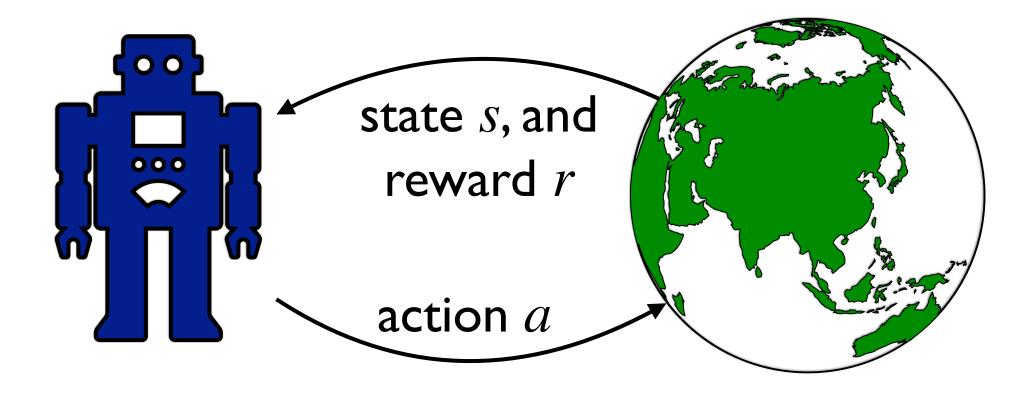


### Overview

### Policy learning from interaction is challenging Videos can aid

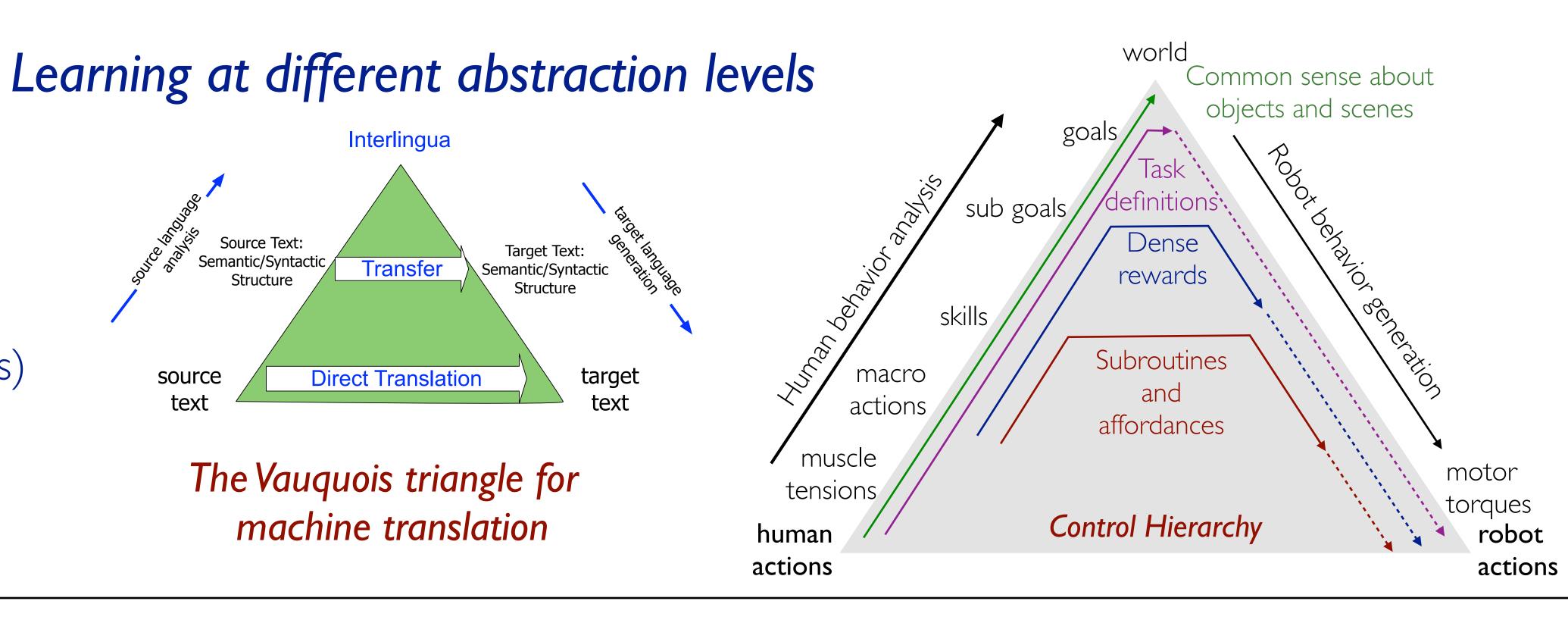


- Challenging to specify reward functions
- Impractically large sample complexity
- Learning signal derived solely from interaction
- Poor generalization due to lack of visual diversity in training, sim2real transfer

### However,

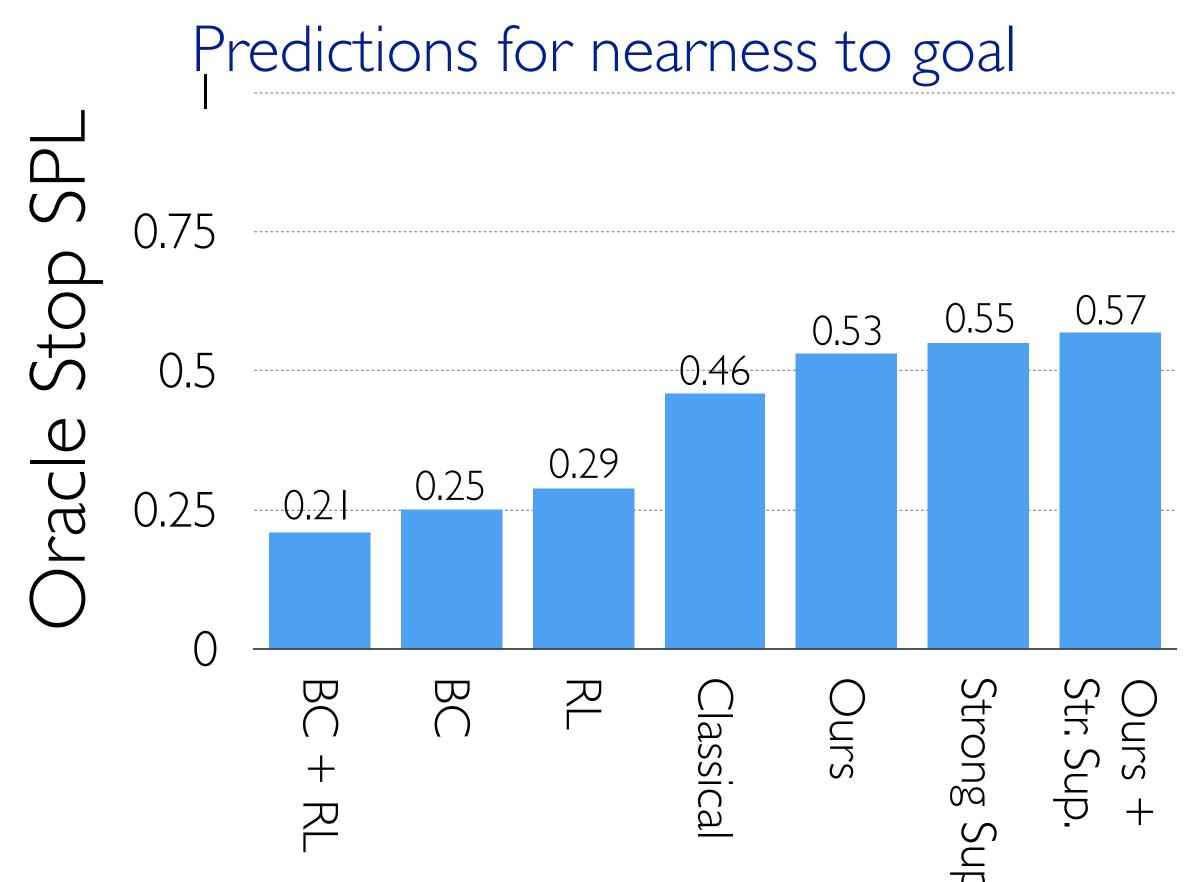
- Videos don't come with action labels
- Goals and intents are not known
- Depicted trajectories may be sub-optimal
- Embodiment gap (sensors / actions / capabilities)
- Only showcase positive data

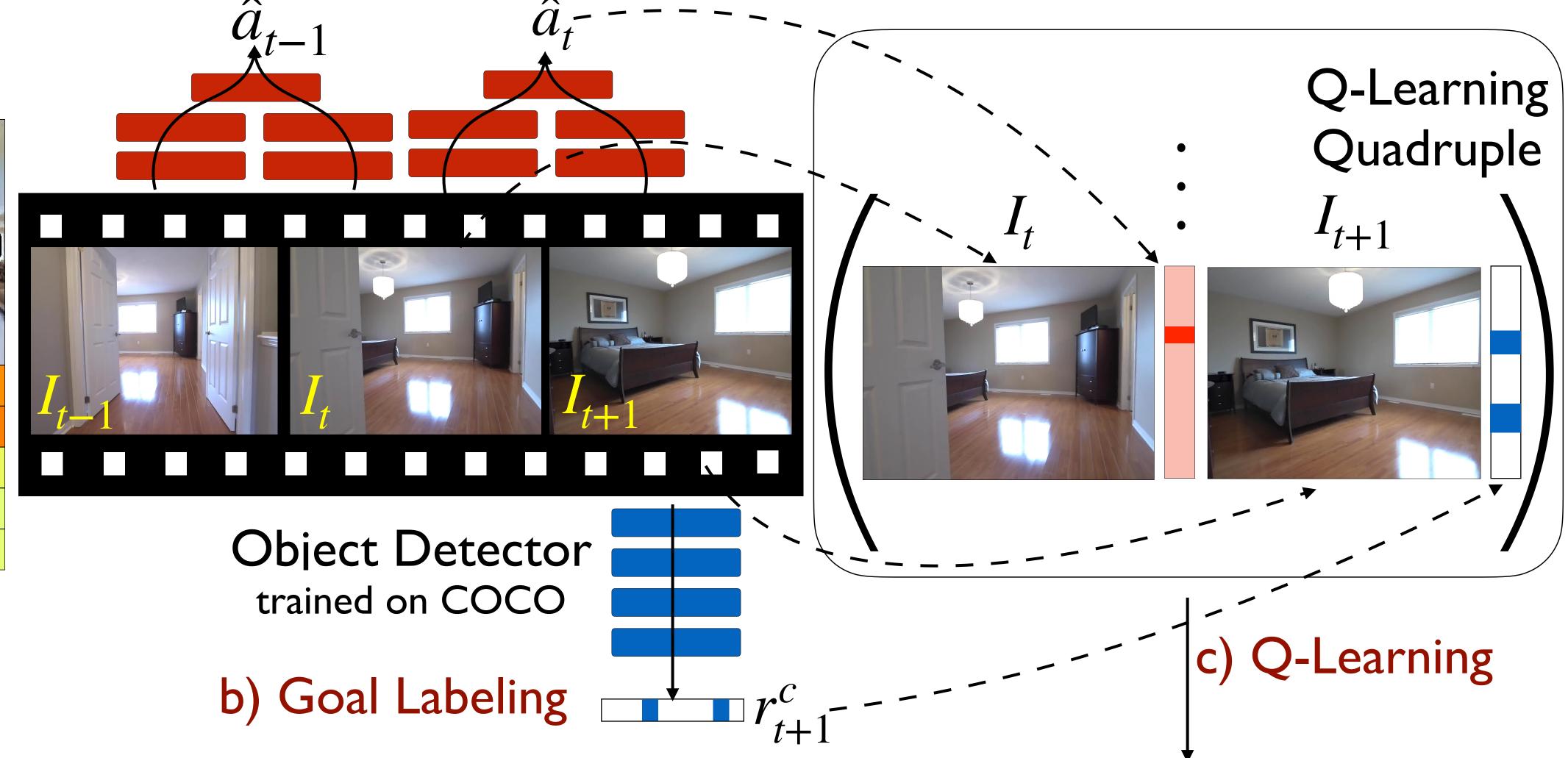




## A.Value Learning from Videos







- Better than strong exploration baselines

# Scaling up Robot Learning by Understanding Videos Saurabh Gupta

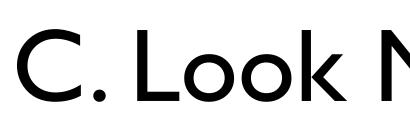
## B. Human Hands as Probes for Interactive Object Understanding

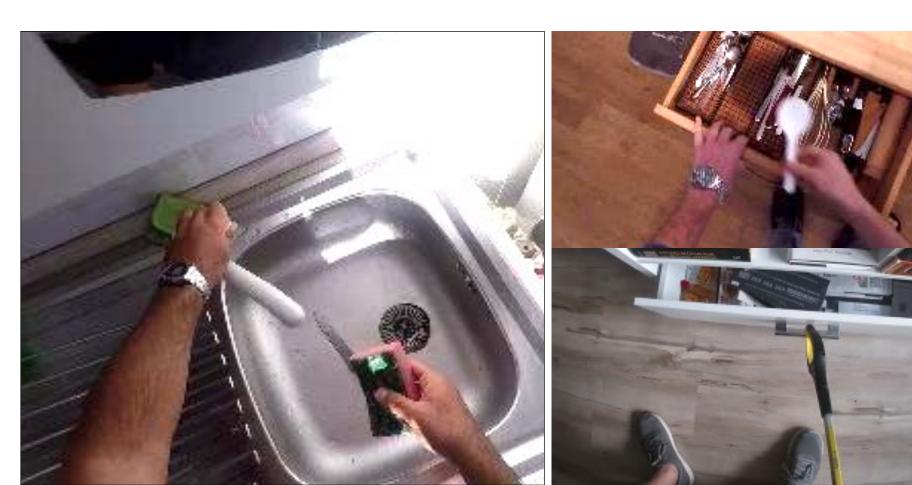
We develop interactive object understanding from the natural ways in which people interact with objects in egocentric videos.

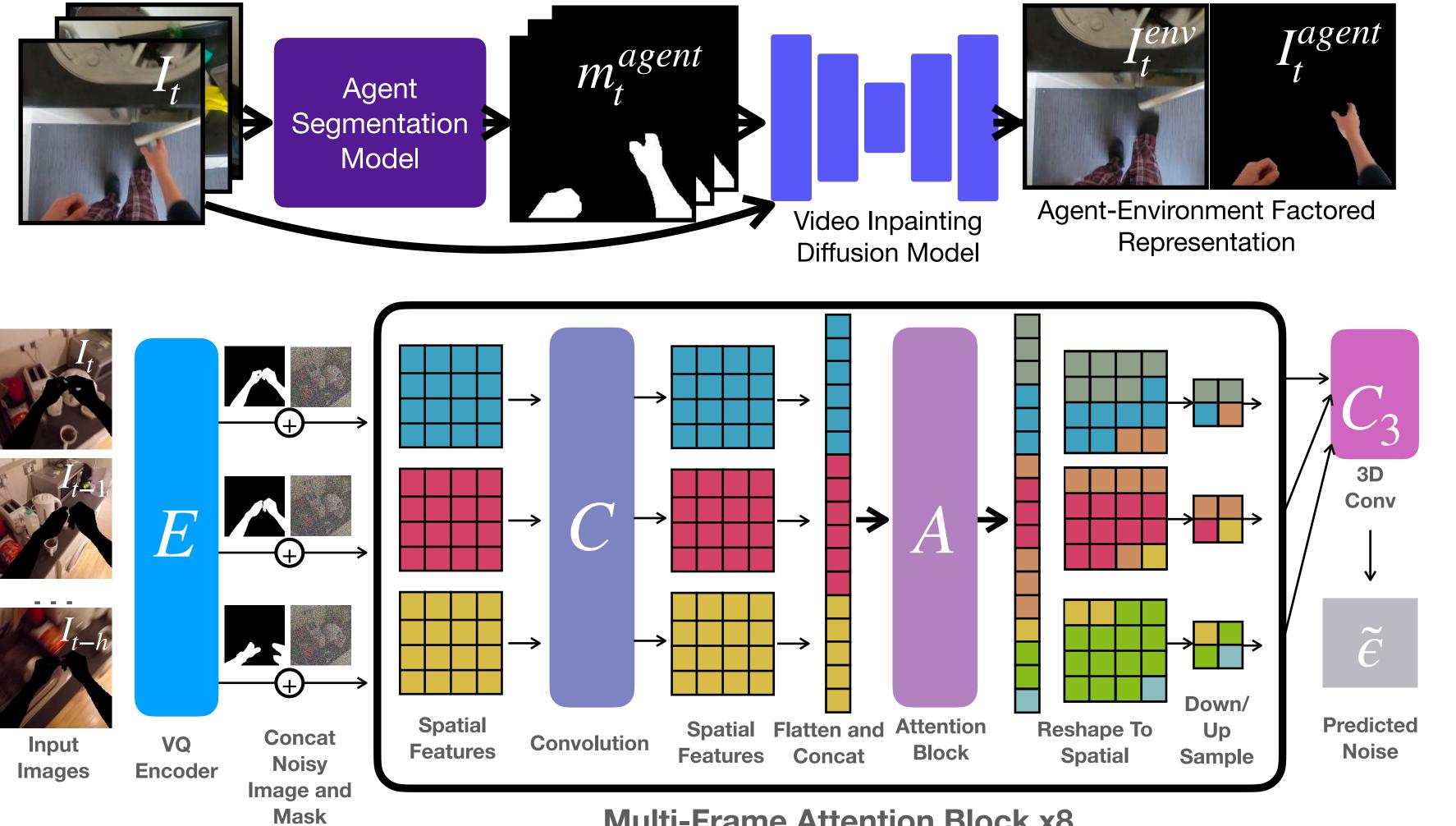


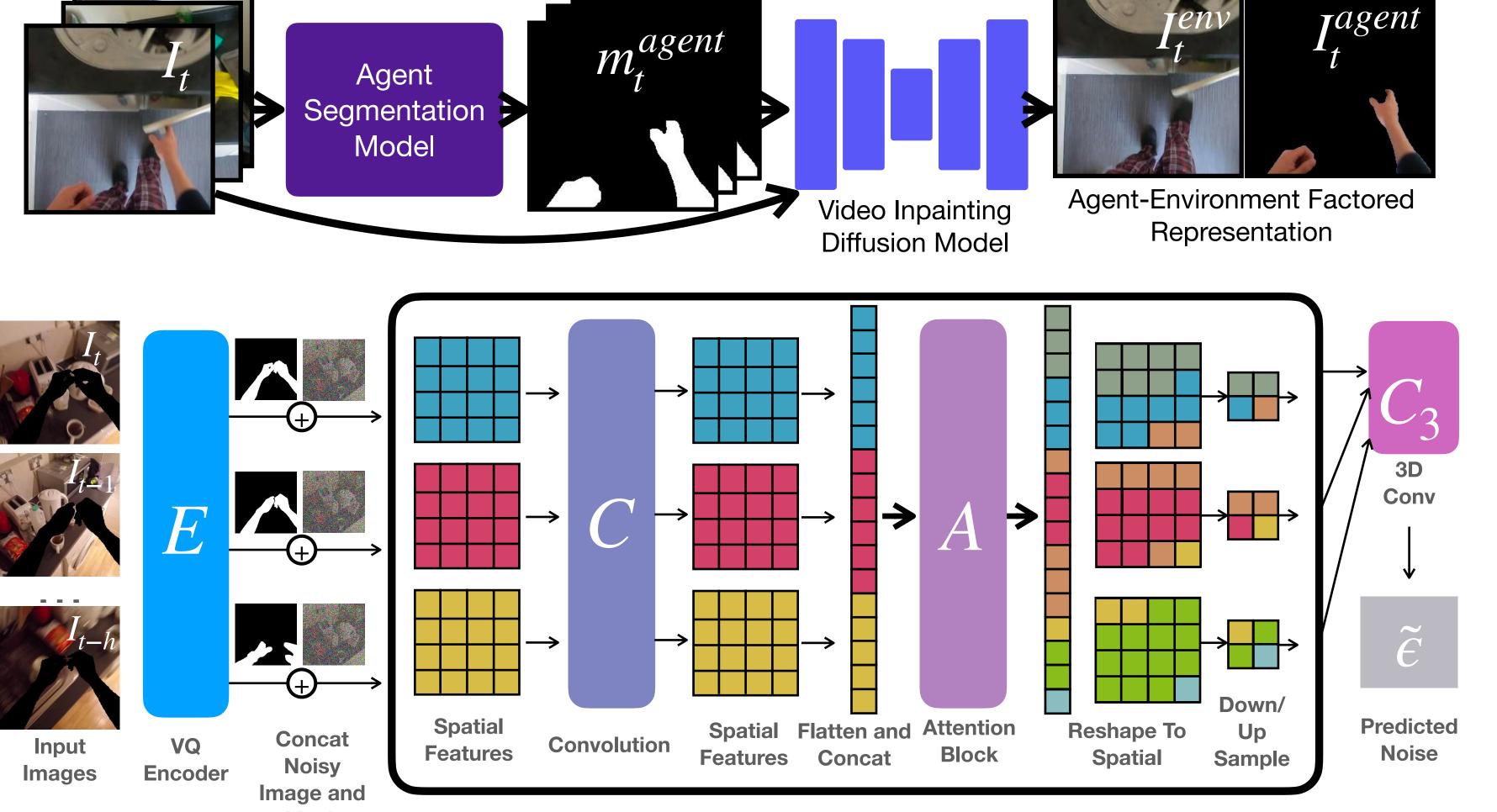
- I. Which sites can we interact at? -(cupboard handles) 2. How to interact with those sites (using adducted thumb grasp) 3. What happens when we do? 1
- (the cupboard opens)





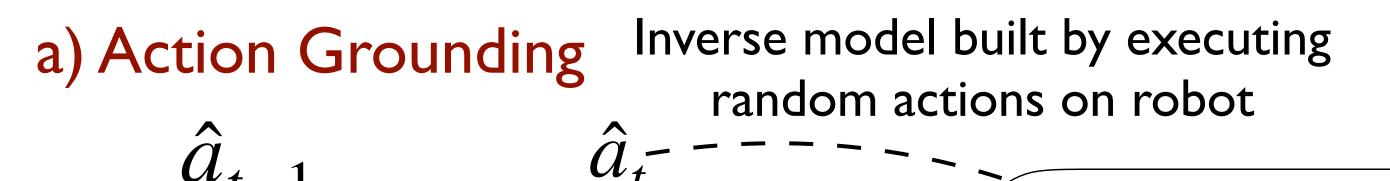








• Large diversity may provide good generalization. • Demonstrations may directly show how to solve long horizon tasks. • Depict what the world is like, and how it works.

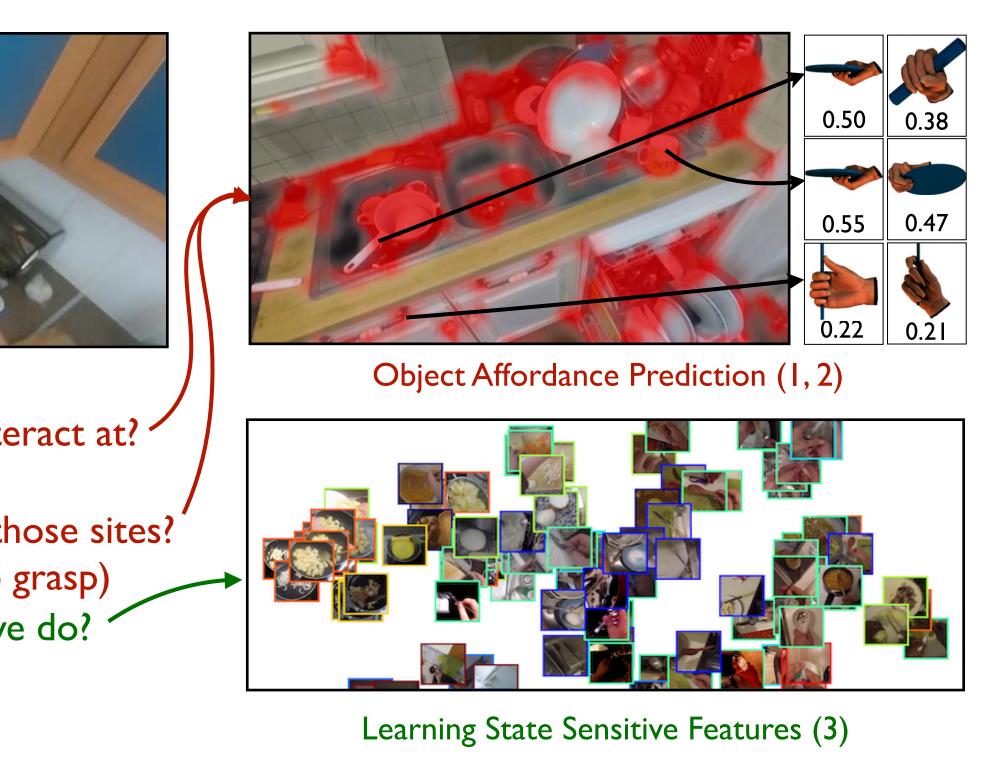


Value function that predicts nearness to goal:  $f(I, c) = \max Q^*(I, a, c)$ 

• Stronger than behavior cloning on videos and BC + RL

• Stronger than even RL methods trained with dense rewards with 250x more interaction samples and 6x more environments with direct interaction access

Improves performance when combined with strongly supervised model

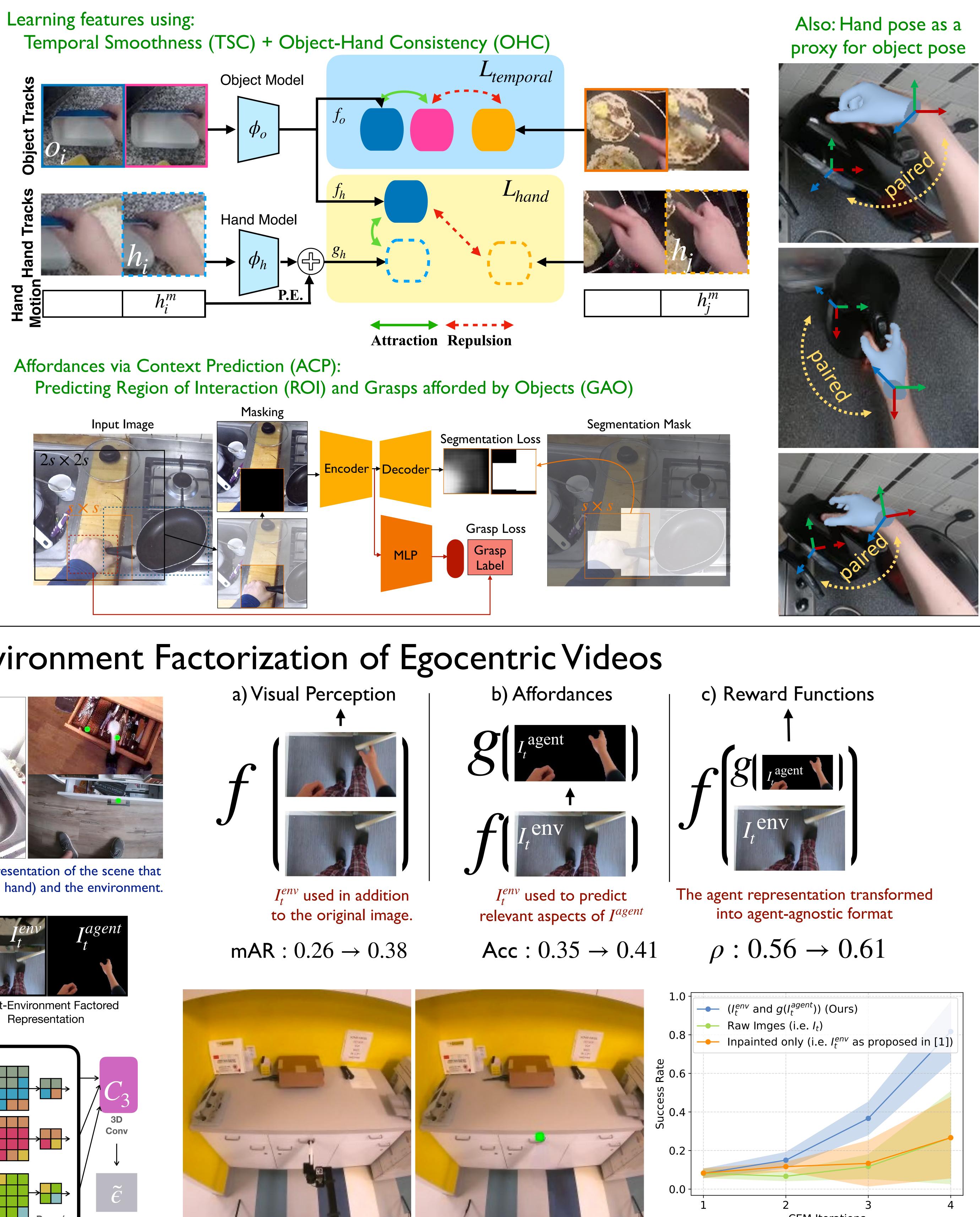


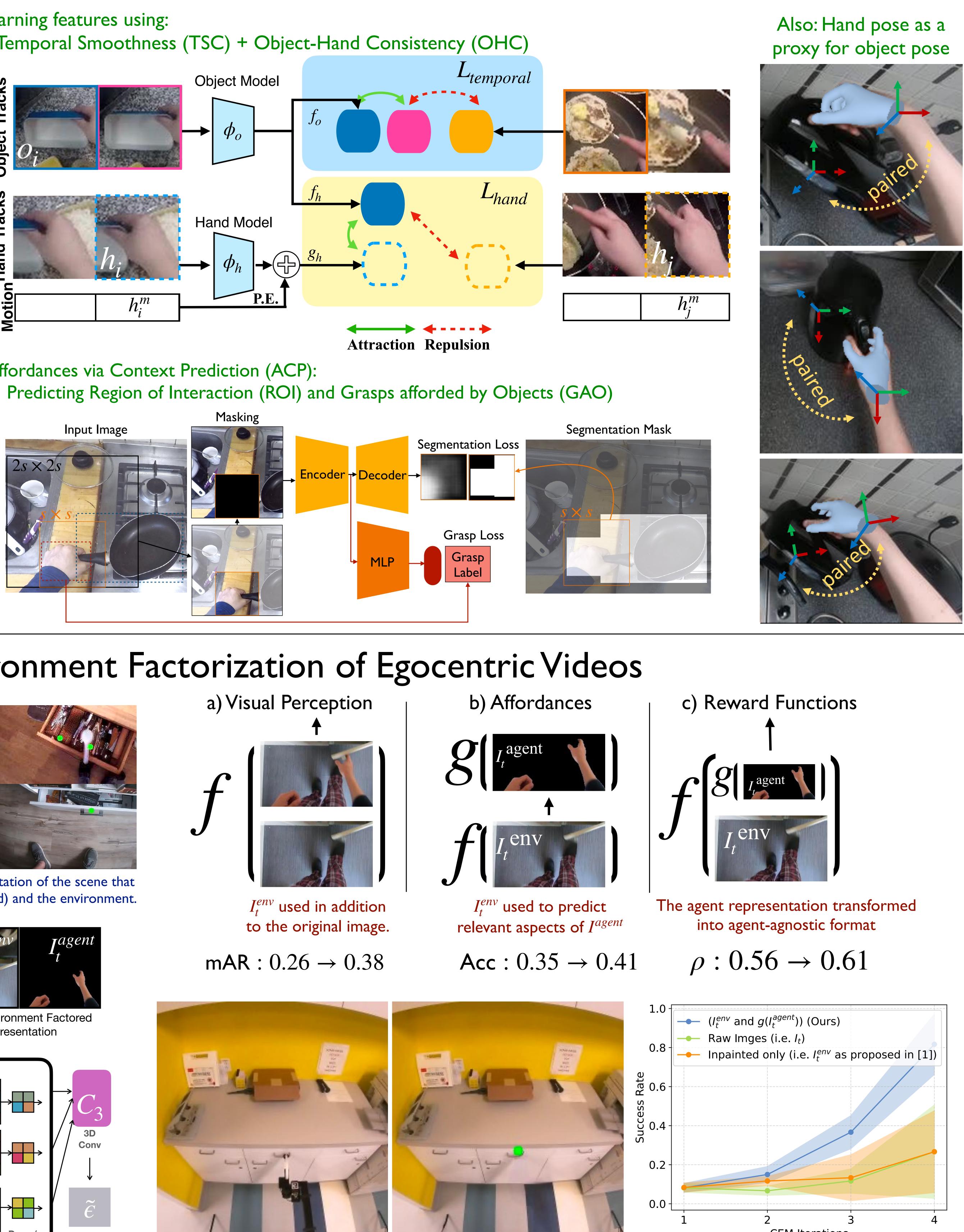
Human hands provide cues for interactive object understanding

Attending to hands localizes and stabilizes active objects.

• Hands show where all we can interact in the scene.

• Analyzing hands reveals information about objects: state and how to interact. • Hand pose may provide hints for object pose.

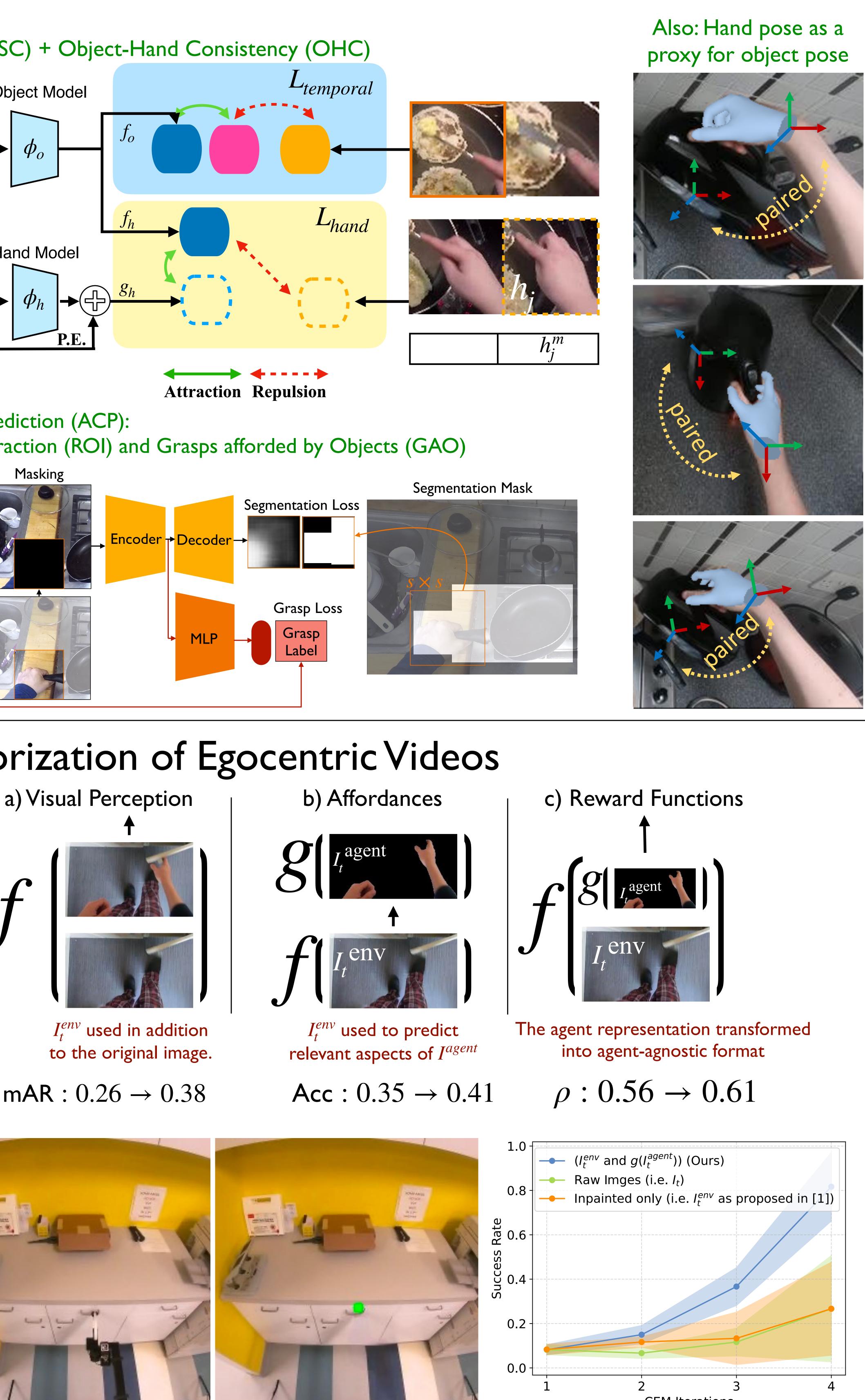


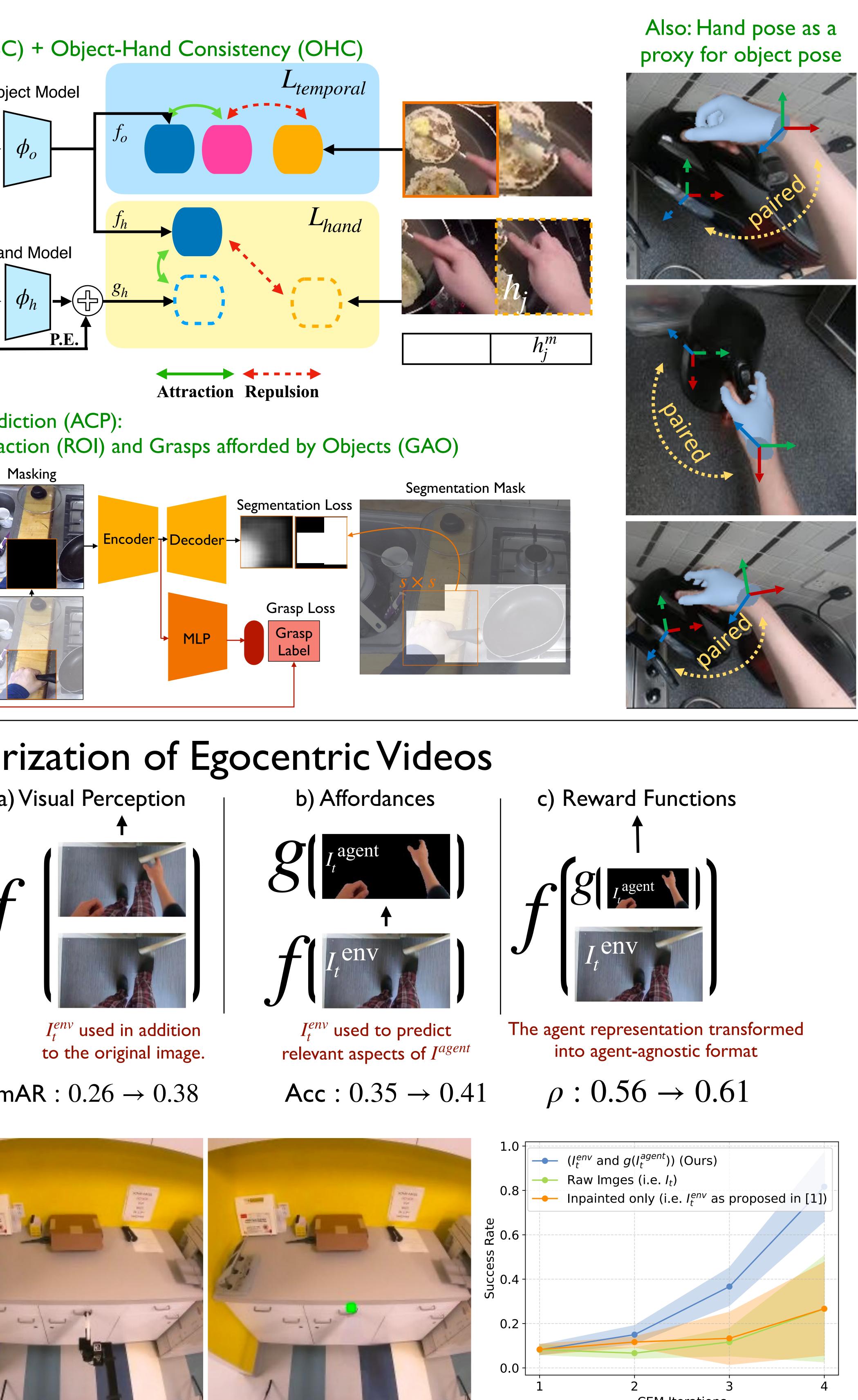


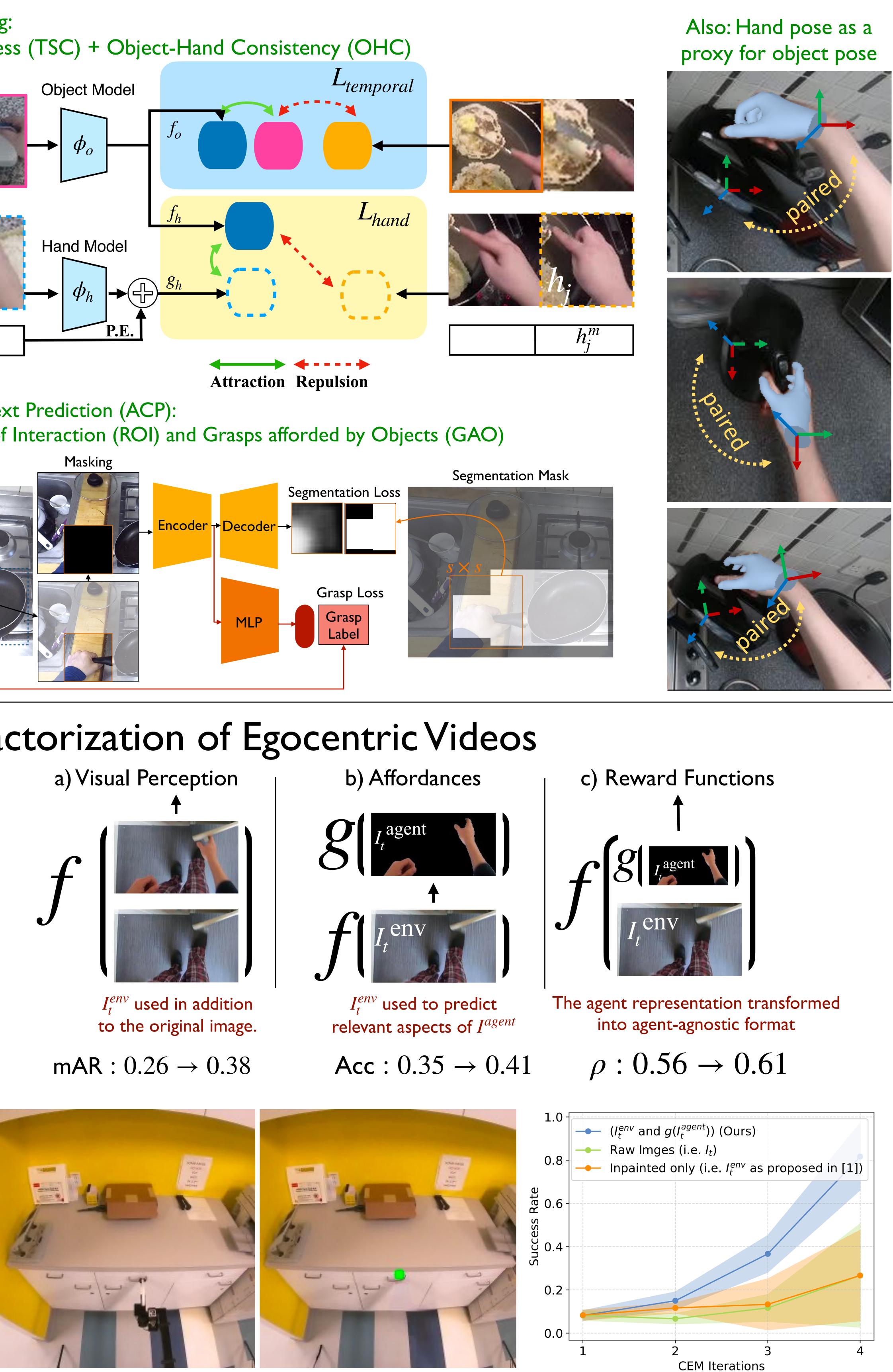
### C. Look Ma, No Hands! Agent-Environment Factorization of Egocentric Videos



We extract a factored representation of the scene that separates the agent (human hand) and the environment.







**Multi-Frame Attention Block x8** 











Sahil Modi



Matthew Chang Arjun Gupta Aditya Prakash Mohit Goyal

Rishabh Goyal

Drawer opening policy trained in the in the real world. We find that using the agentenvironment factored representation improves learning speed.