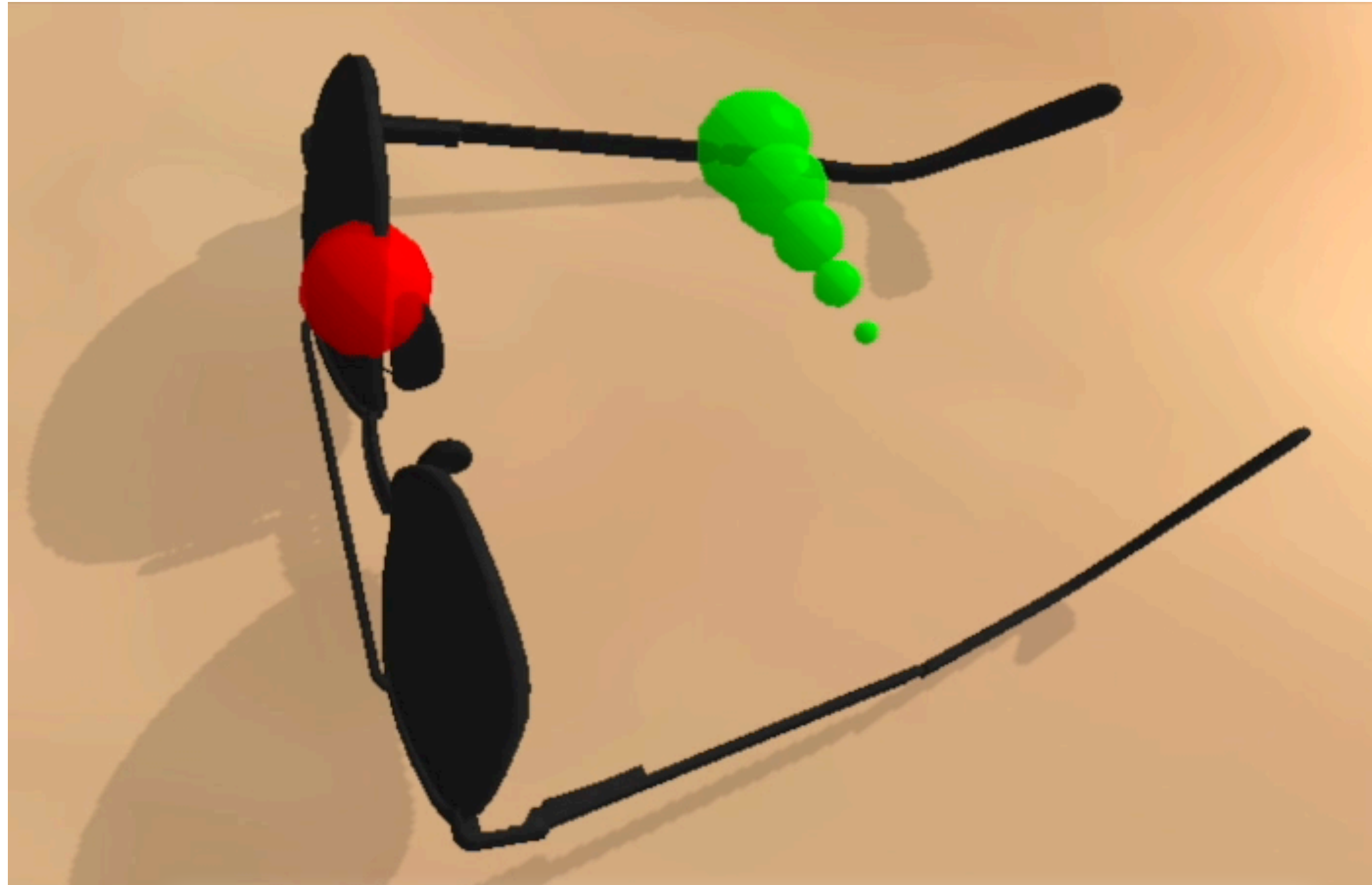


Keys  
(Three Link)

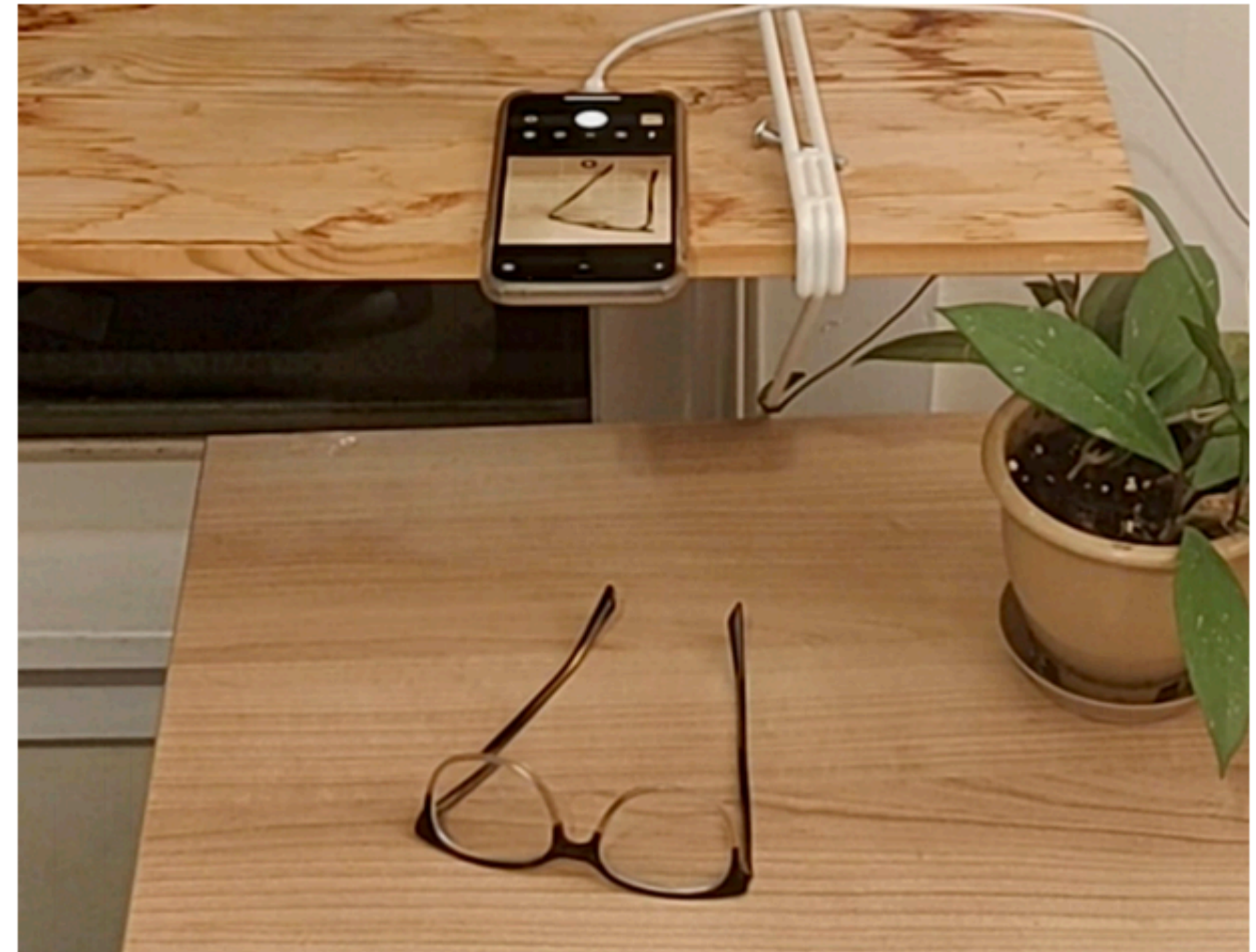
*Two or three links with different kinematic structure*



# Experiment Results



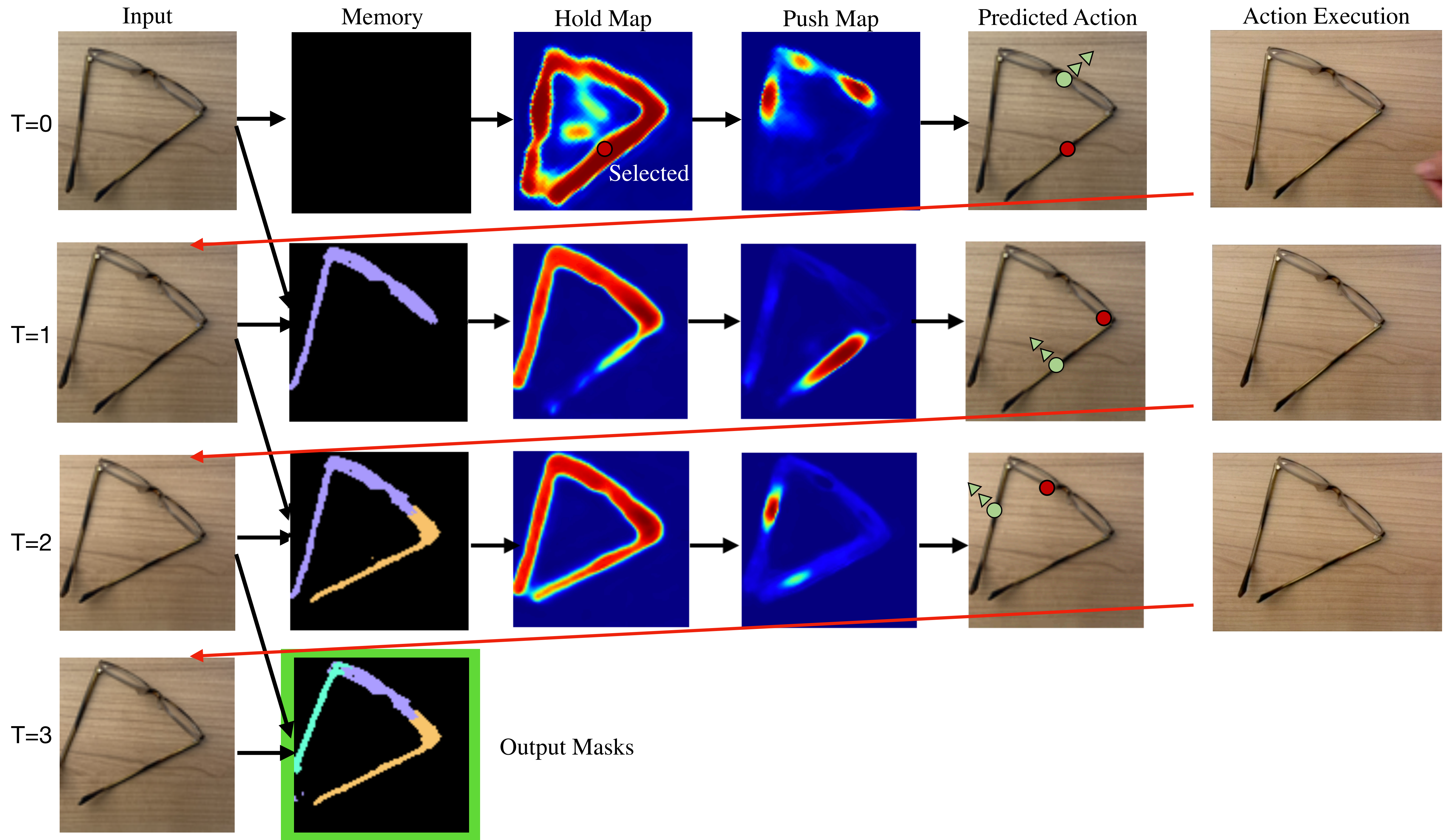
**Train in Simulation**



**Test on Real World Objects\***

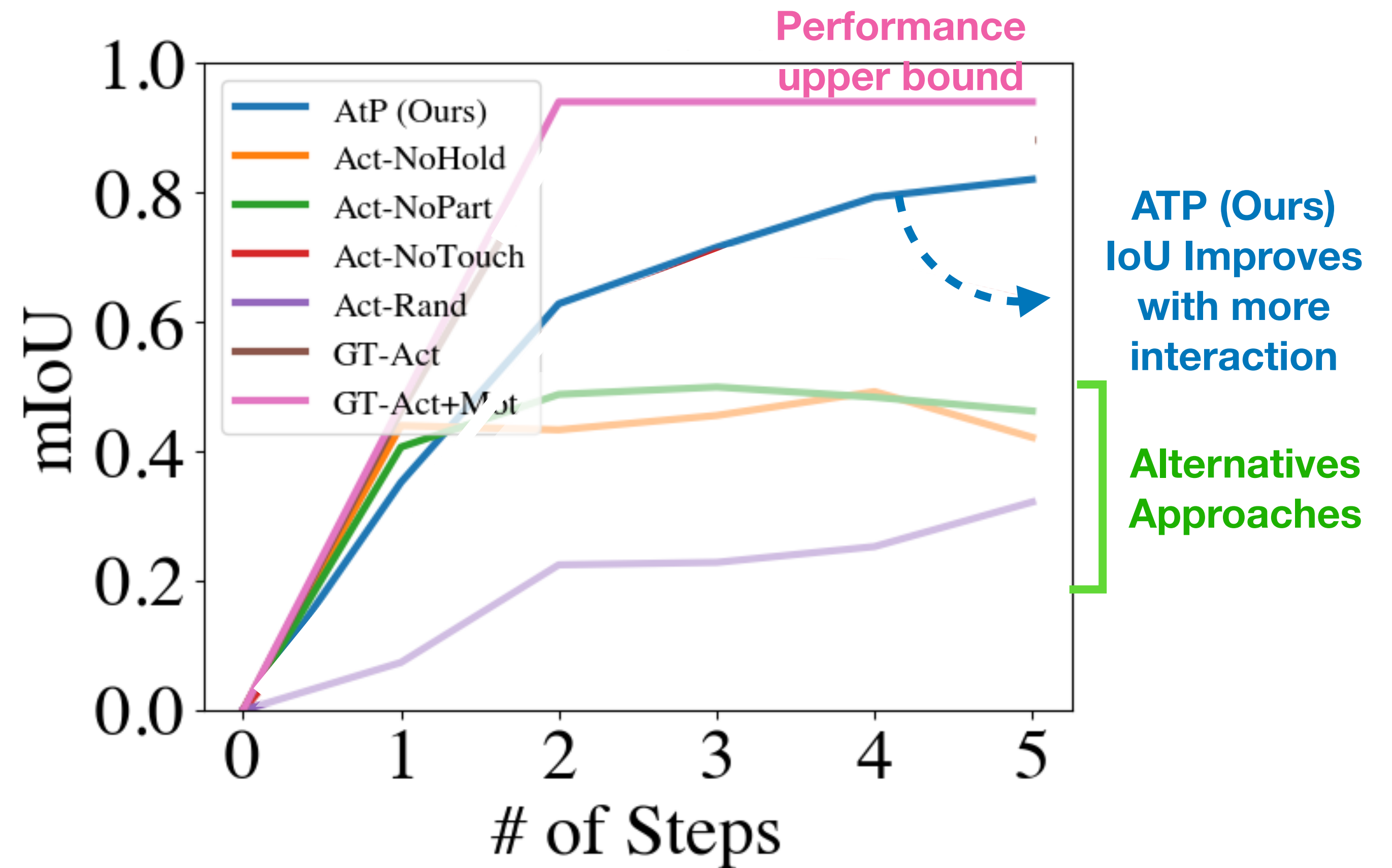
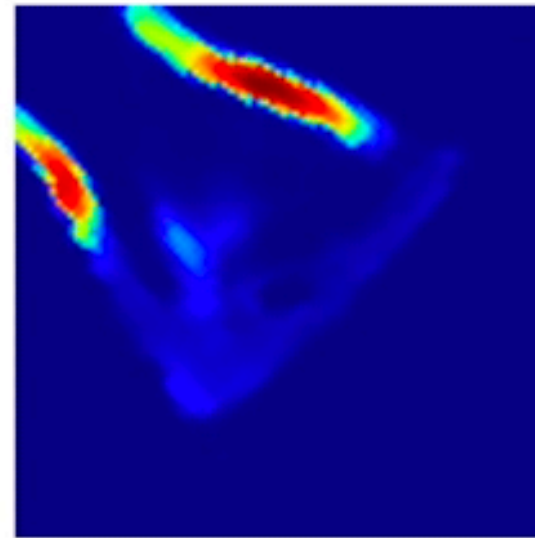
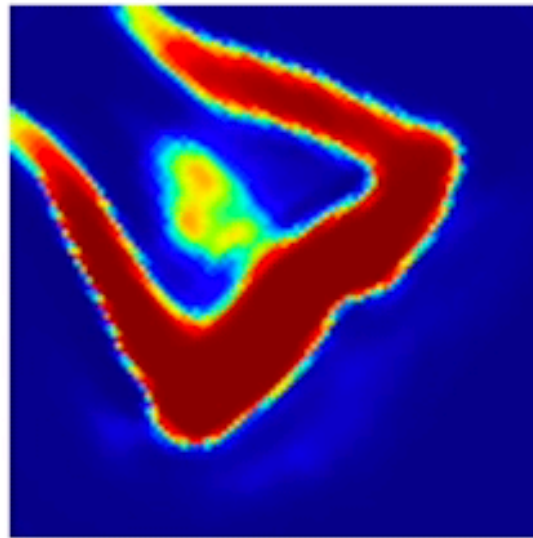
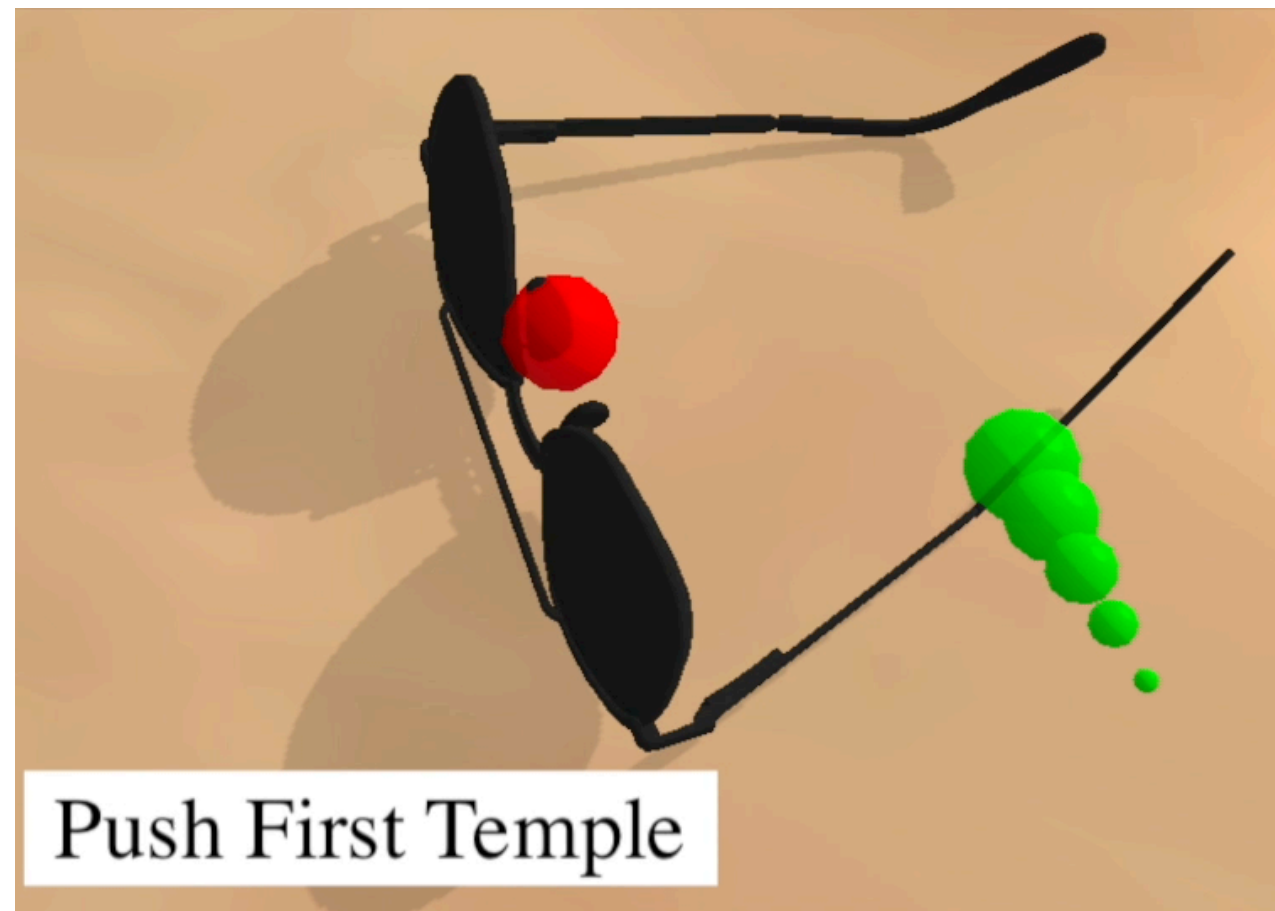
*\* working from home setup :)*





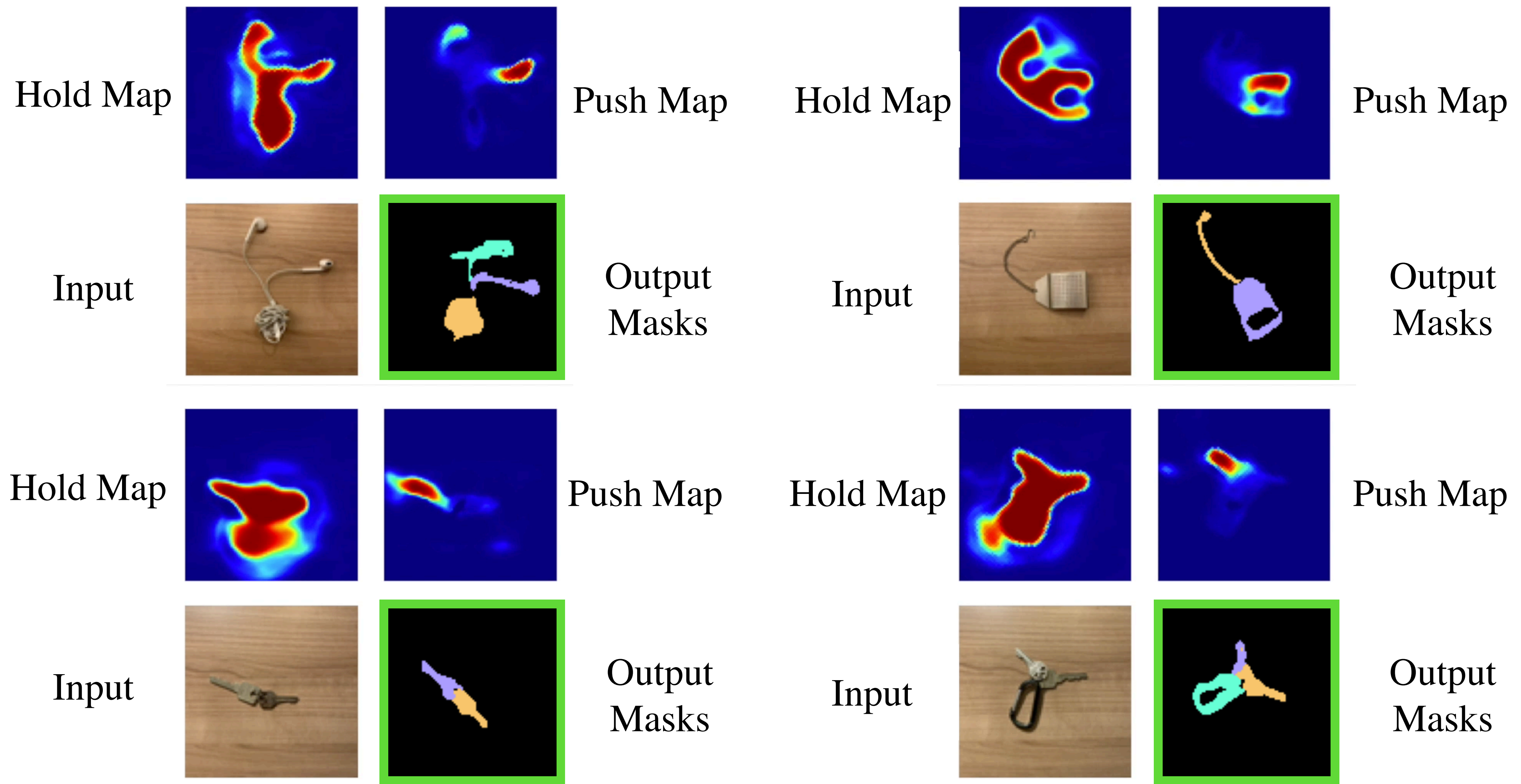


# Qualitative Results in Simulation





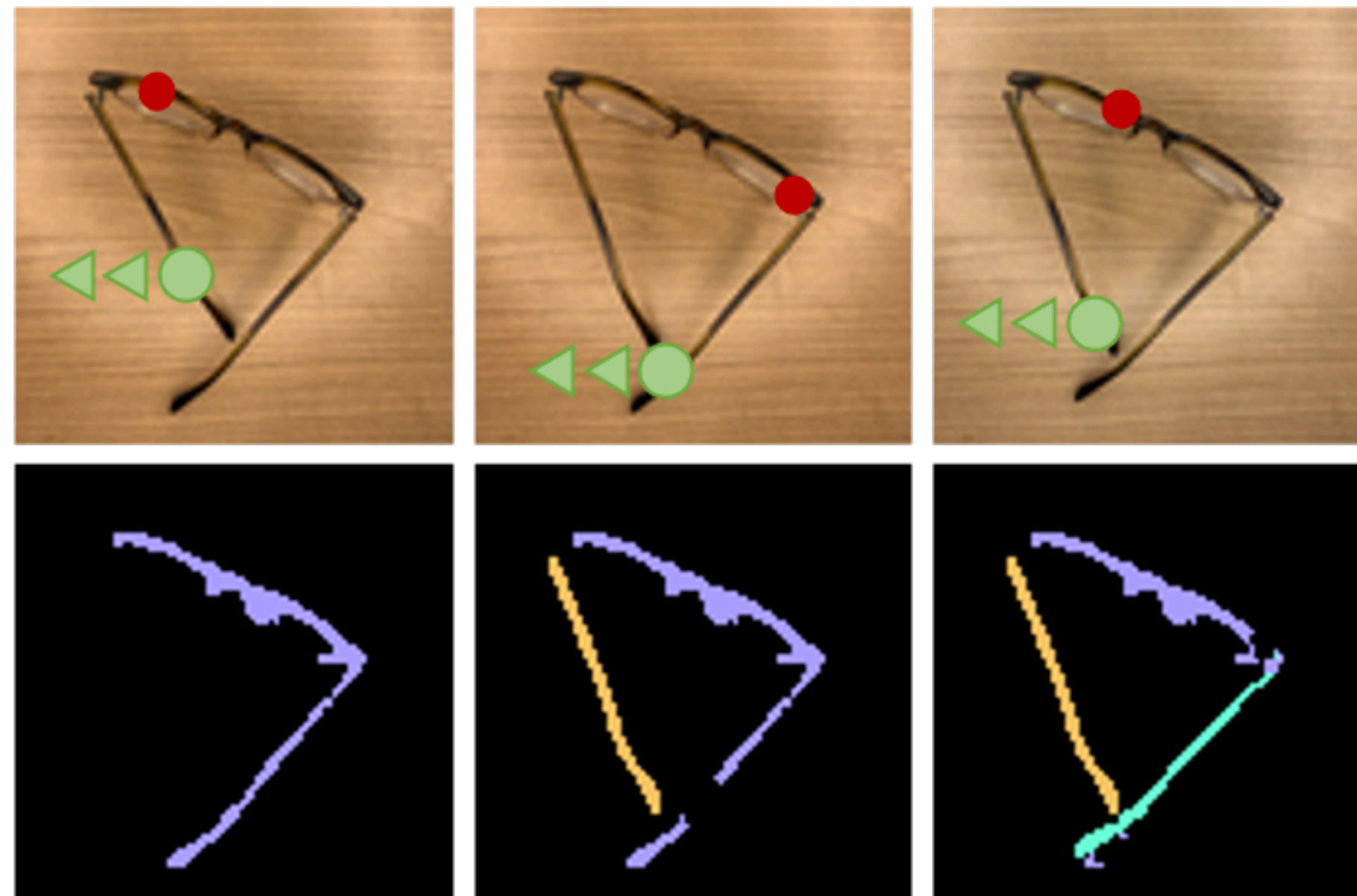
# Results: Real World Generalization



Webpage: <https://atp.cs.columbia.edu/> with online Demo!

# Structure from Action

## Act the Part



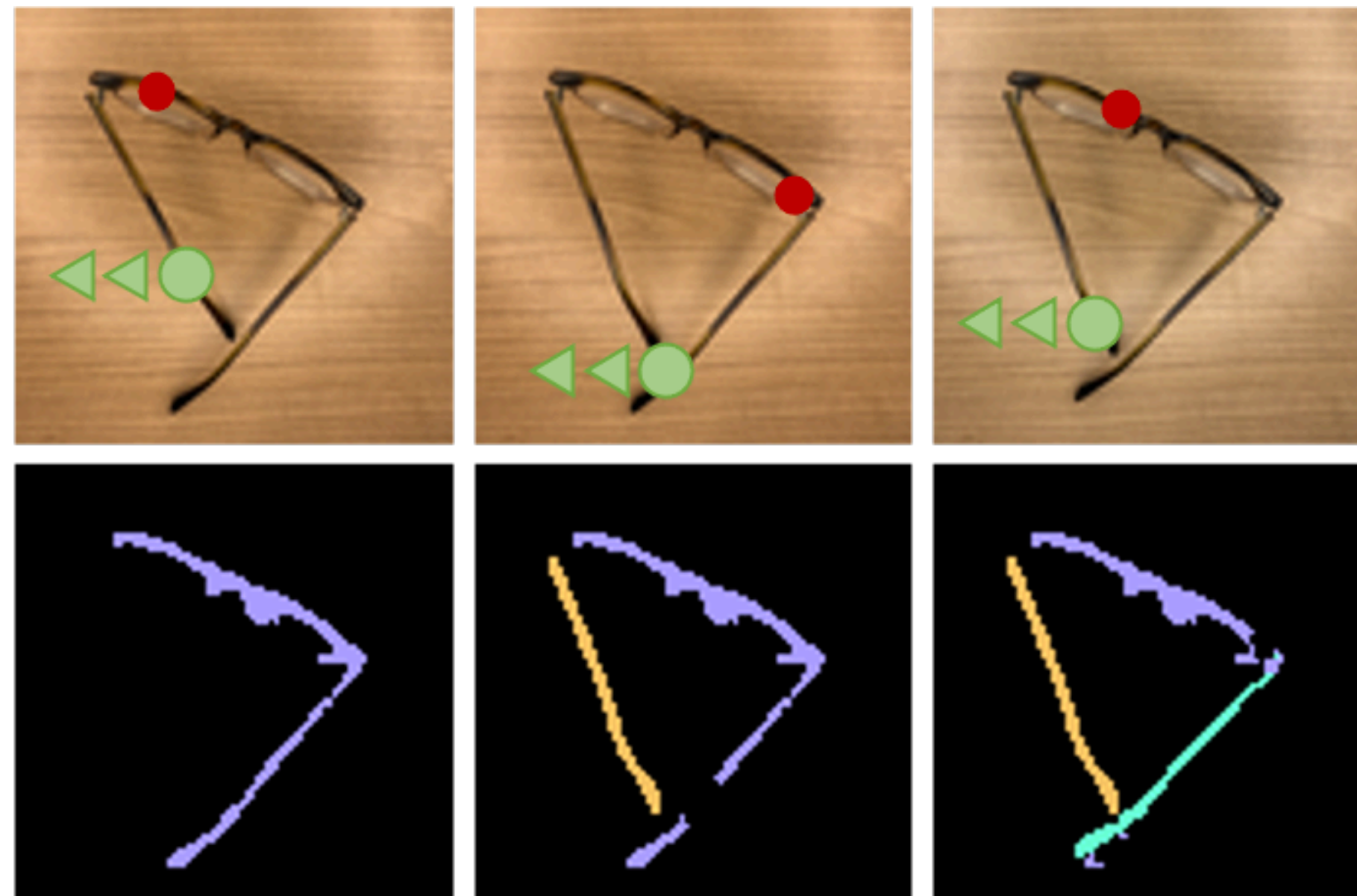
## Summary:

- AtP is able to learn effective interaction strategies isolating and discovering parts
- It is able to generalize to novel categories of objects with unknown and unseen number of links

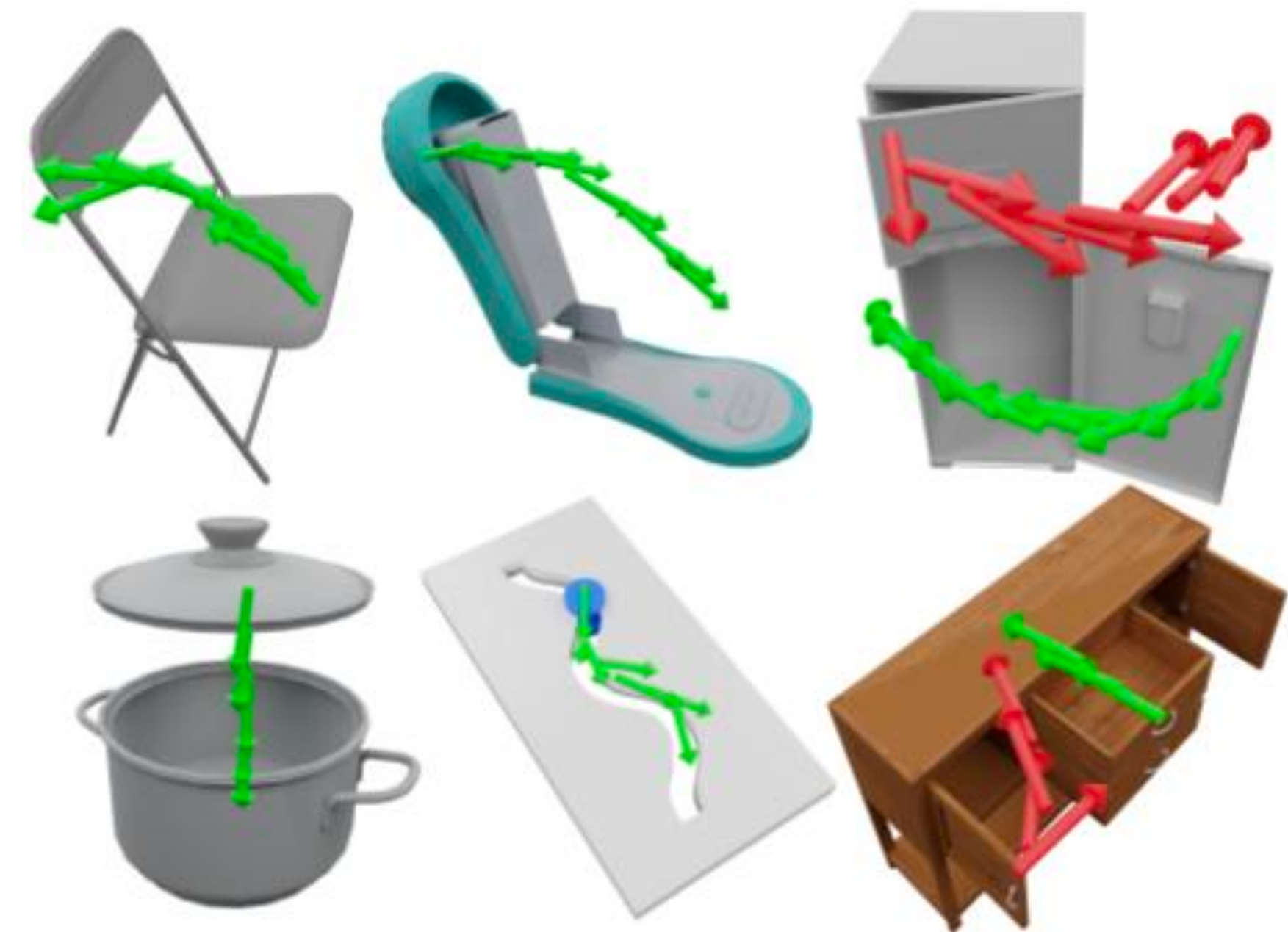


# Structure from Action

## Act the Part



**Simple 2D Action**  
**with discrete action direction**

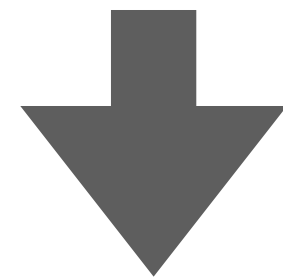


**x Closed-loop, 3D Actions**  
**x Goal-conditioned Manipulation Tasks**

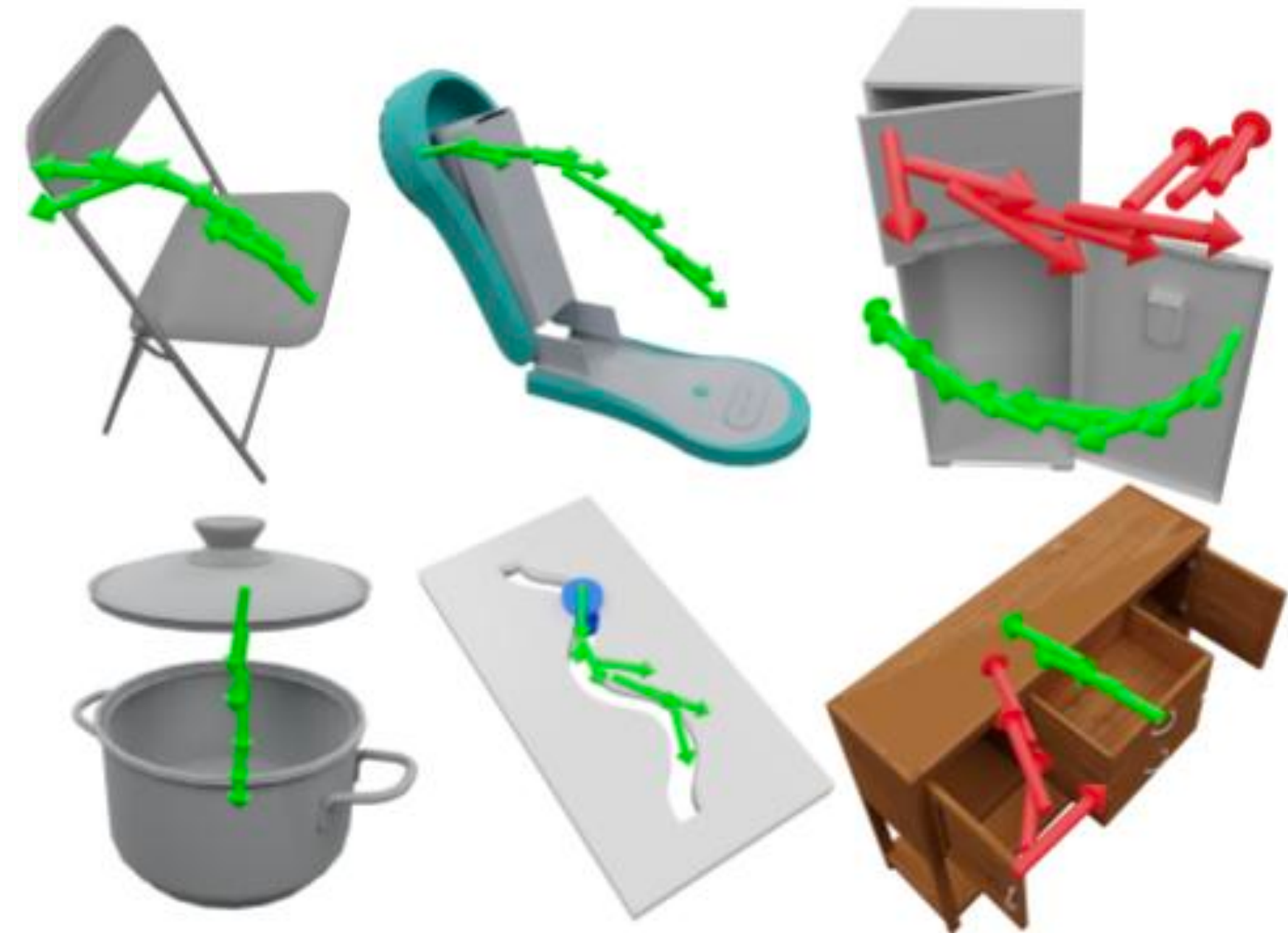
# Structure from Action

## Upgrade the Interaction Policy

- General Action Representation (Continuous Action in  $SE(3)$ )
- Closed-loop Action Sequence
- Goal-Conditioned Manipulation Tasks



**Improve the Understanding of the Object Structure**

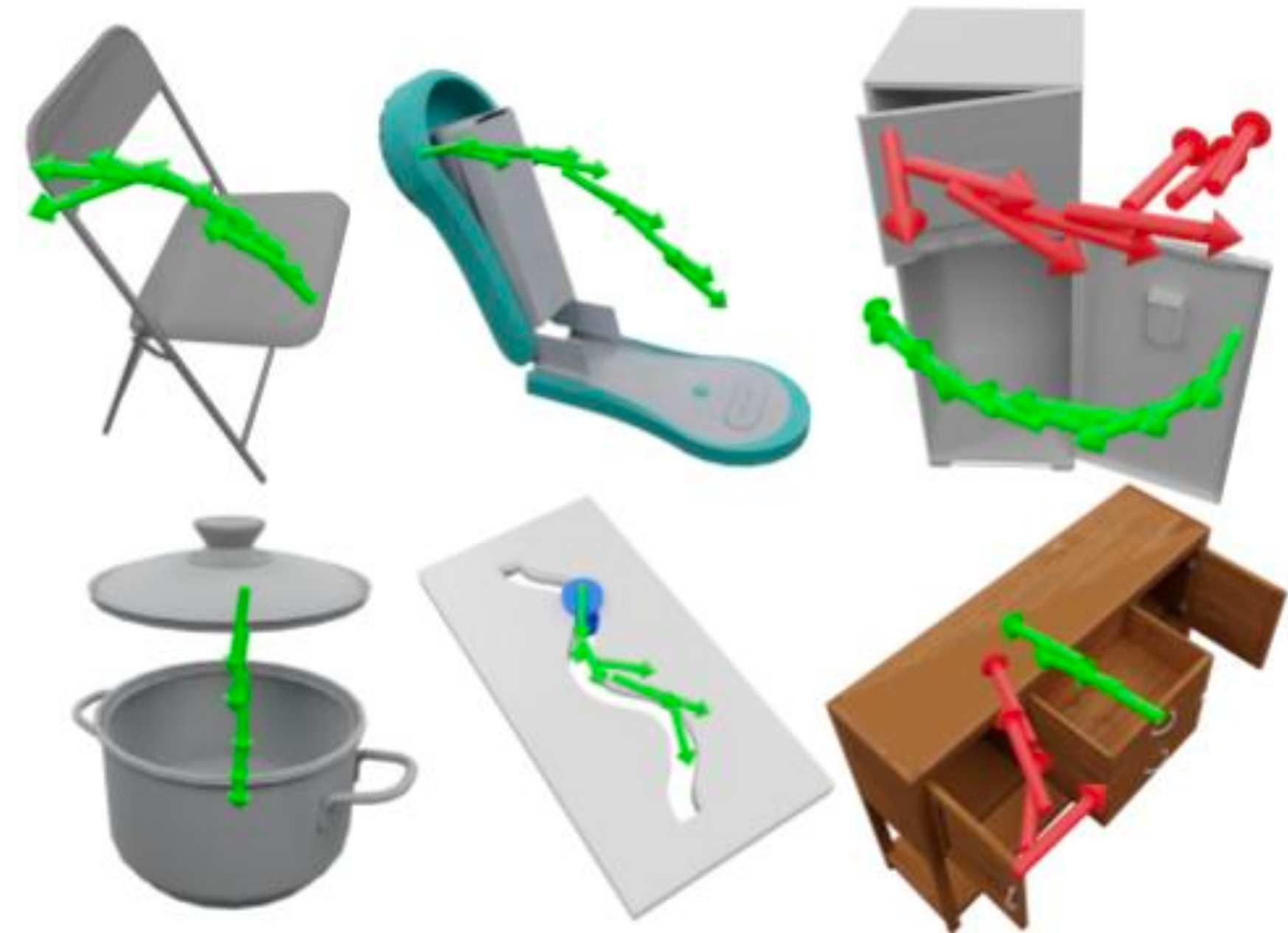




# Structure from Action

## UMPNet: Universal Manipulation Policy Network for Articulated Objects

Zhenjia Xu, Zhanpeng He, Shuran Song



# Universal Manipulation Policy



Action trajectories may vary drastically due to objects kinematic structures and geometry.

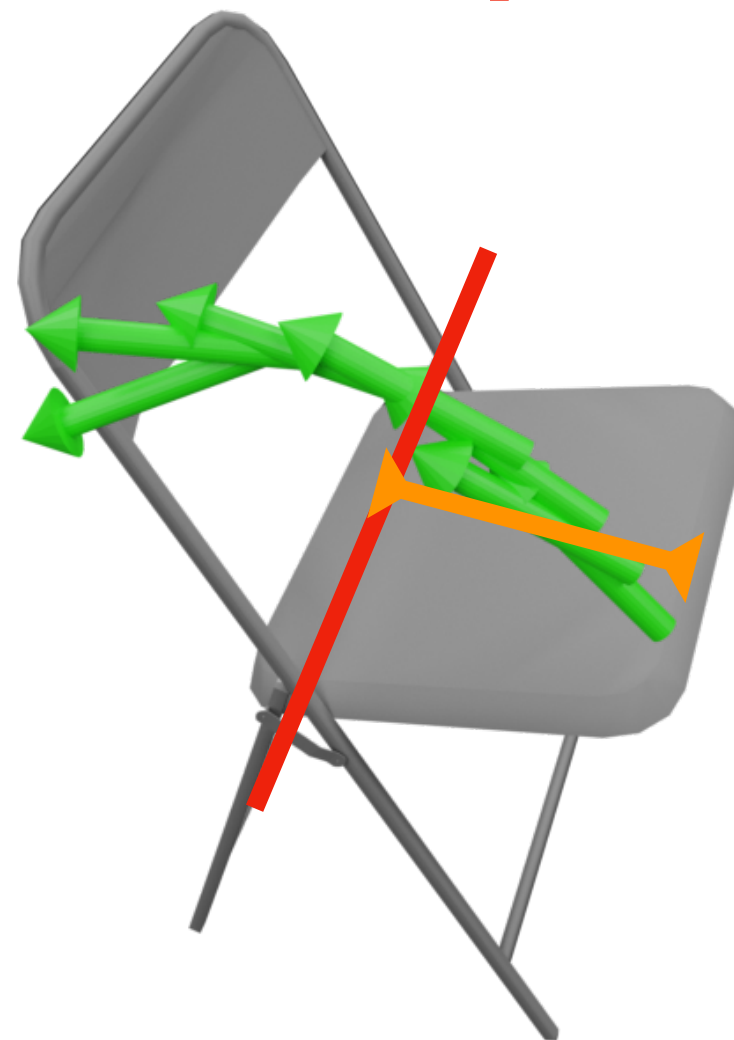
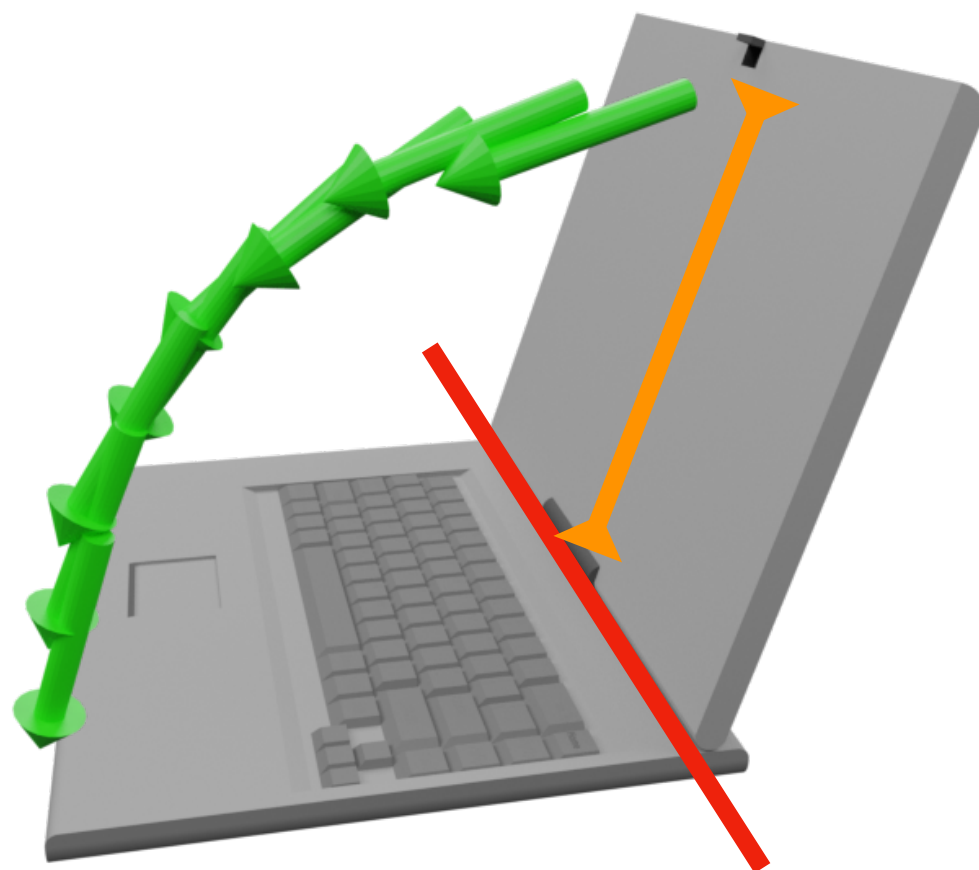
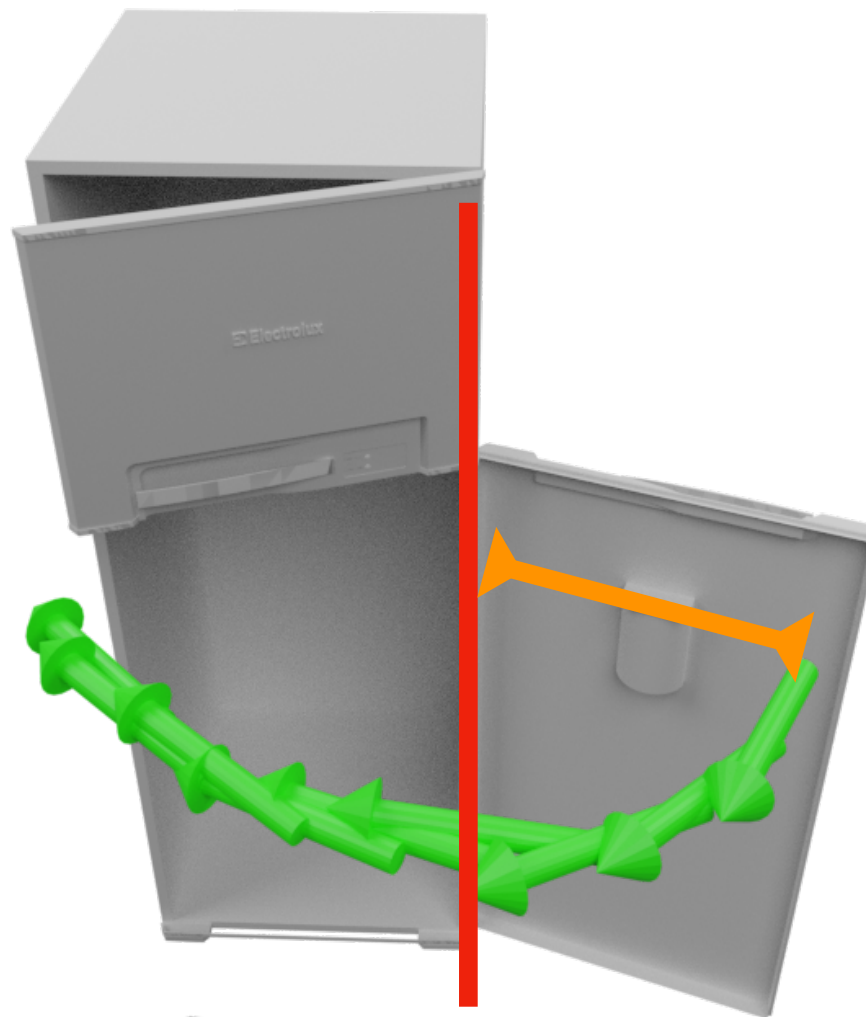
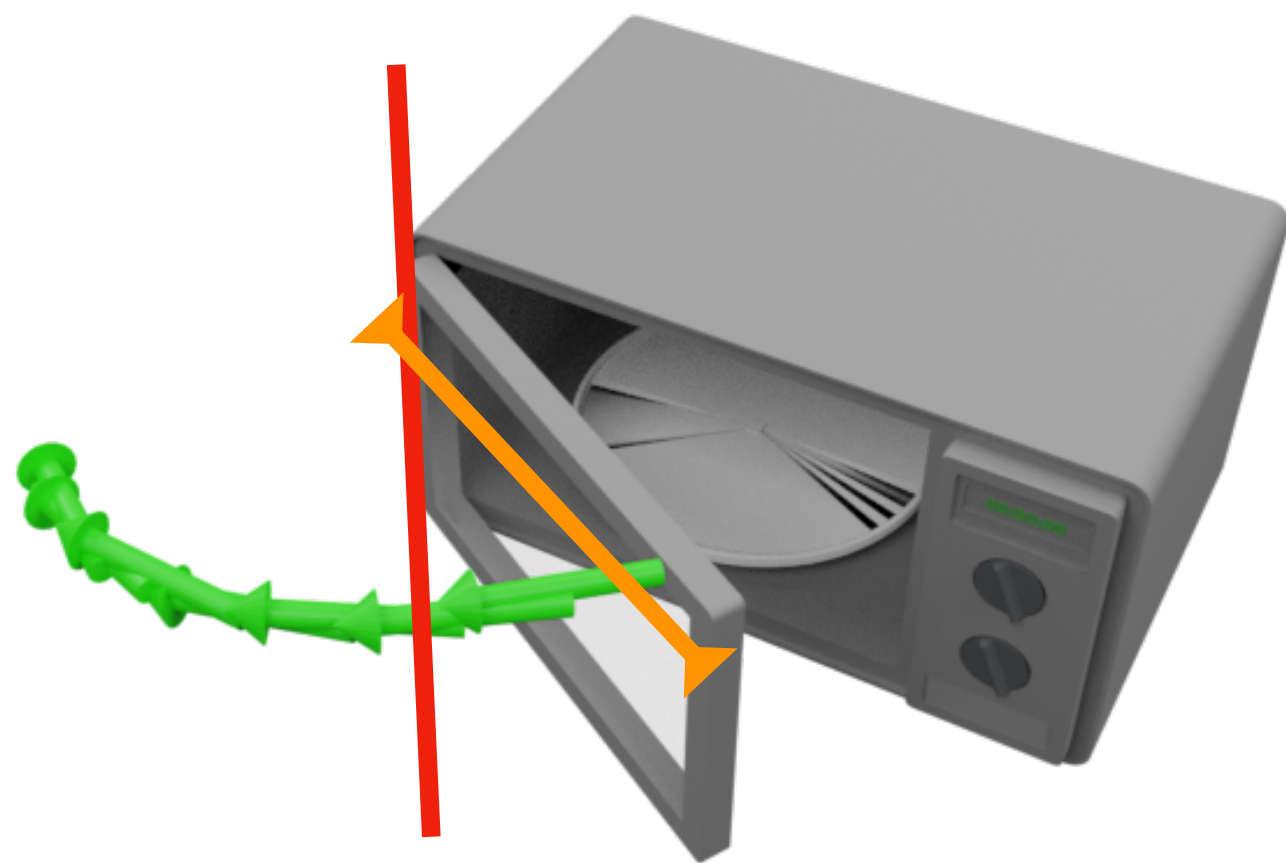


Goal: learn a single manipulation policy to handle all these objects from visual observations.



# Universal Manipulation Policy

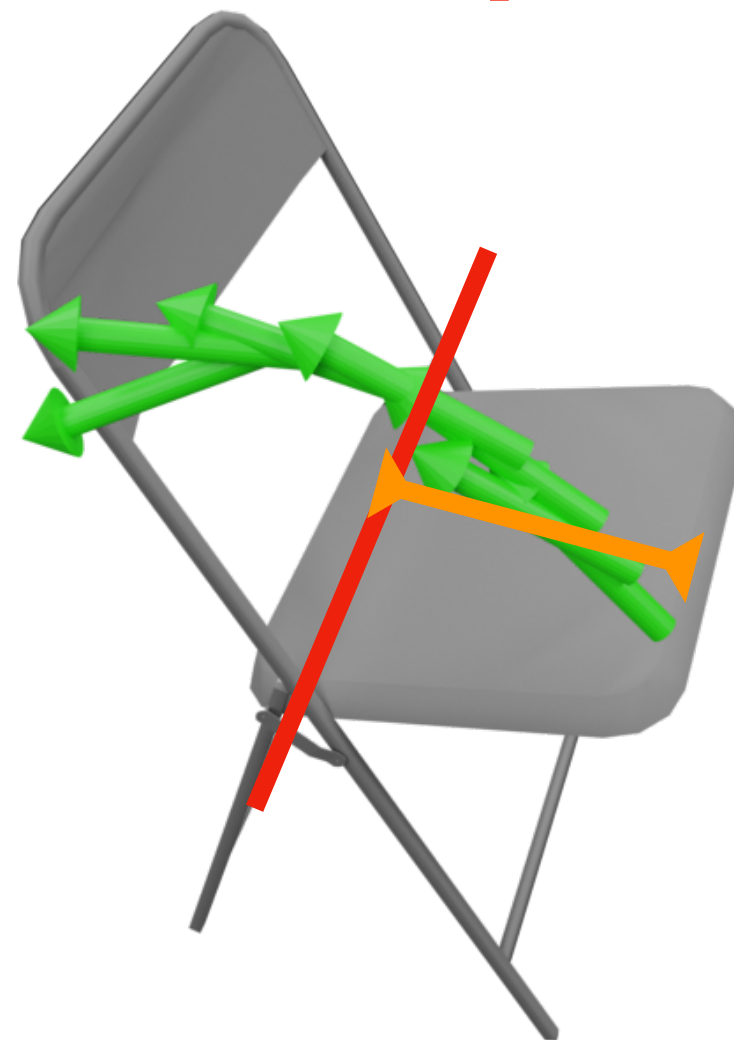
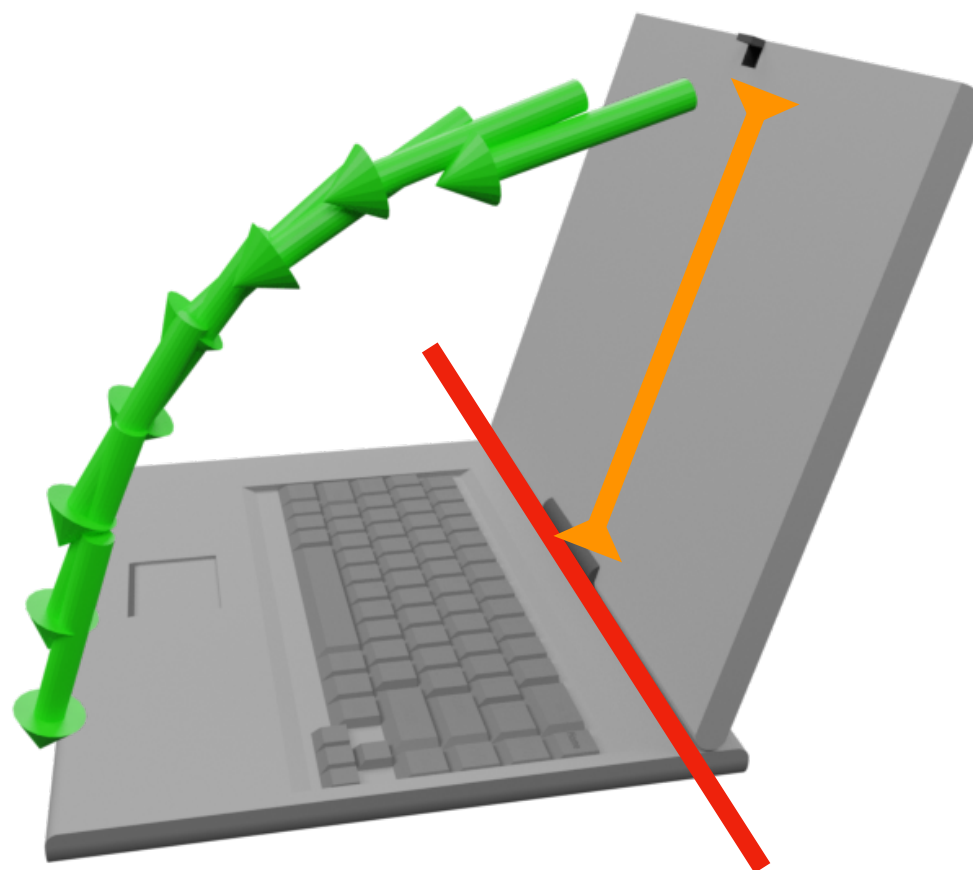
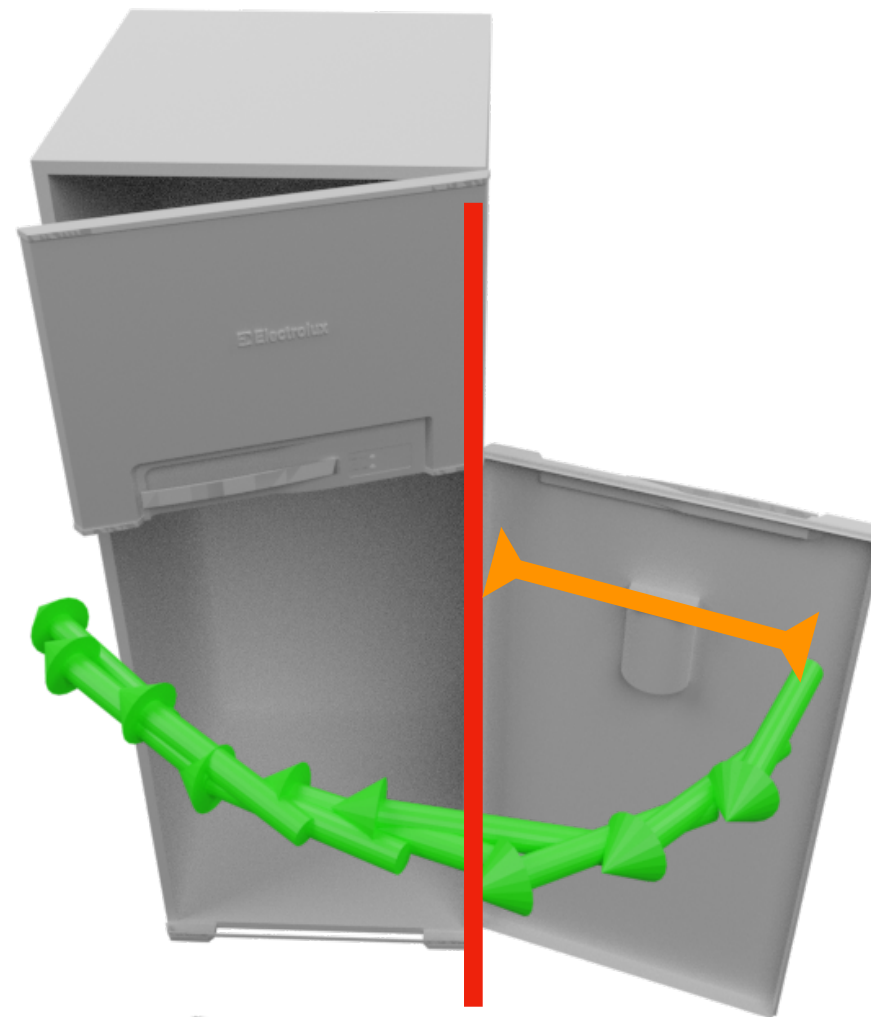
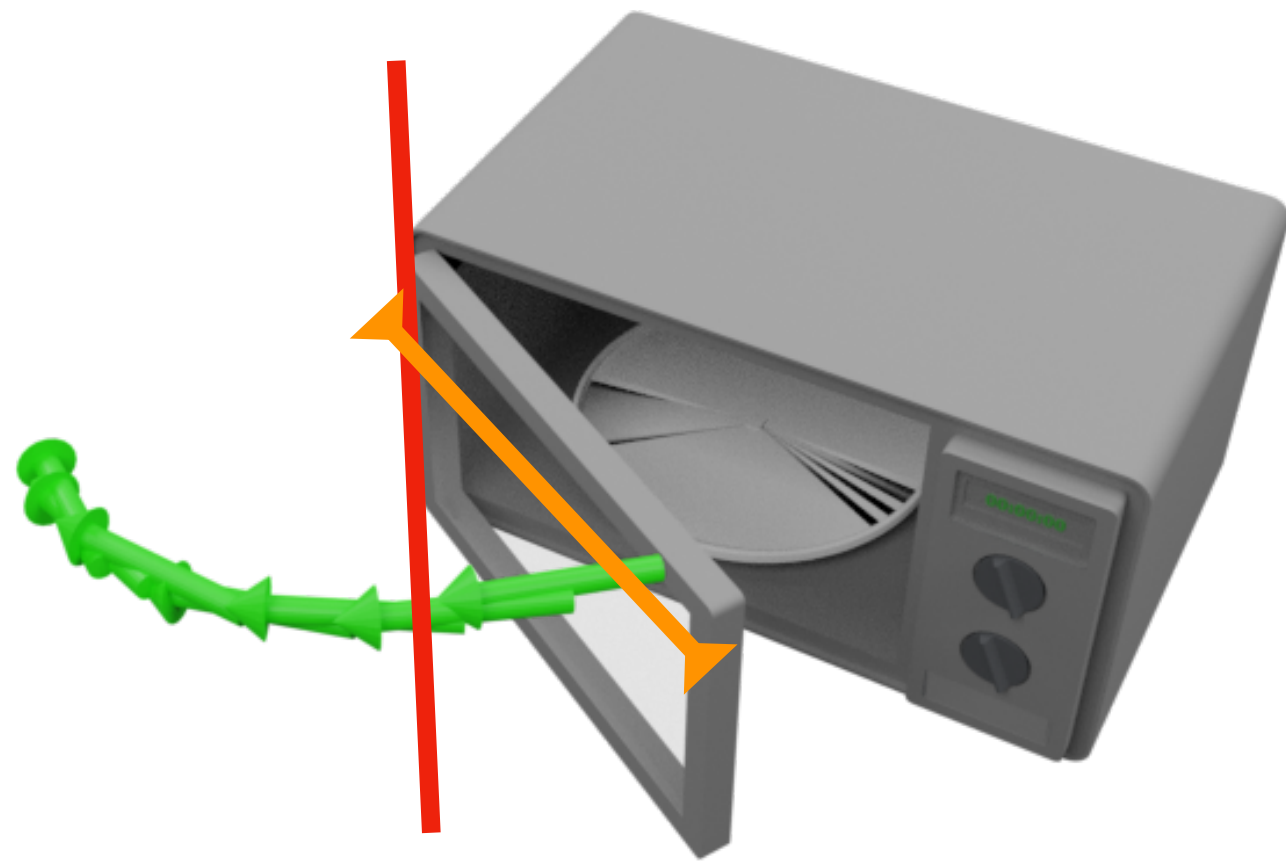
Why it is Possible ??



Can be summarized by a similar high-level function conditioned on the objects' underlying kinematic structure.

# Universal Manipulation Policy

Why it is Possible ??



Learning to interact with a diverse set of articulated objects

Acquire Generalizable Knowledge on:

- Articulation structure
- How these structures would react to different actions.

Beyond a specific object instance or category



# What is an effective interaction policy?

- Where to interact

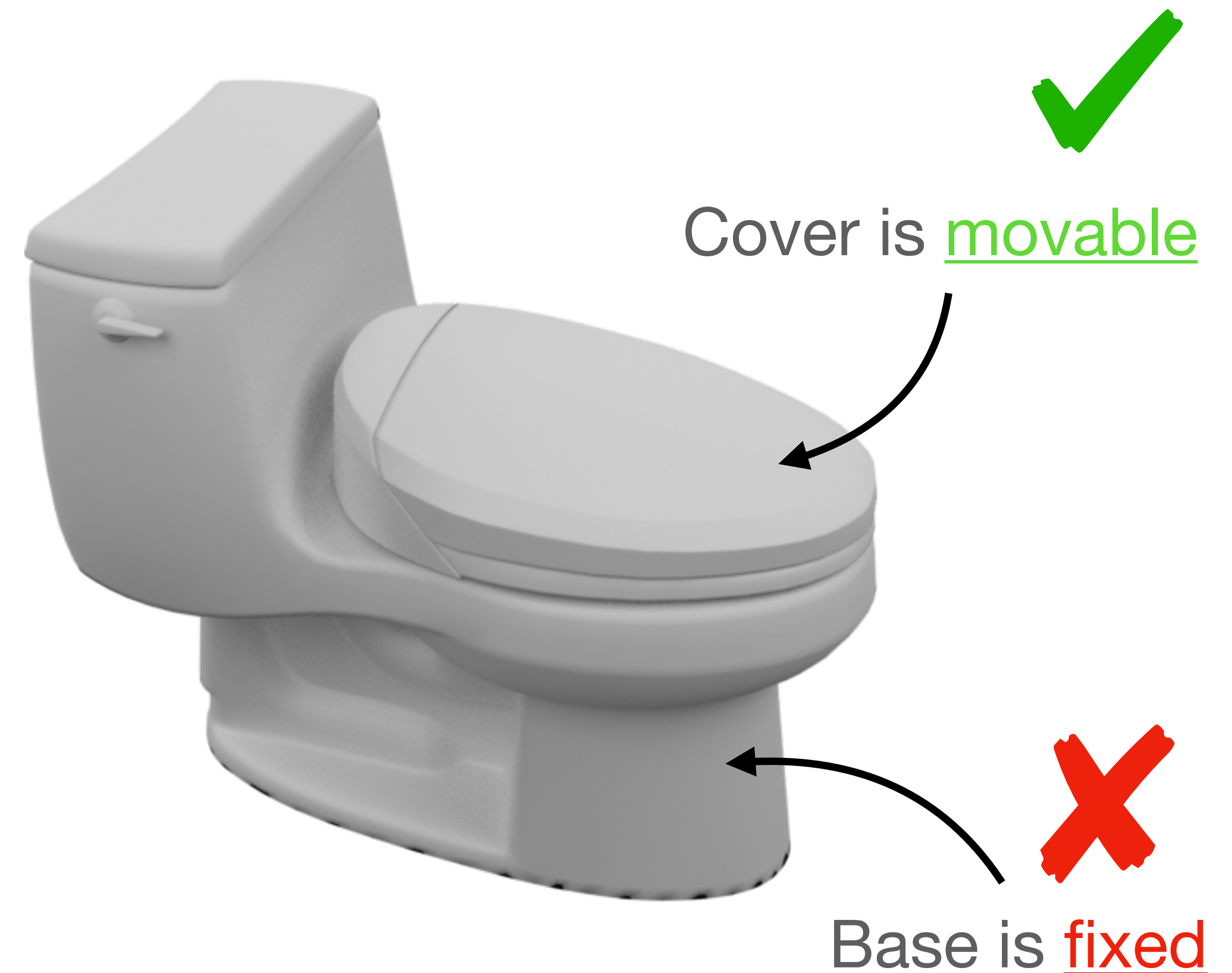
Interacting with the movable cover instead of the base

- Effective action direction

Pulling up instead of pushing down

- Arrow-of-Time awareness

Keeping pulling up the cover to visit novel states instead of moving up-and-down



# What does the policy need to learn?

- Where to interact

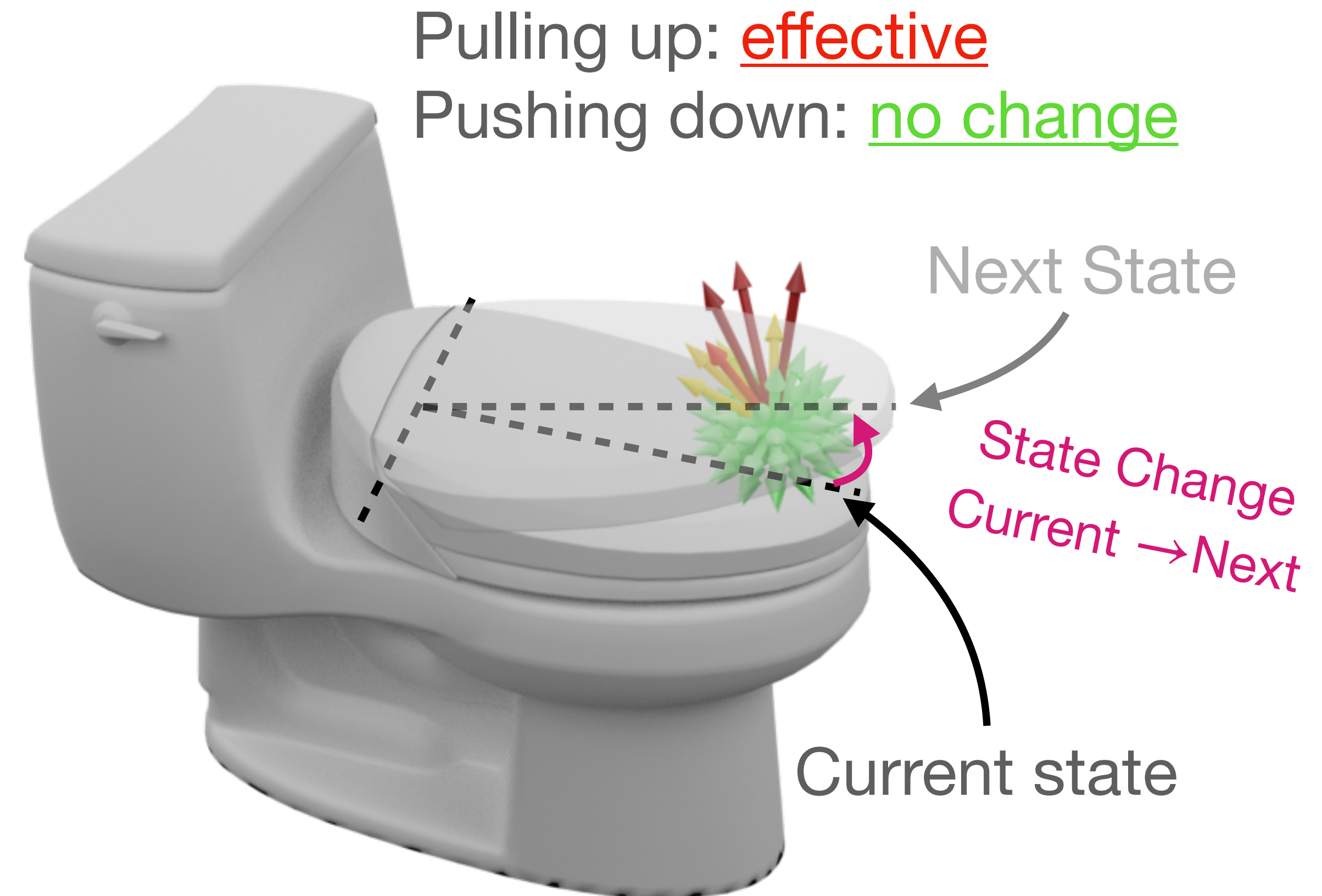
Interacting with the movable cover instead of the base

- Action direction

Pulling up instead of pushing down

- Arrow-of-Time awareness

Keeping pulling up the cover to visit novel states instead of moving up-and-down





# What does the policy need to learn?

- Where to interact

Interacting with the movable cover instead of the base

- Effective action direction

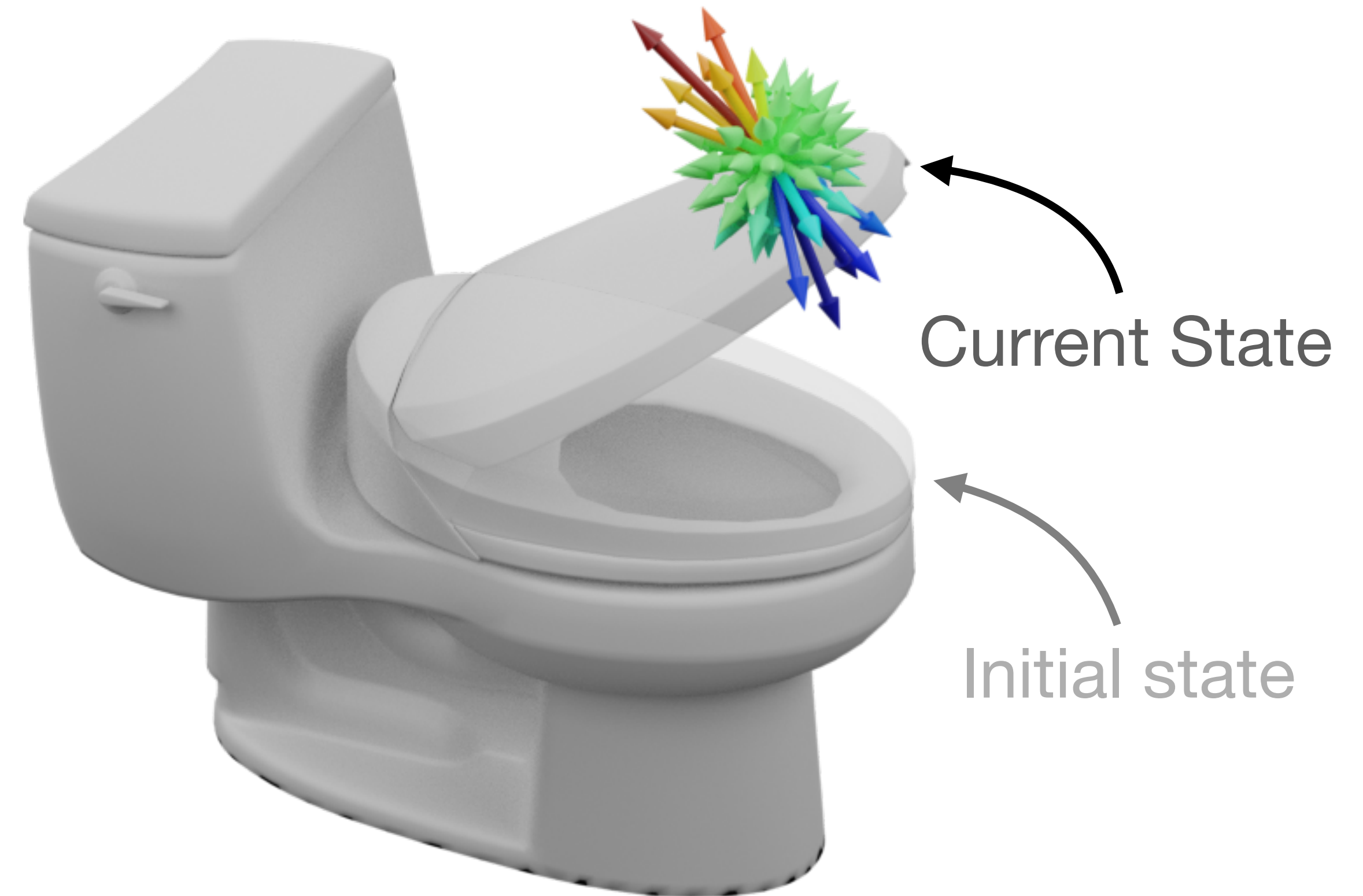
Pulling up instead of pushing down

- Consistent action direction

Keeping pulling up the cover to visit novel states instead of moving up-and-down. (Arrow-of-Time awareness)

Pulling up: toward novel state

Pushing down: back to past



# What does the policy need to learn?

- Where to interact

Interacting with the movable cover instead of the base

- Effective action direction

Pulling up instead of pushing down

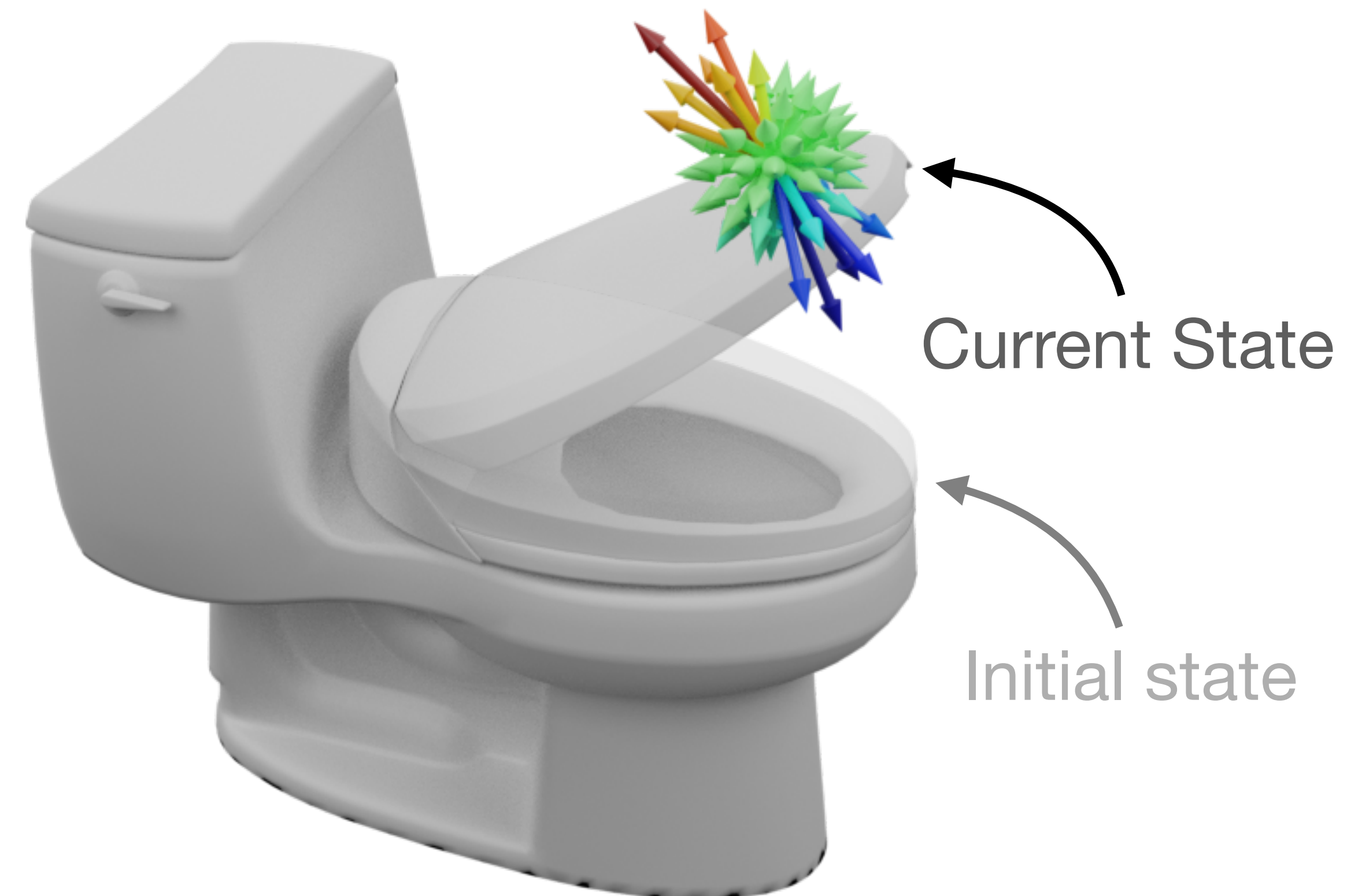
- Consistent action direction

Keeping pulling up the cover to visit novel states instead of moving up-and-down. (Arrow-of-Time awareness)

↓  
Goal conditioned manipulation, without goal-conditioned training

Pulling up: toward novel state

Pushing down: back to past

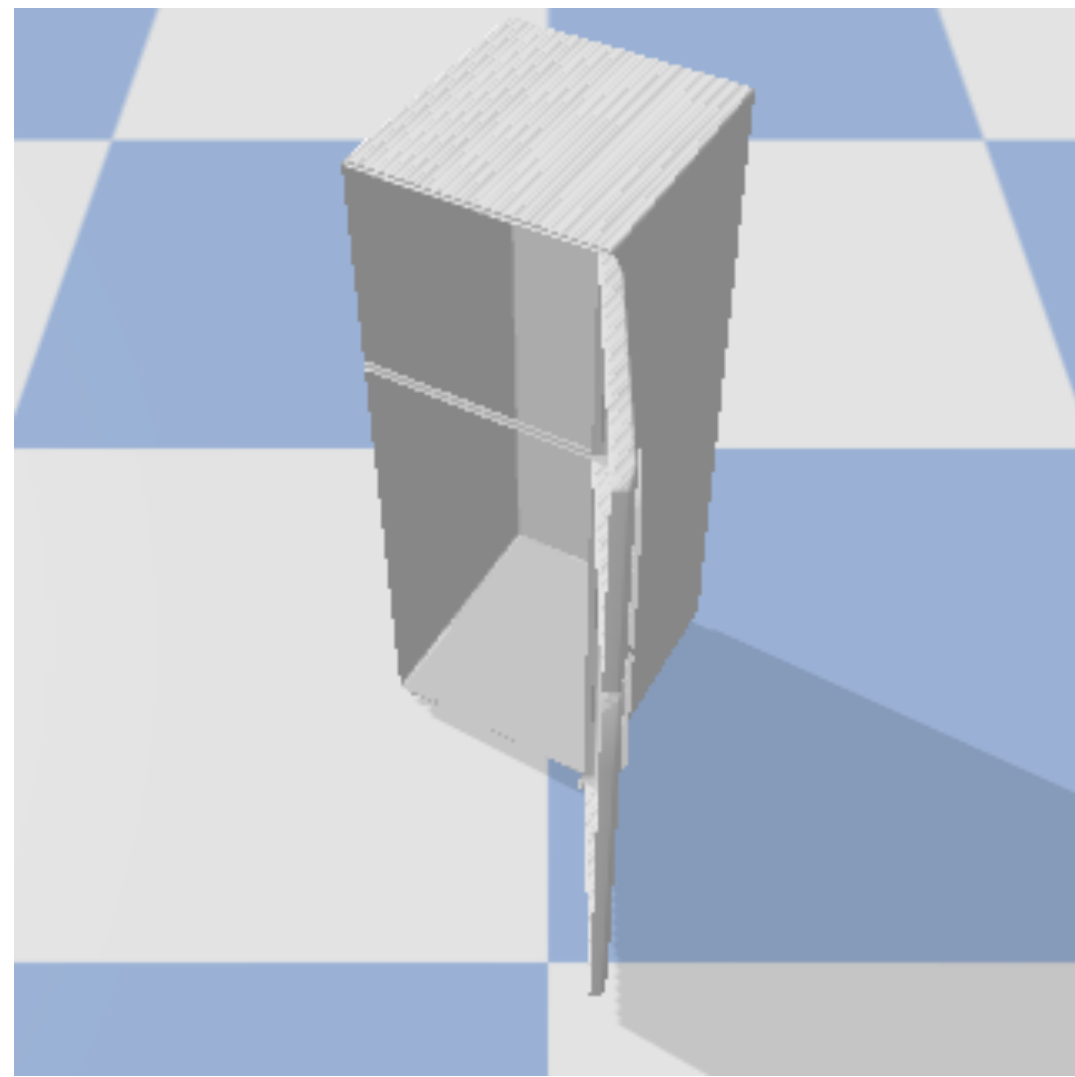




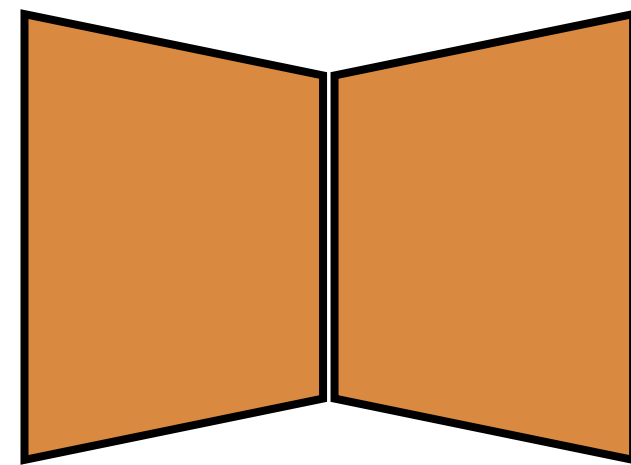
# Universal Manipulation Policy Network

# Universal Manipulation Policy Network

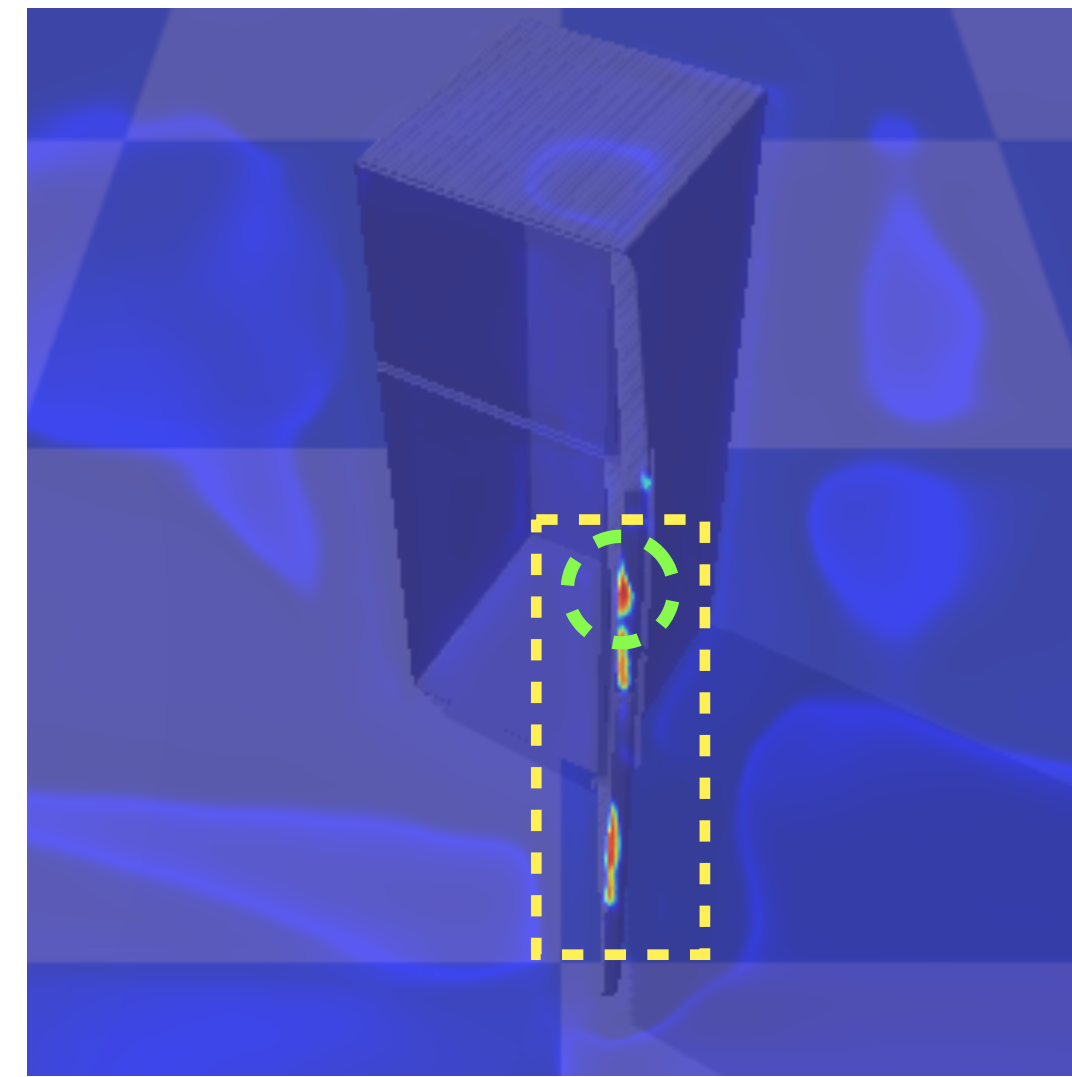
Interaction Position Inference



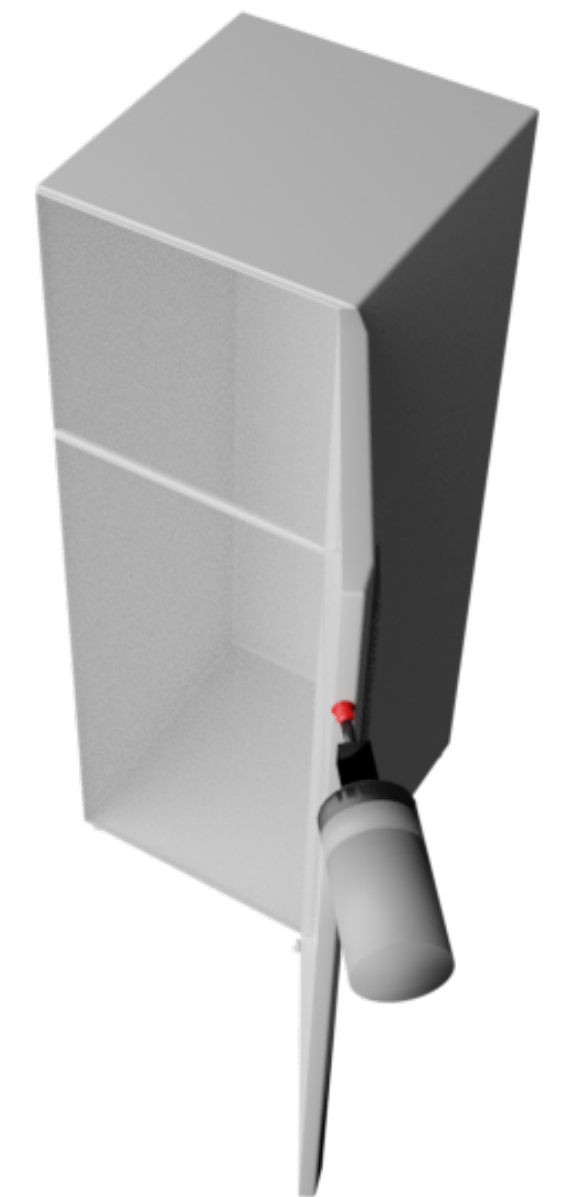
Visual Observation



Position Net



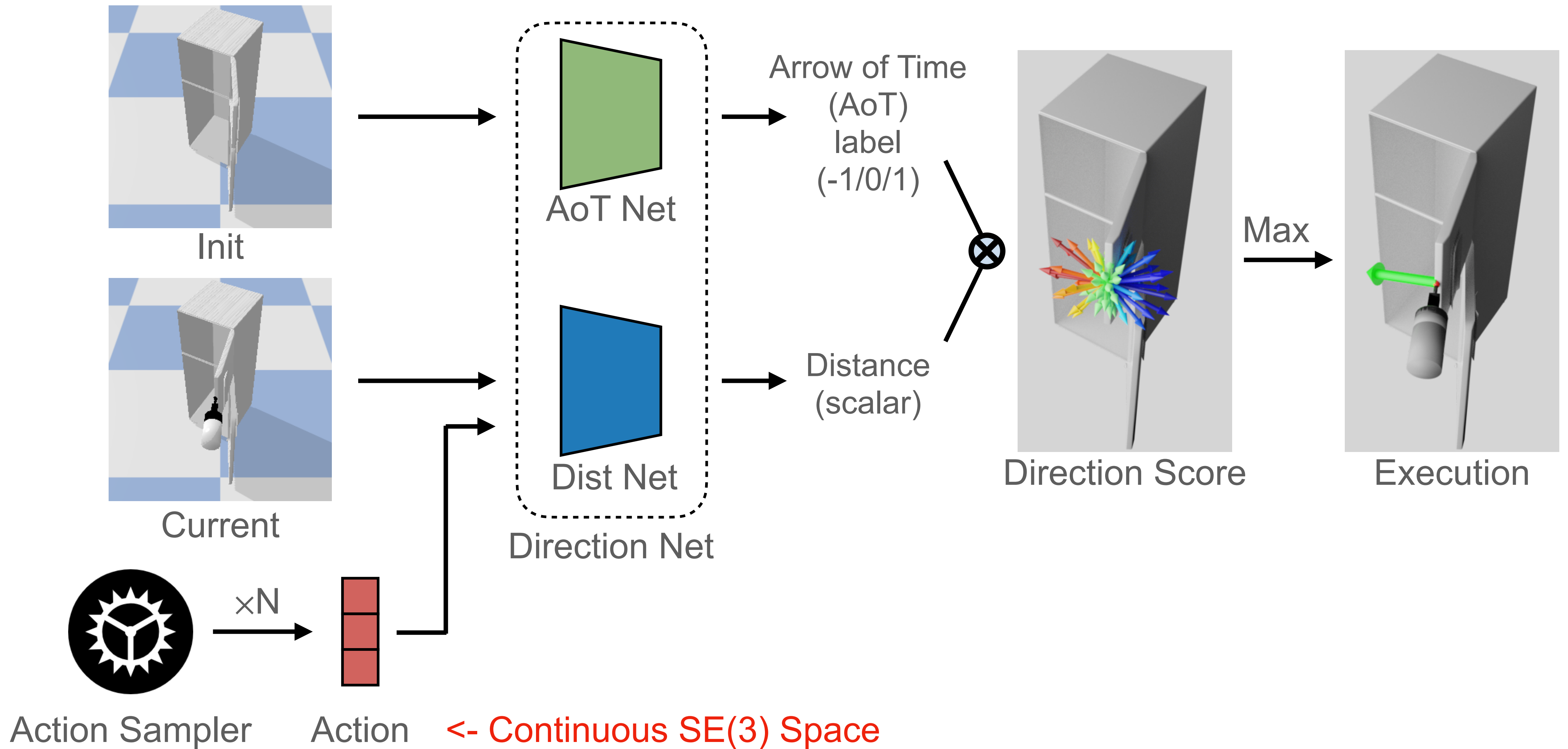
Position Affordance



Execution



# UMPNet: Direction Inference



# Training and Testing Objects

Training Categories (12)



Testing Categories (10)



**The policy is trained with self-guided exploration without any human demonstrations, scripted policy, or pre-defined goal conditions.**



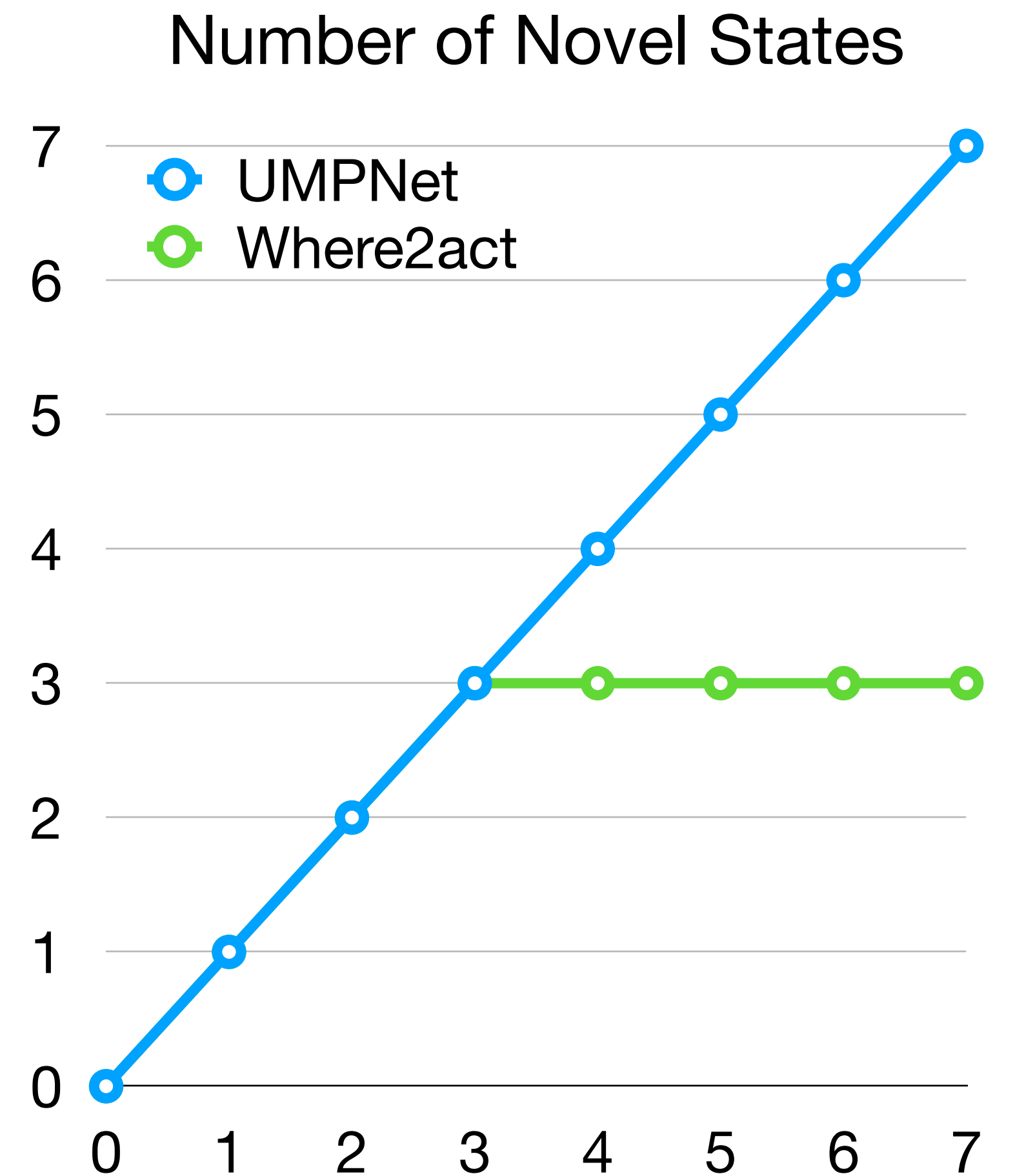
# Open-ended State Exploration



UMPNet

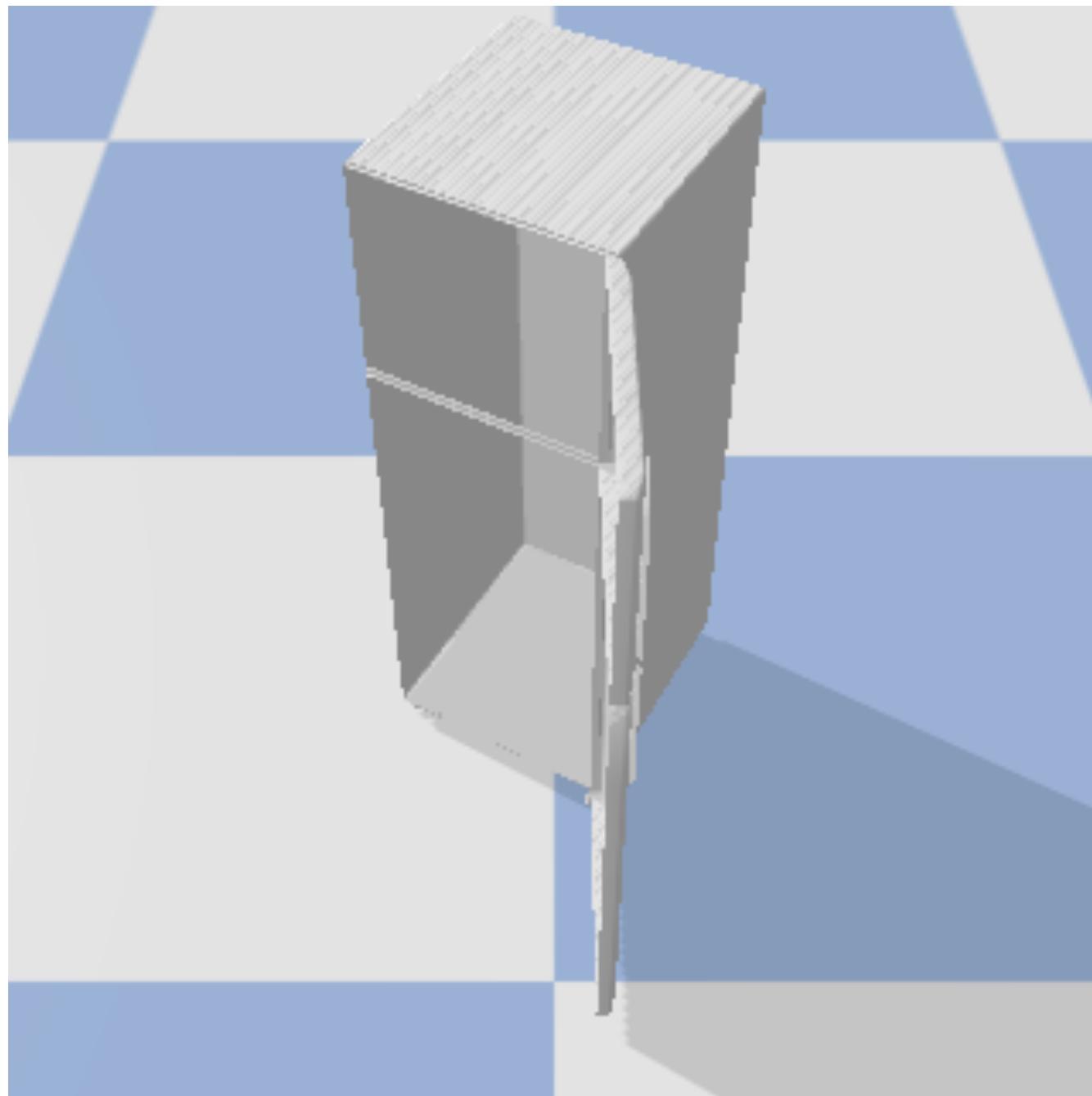


Where2act\*  
[Mo *et al.*]

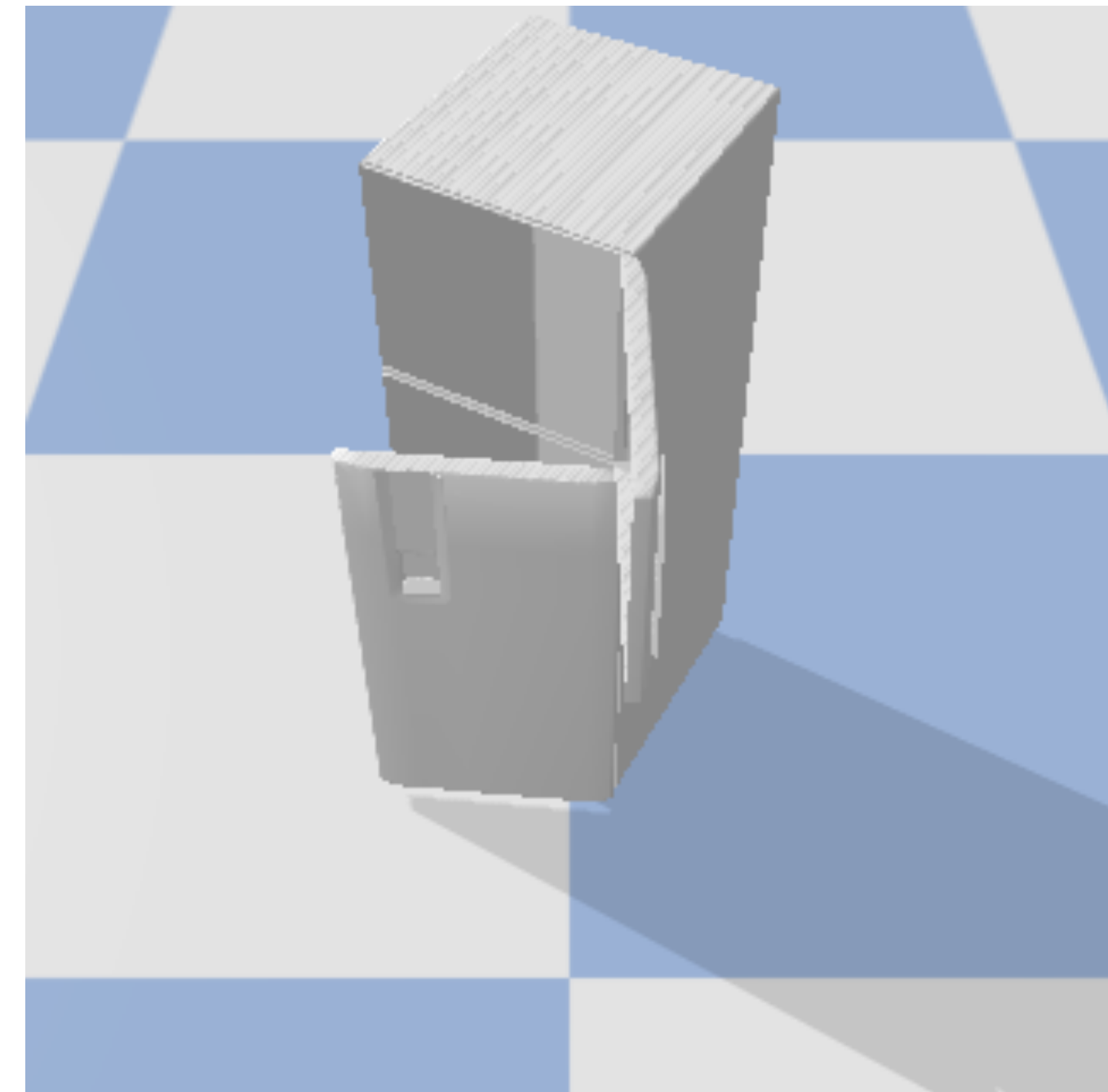
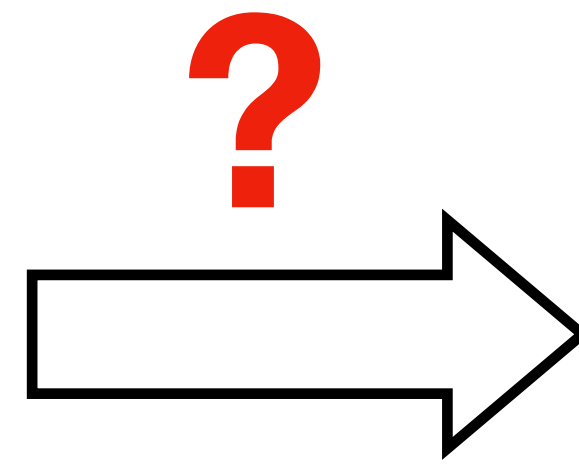


\*Where2act: infers a single-step interaction from current image observation

# Goal Conditioned Manipulation



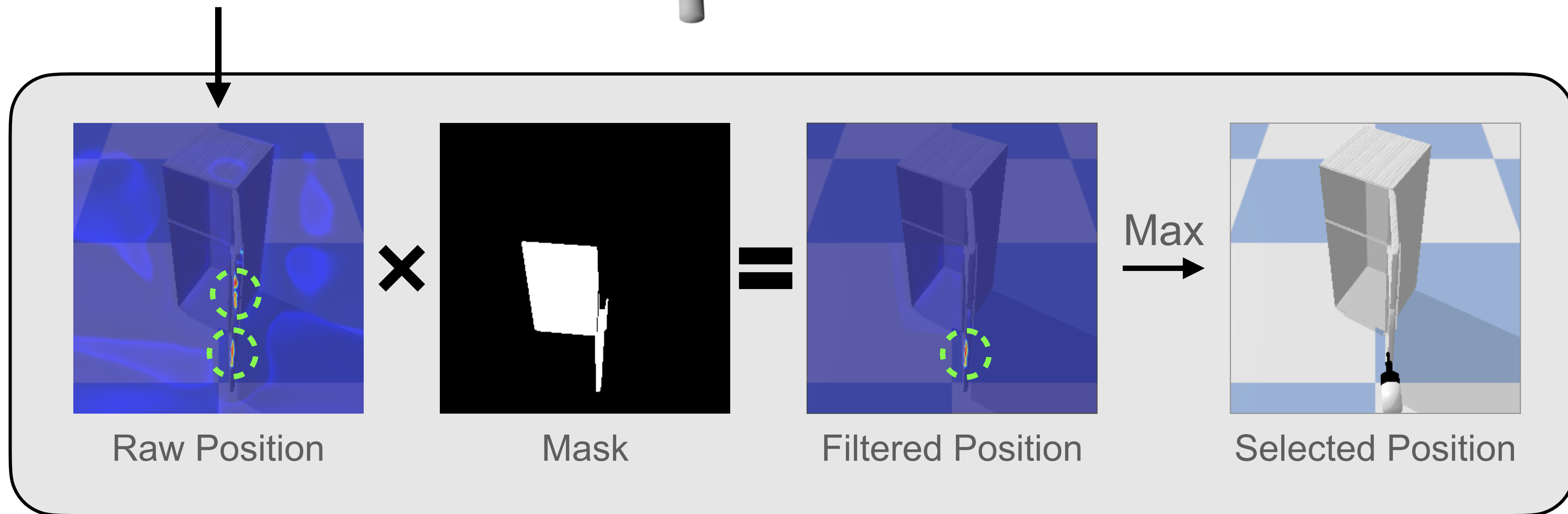
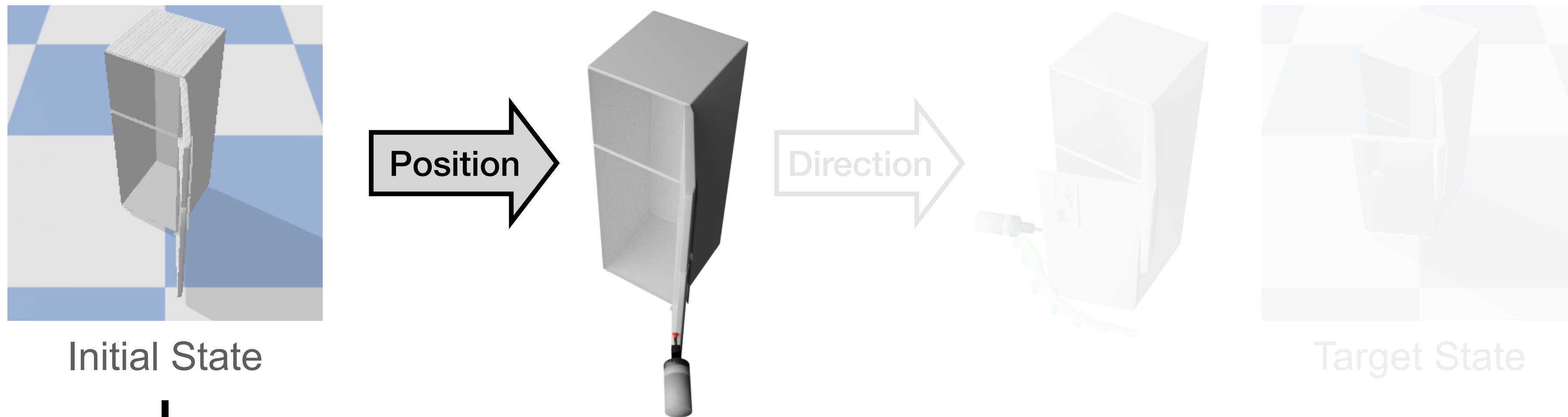
Initial State



Target State

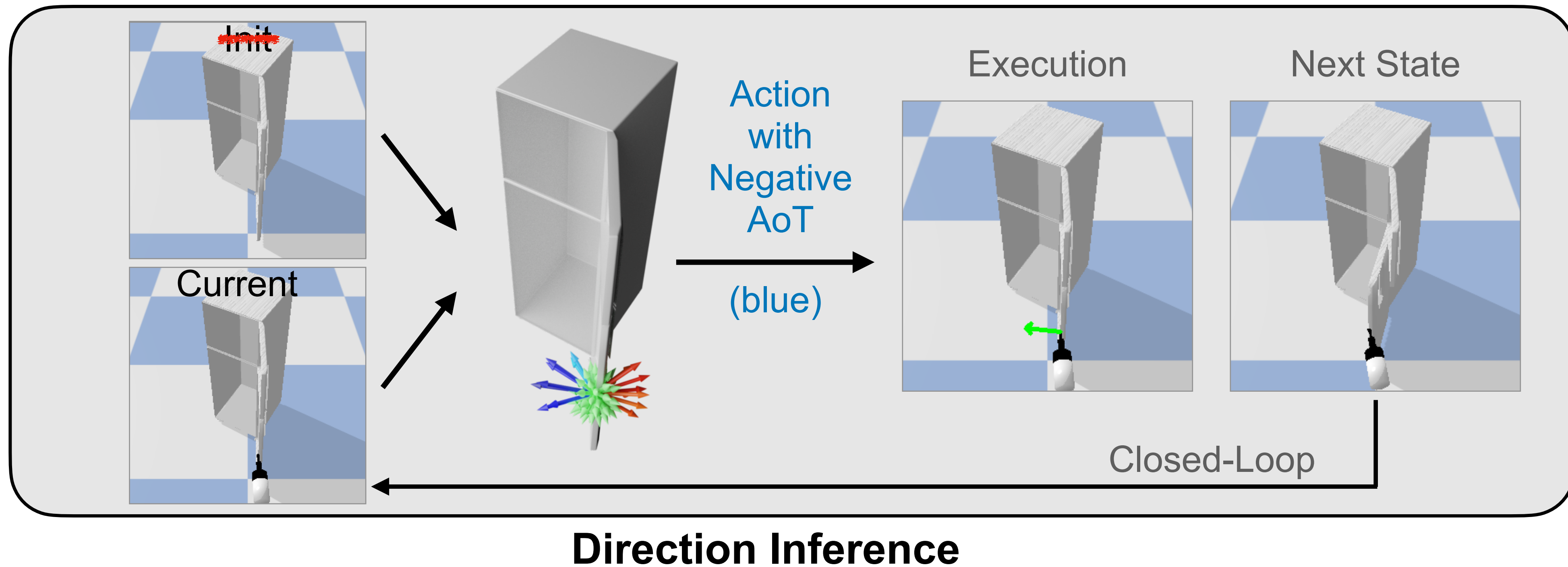
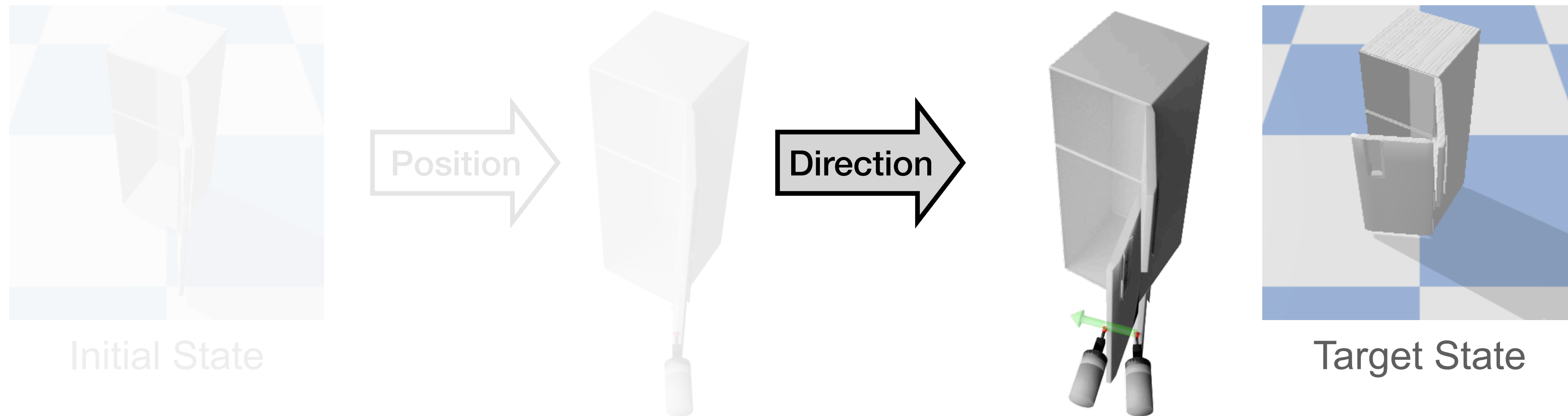


# Goal Conditioned Manipulation



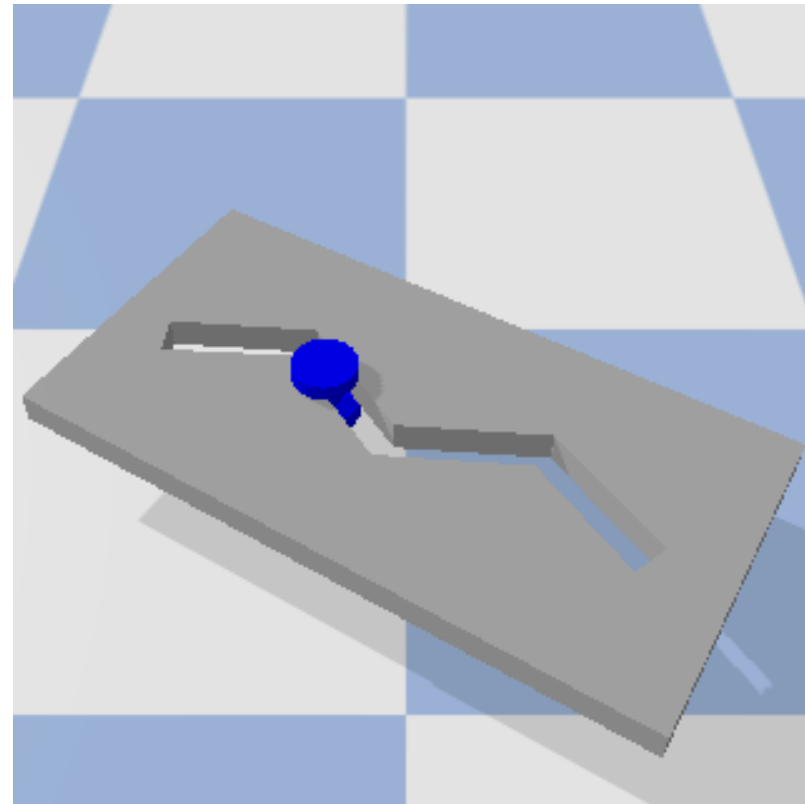
**Position Inference**

# Goal Conditioned Manipulation

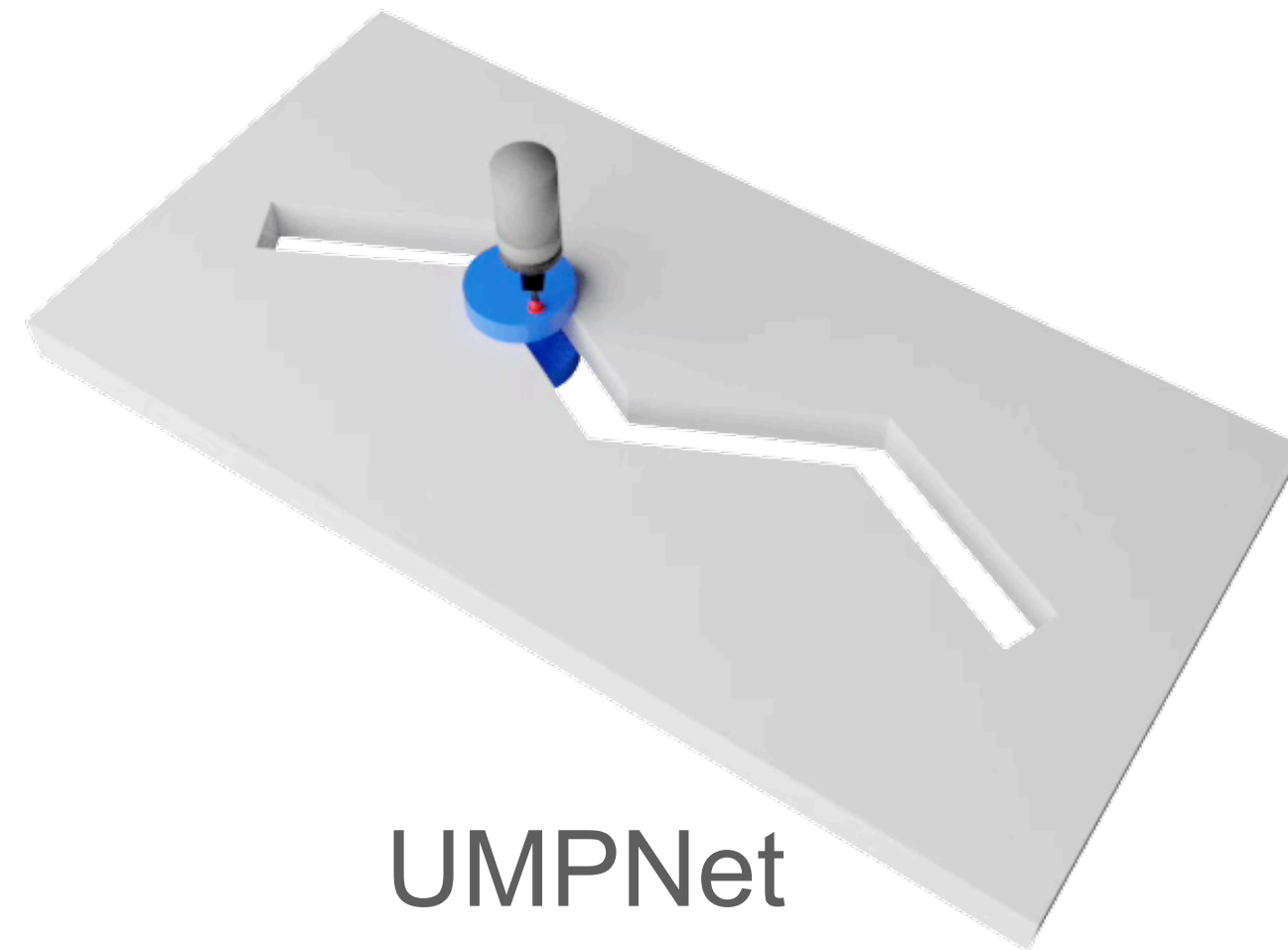




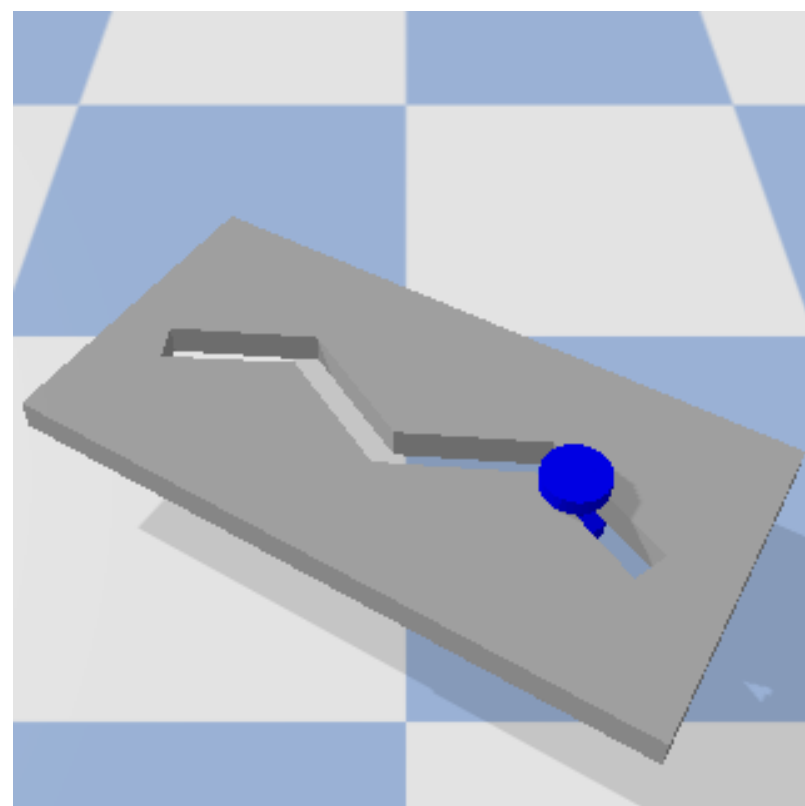
# Goal Conditioned Manipulation: Results



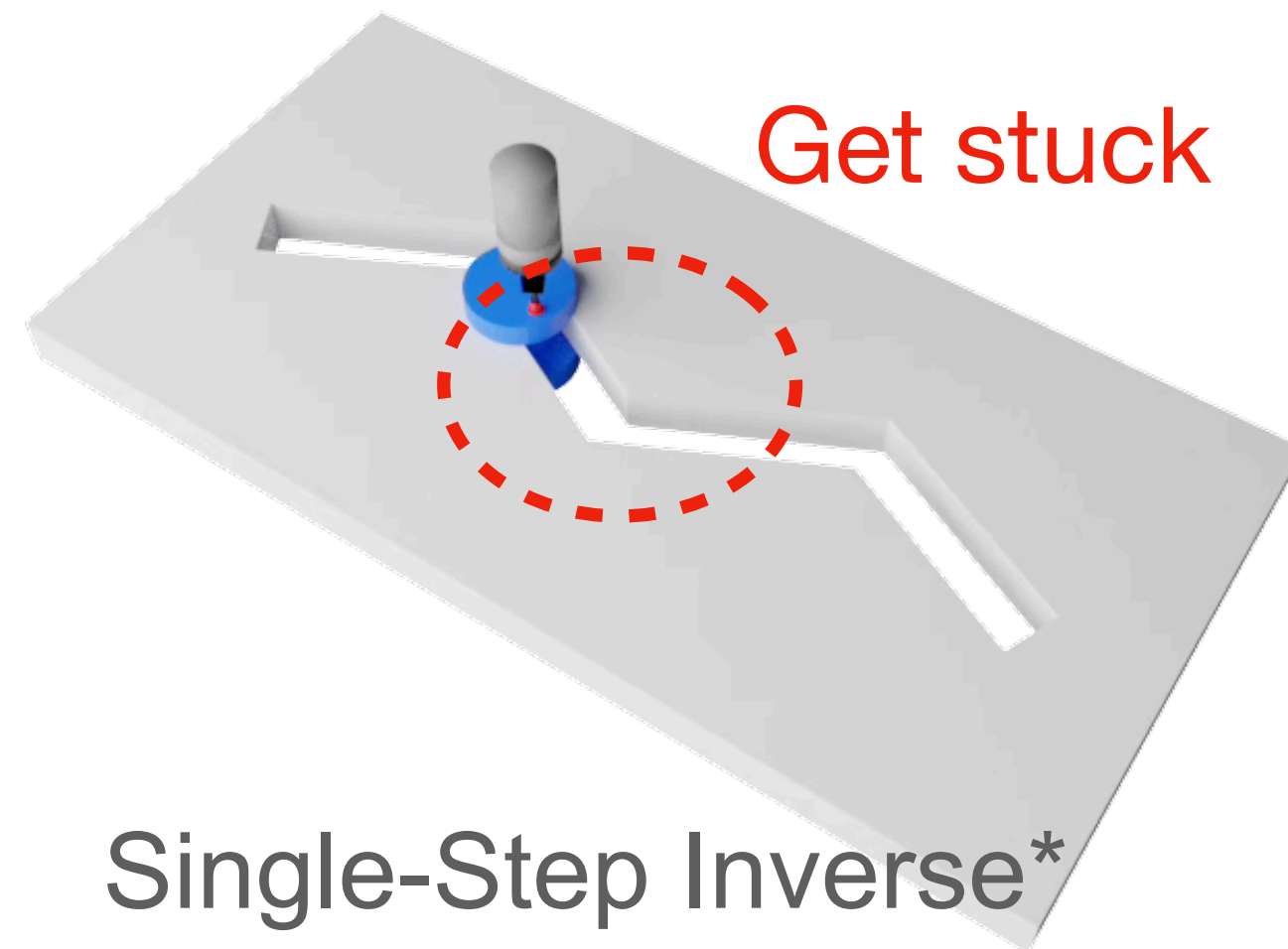
Initial State



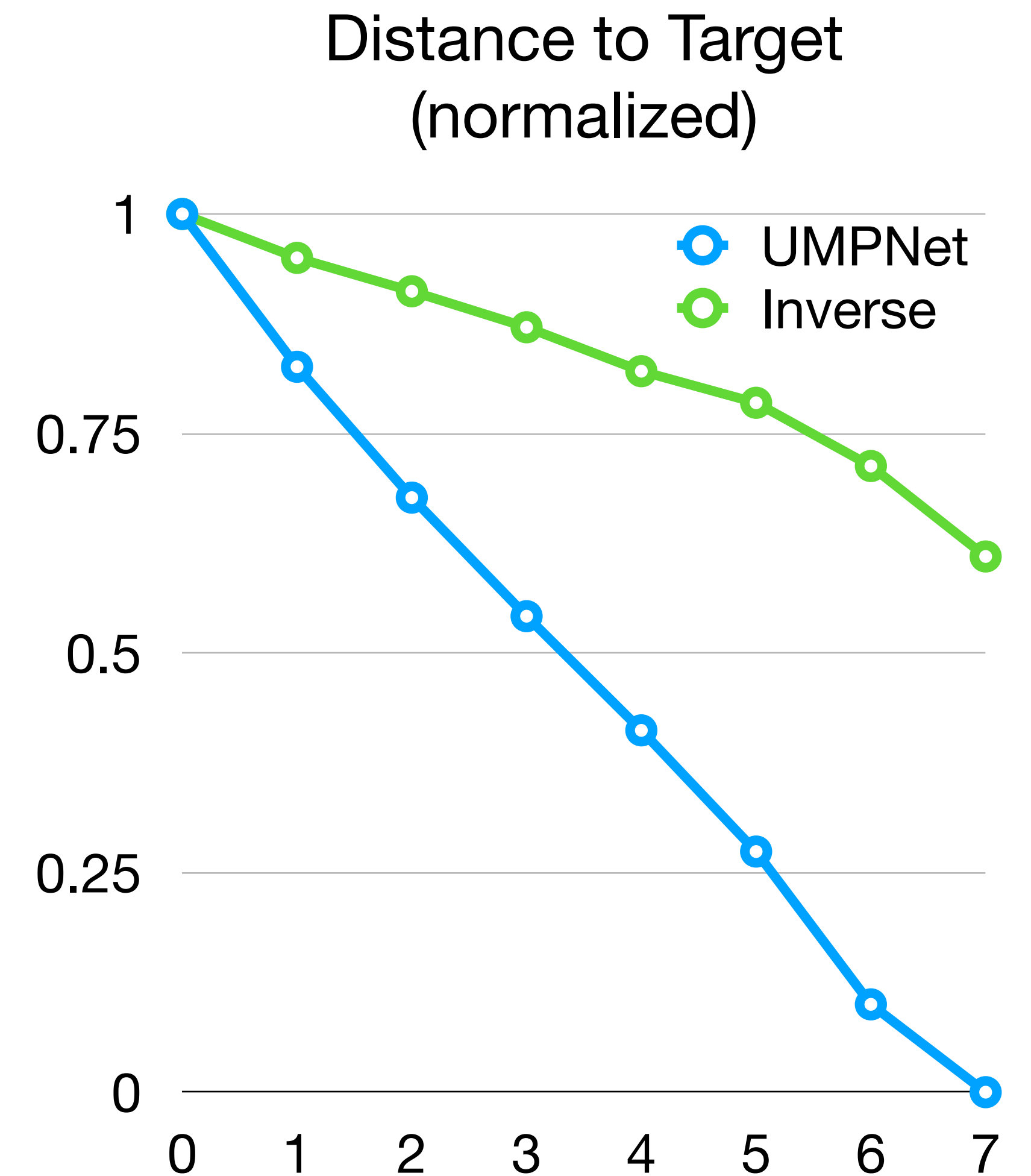
UMPNet



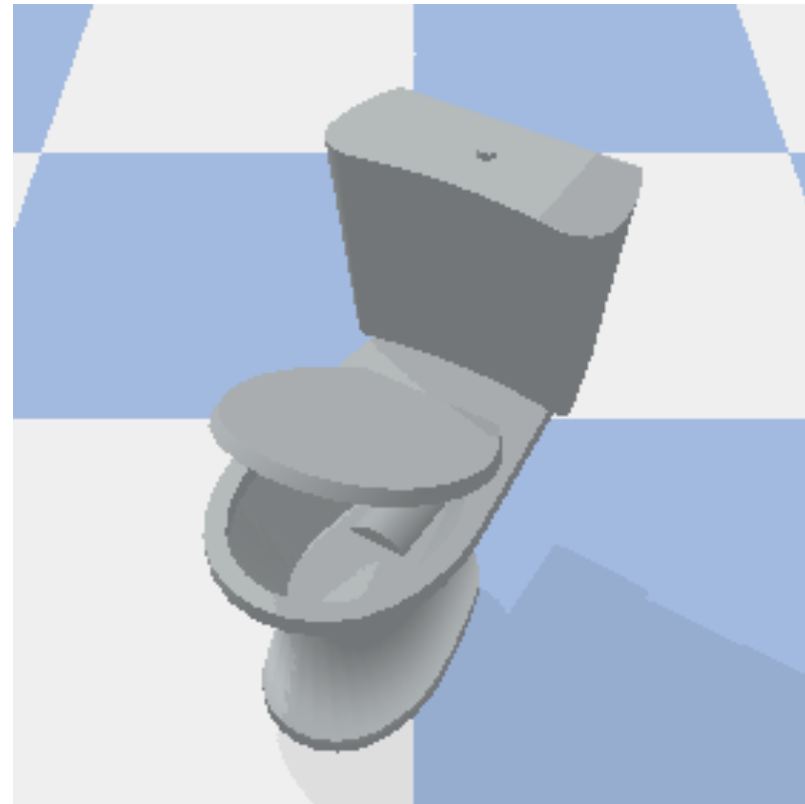
Target State



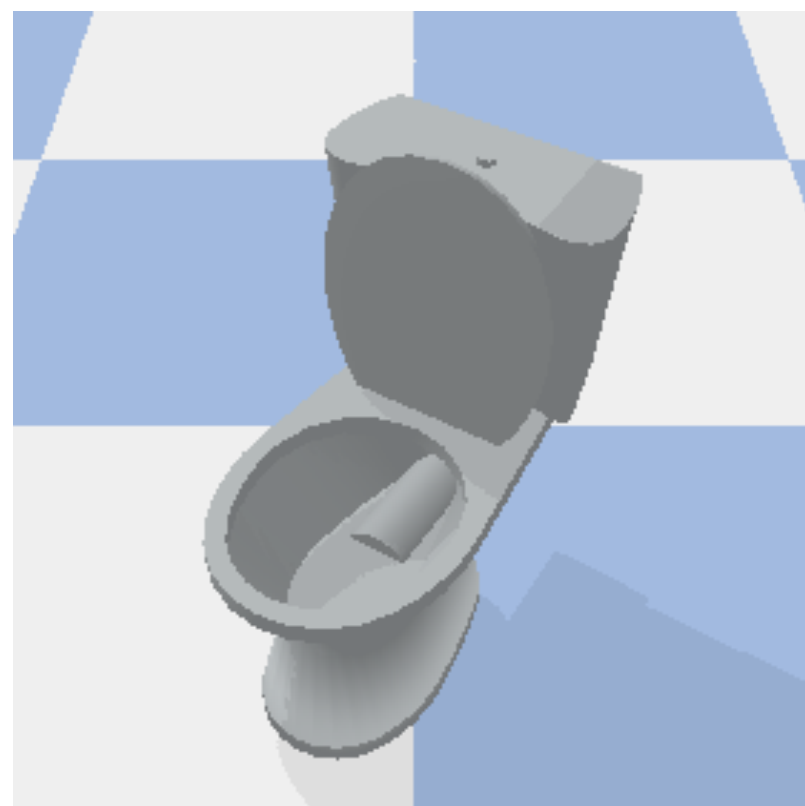
Single-Step Inverse\*



# Goal Conditioned Manipulation: Results



Initial State



Target State

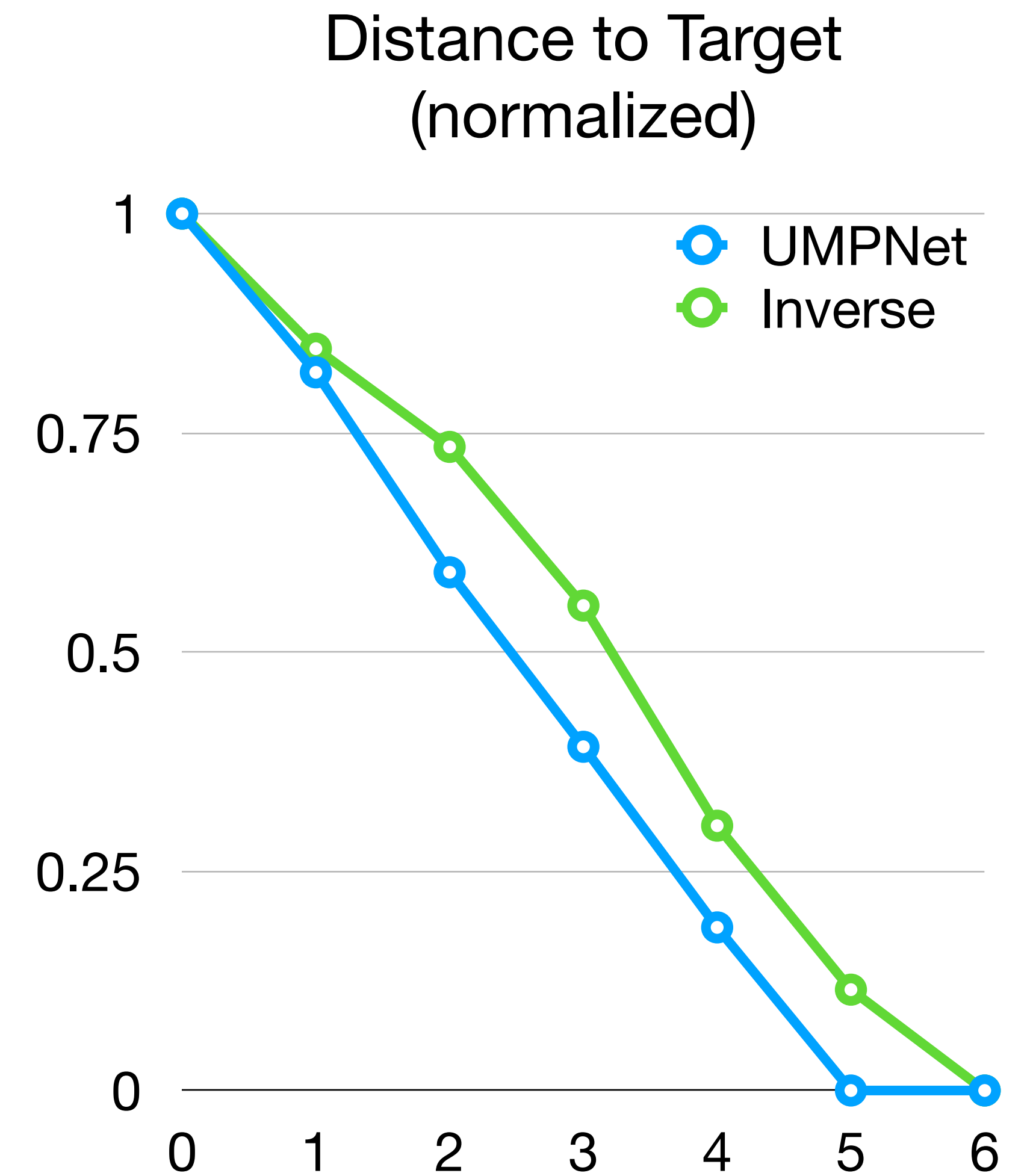


UMPNet

Bad  
direction



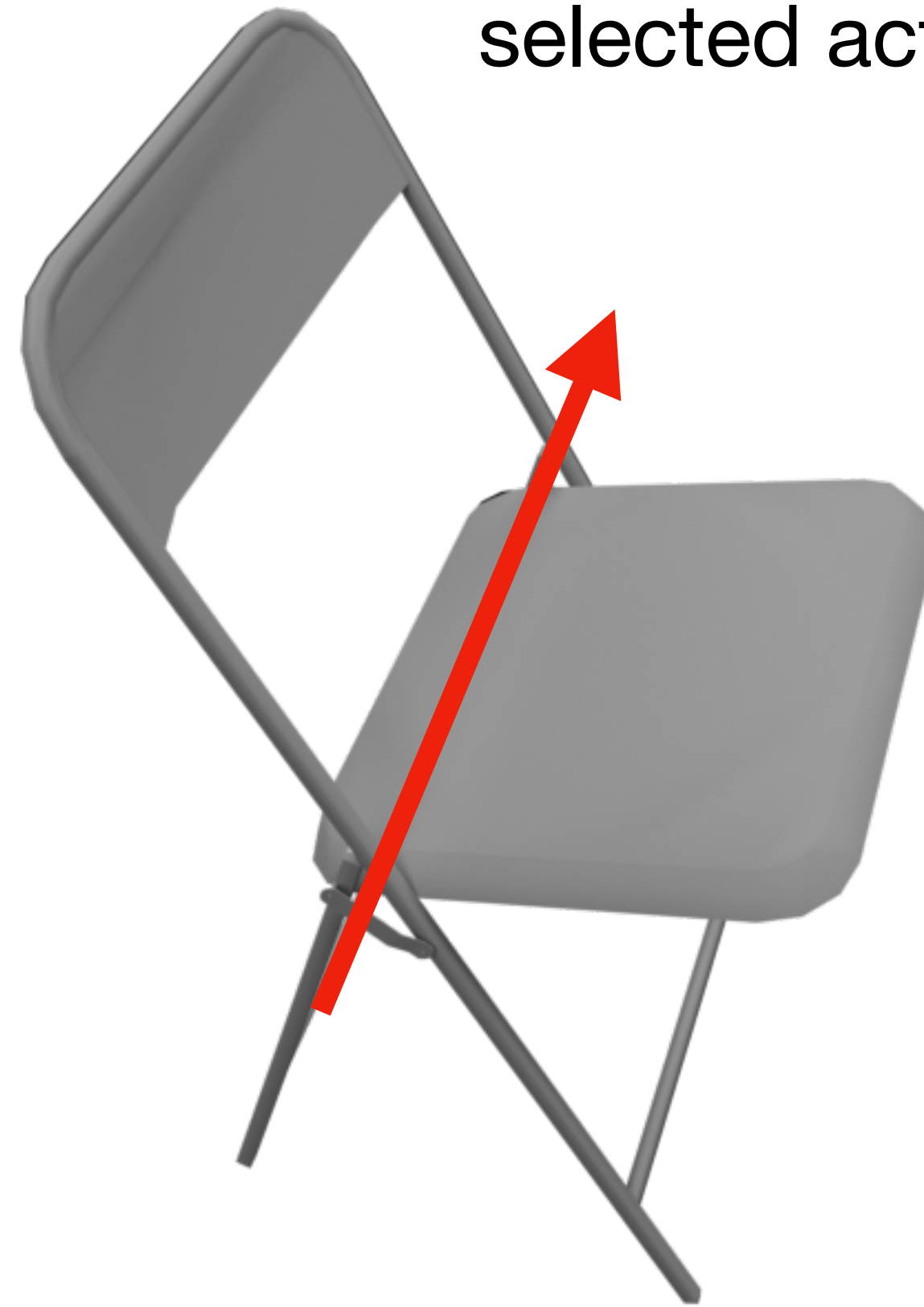
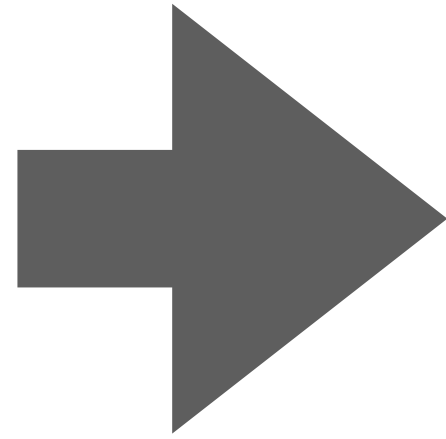
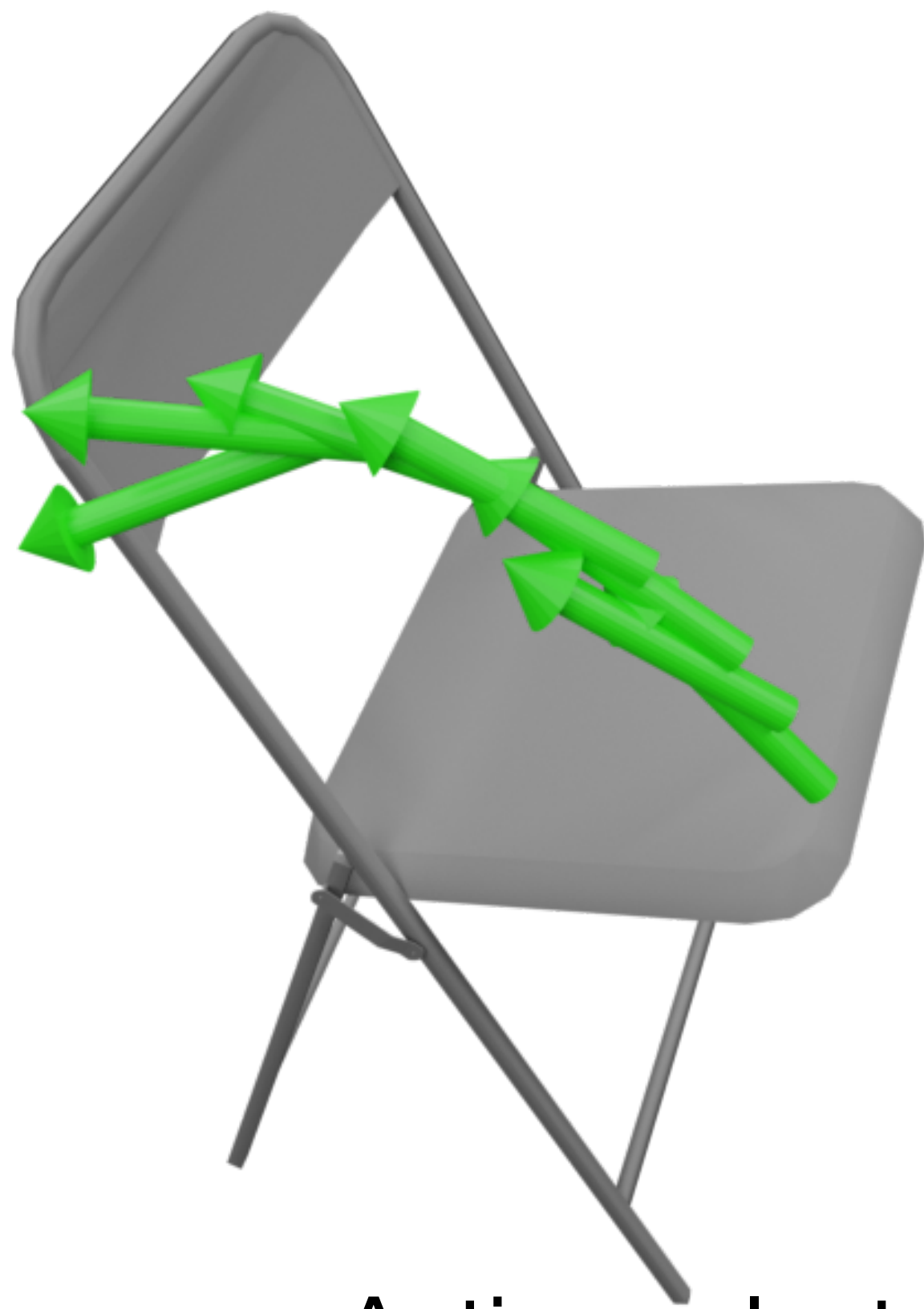
Inverse\*



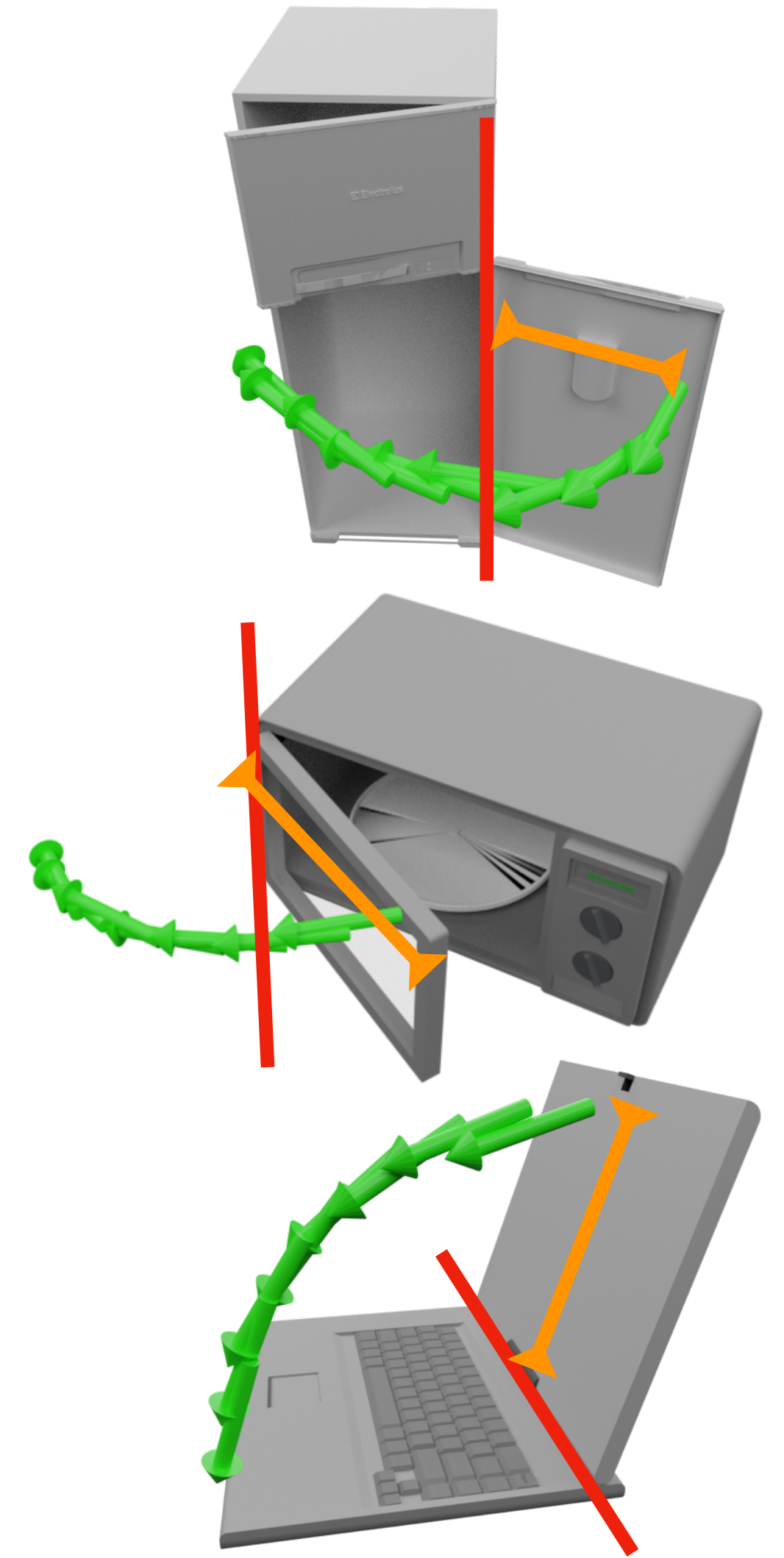


# Action → Structure

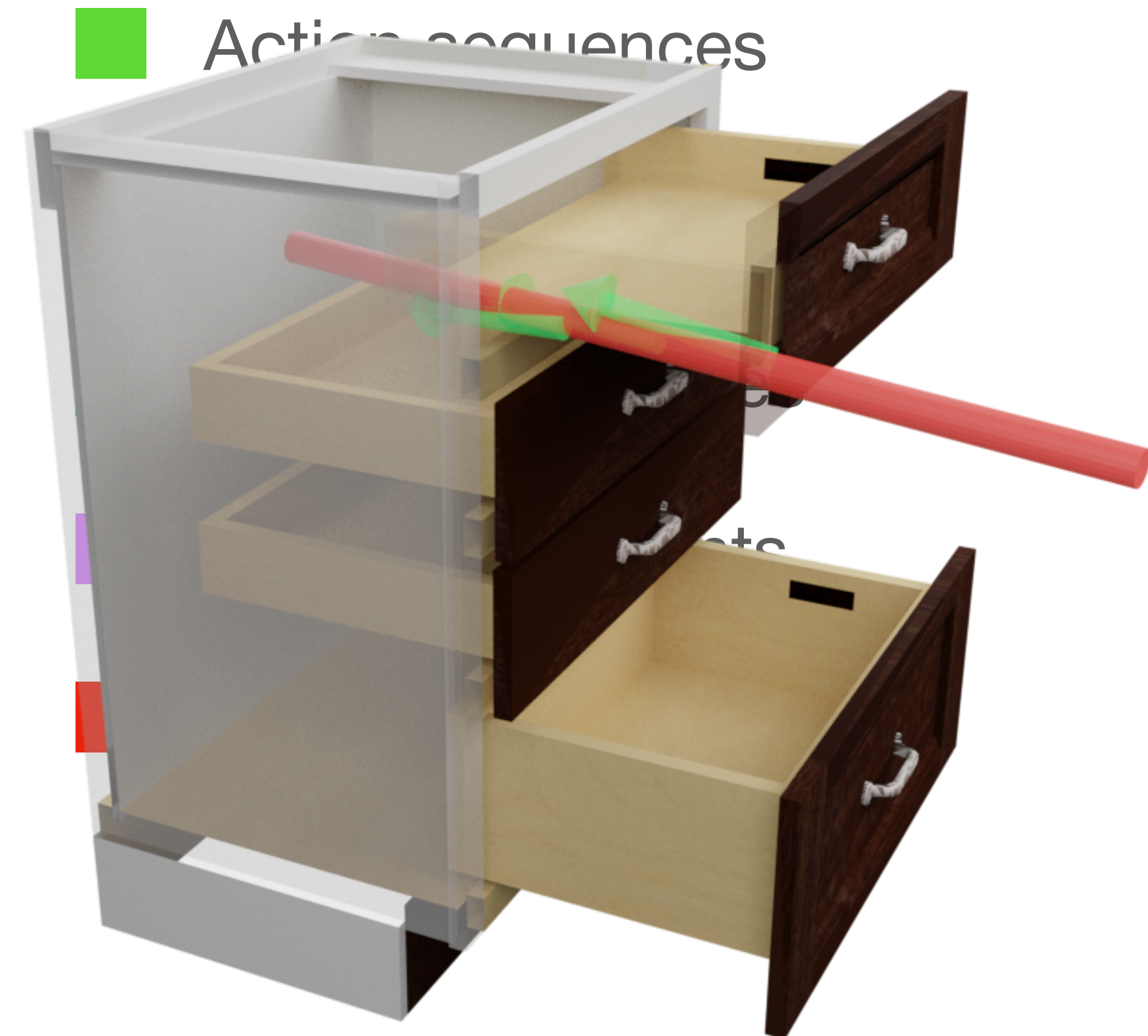
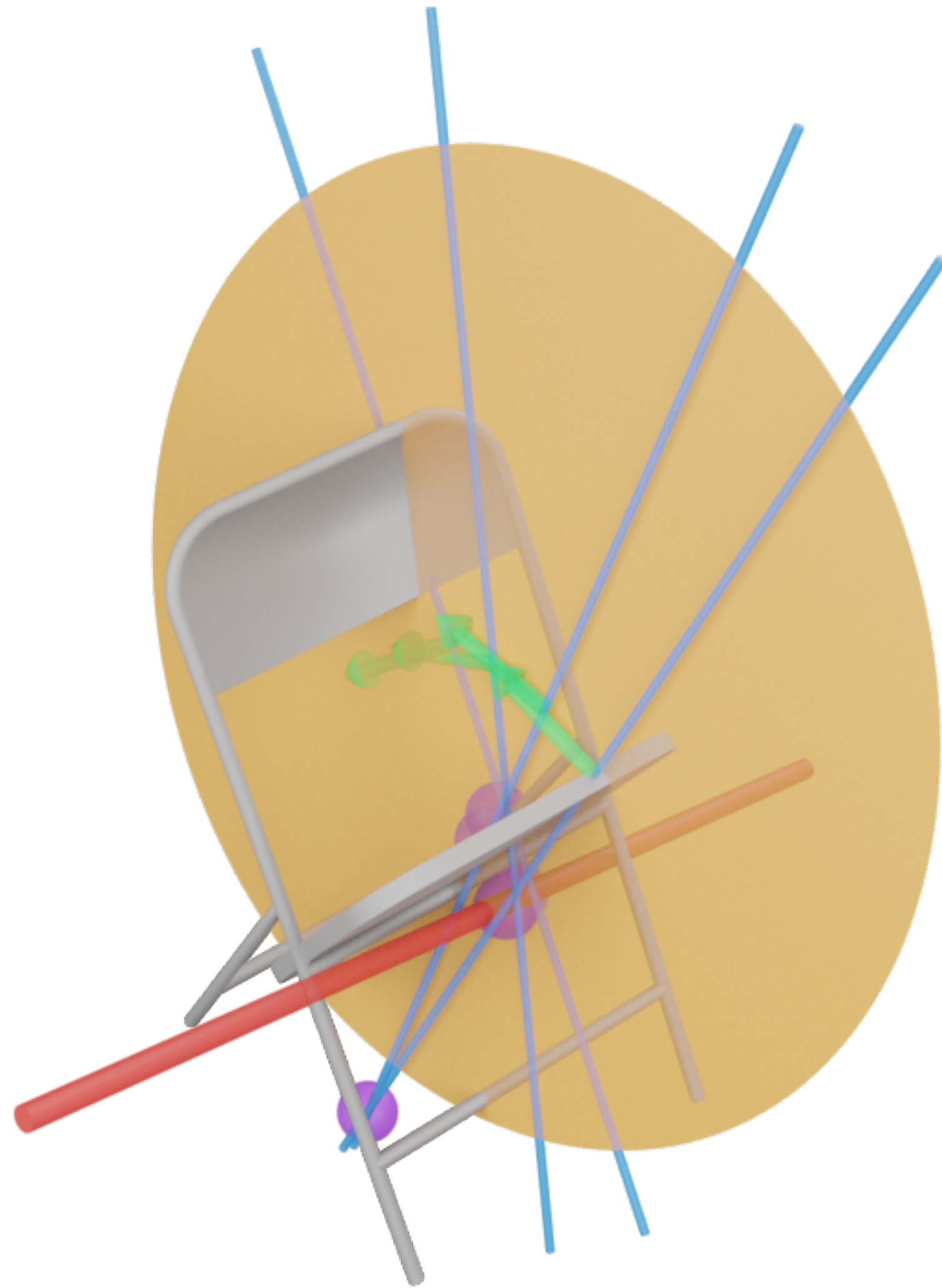
Joint parameters inferred from the  
selected actions



Action selected by the policy should reflect  
its belief on the objects' structure.



# Action → Structure

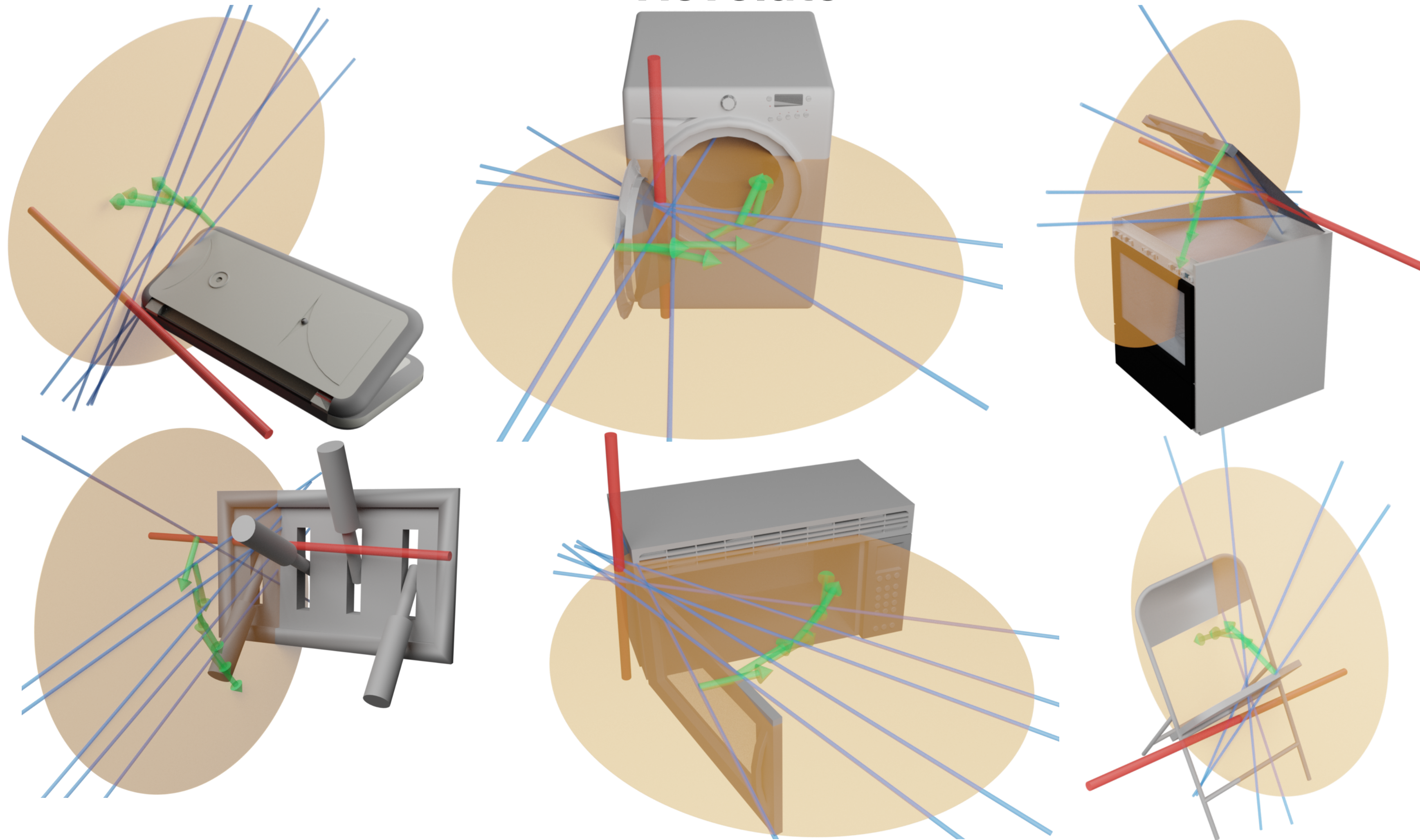


Compute the joint parameters inferred from the actions selected by the policy

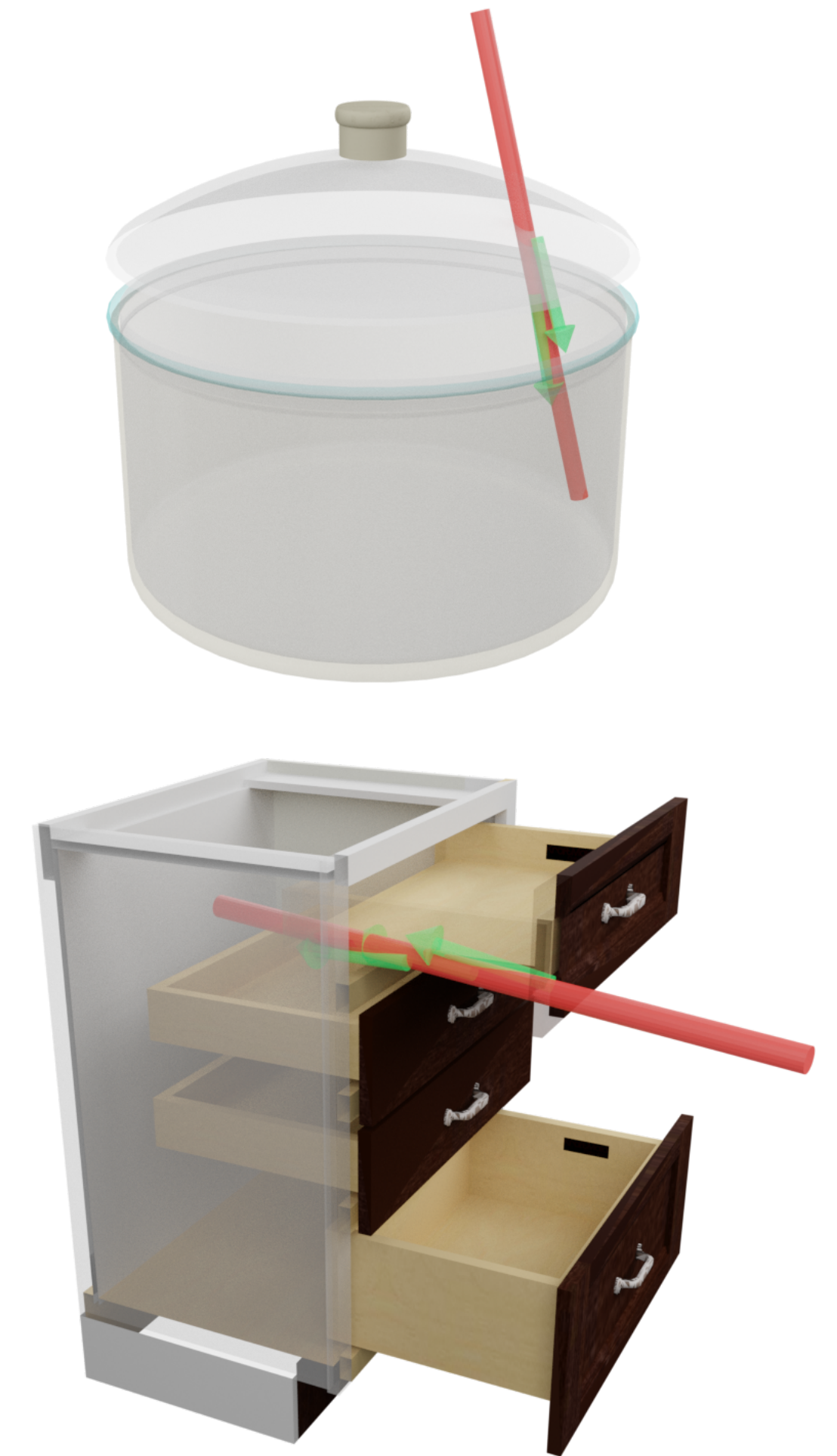


# Articulation Structure Inference

Revolute



Prismatic Joint

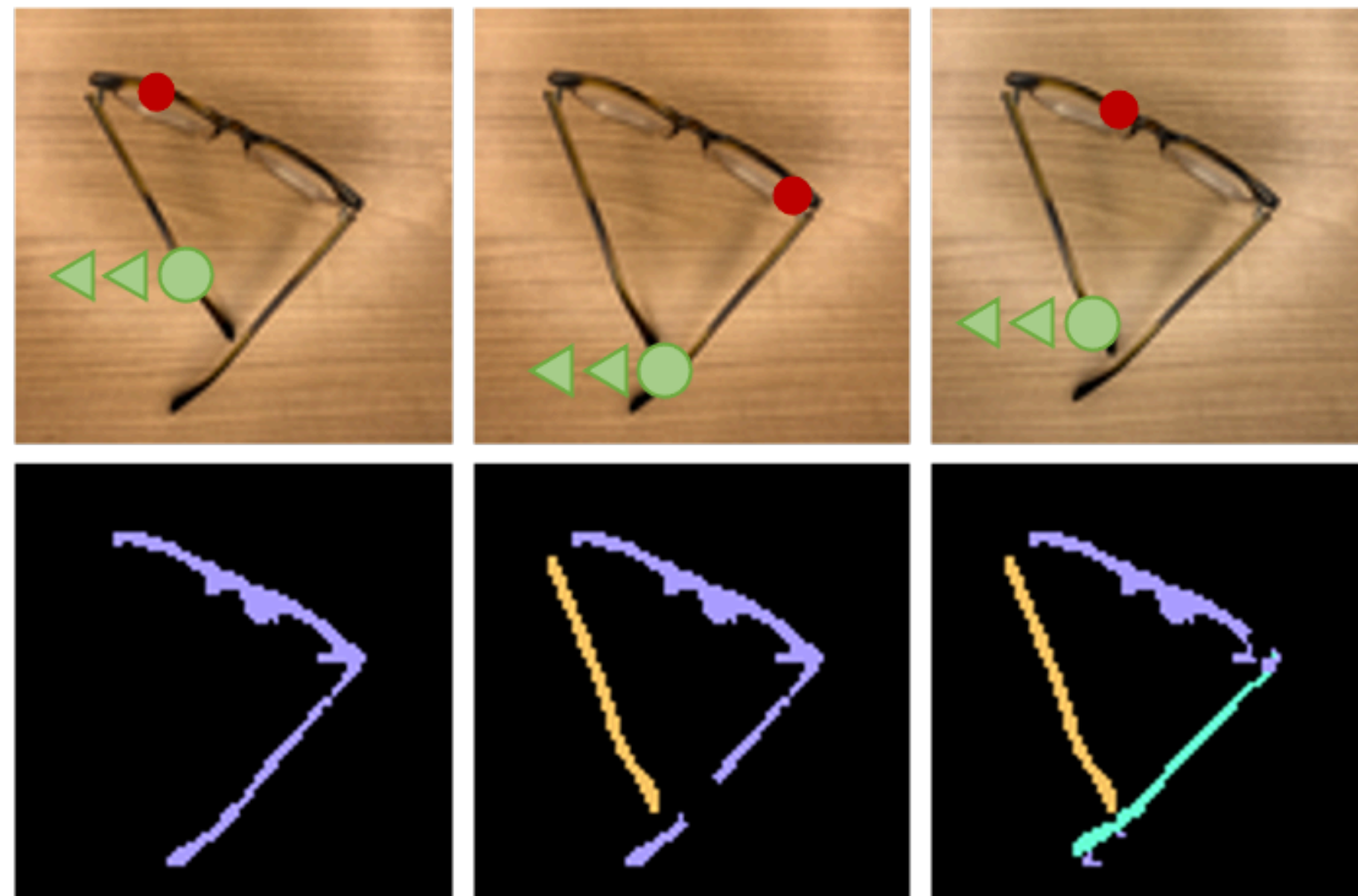


Project Webpage: [ump-net.cs.columbia.edu](http://ump-net.cs.columbia.edu)

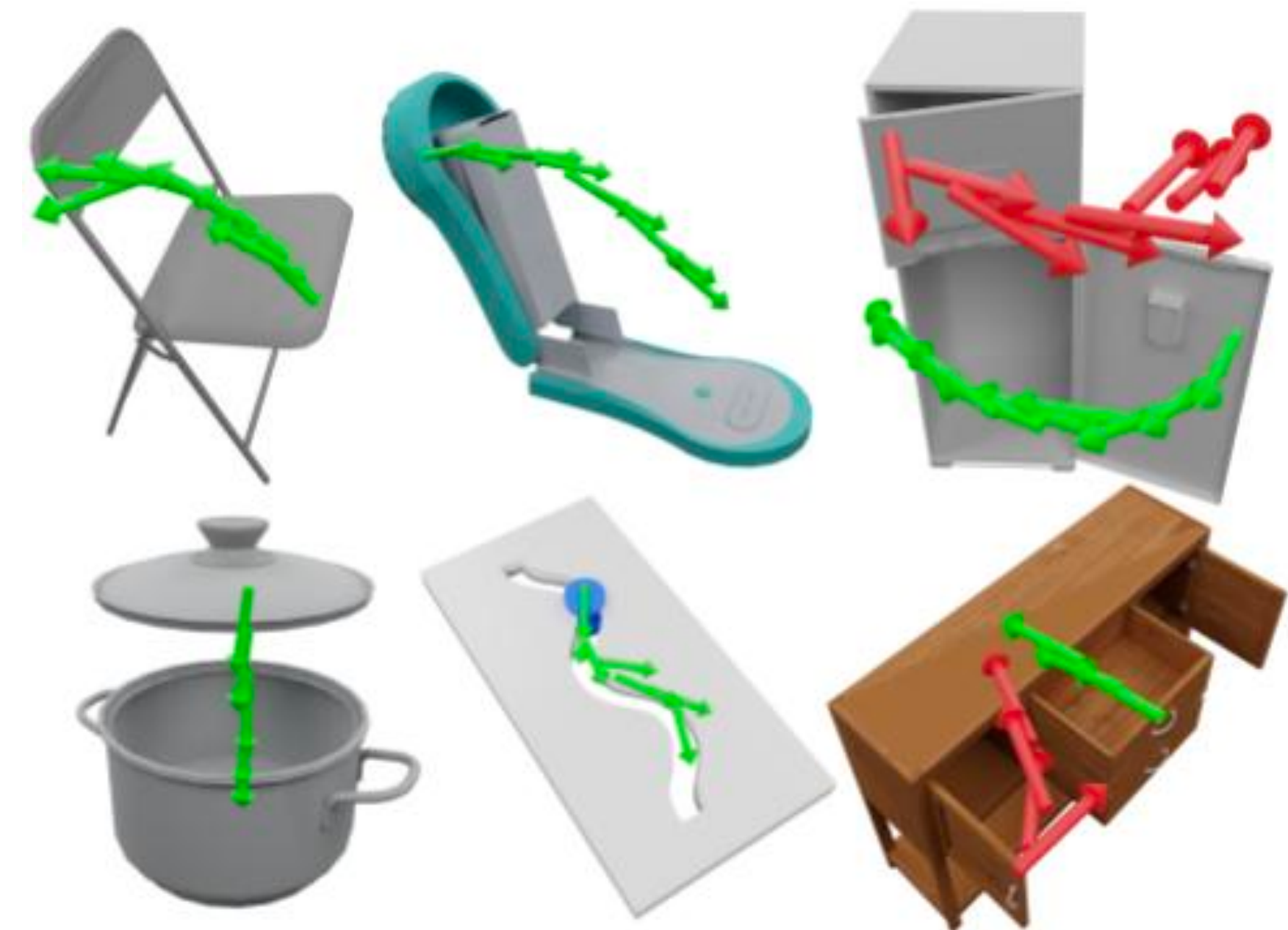


# Structure from Action

**Act the Part:** Learning to Interact to Discover Articulated Object Structure



**UMPNet:** Universal Manipulation Policy Network for Articulated Objects



Underlying structure of object through interaction

Generalize beyond a specific object instance or category



# Acknowledgements

**Act the Part: Learning to Interact to  
Discover Articulated Object Structure**

*Samir Y. Gadre, Kiana Ehsani, Shuran Song*



**Samir Y. Gadre**



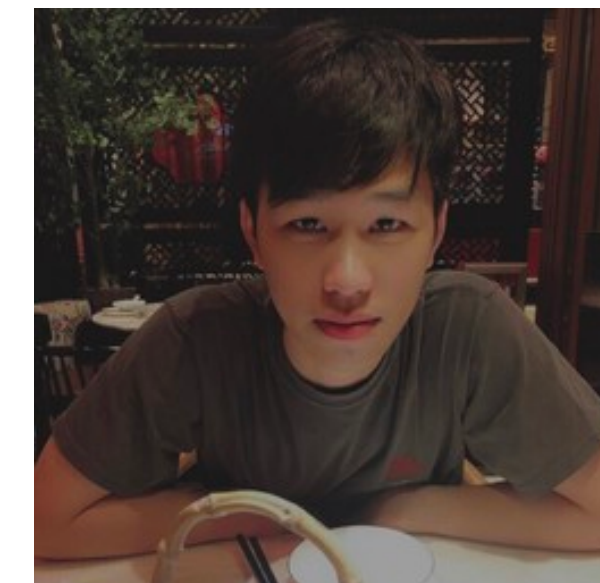
**Kiana Ehsani**

**UMPNet: Universal Manipulation Policy  
Network for Articulated Objects**

*Zhenjia Xu, Zhanpeng He, Shuran Song*



**Zhenjia Xu**



**Zhenjia Xu**



**Google** amazon.com



**Thank You!**

# What's next



