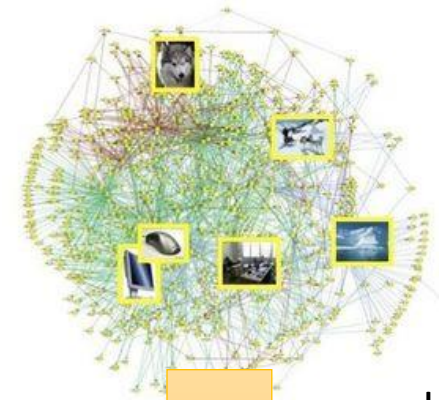
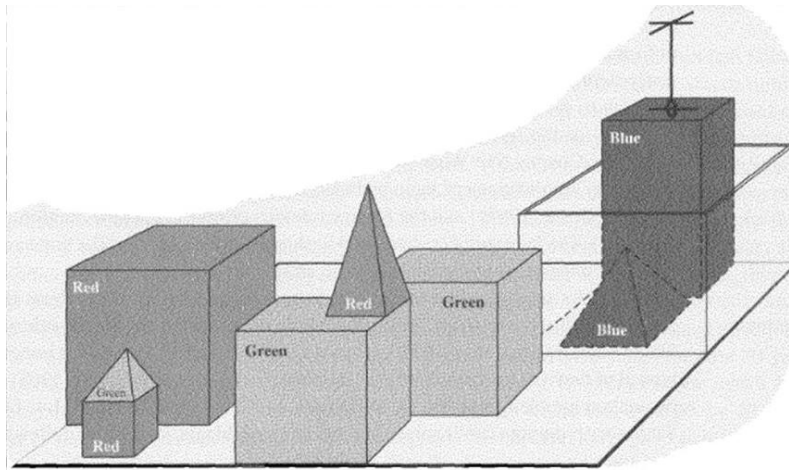


Embodied Implicit Scene Understanding

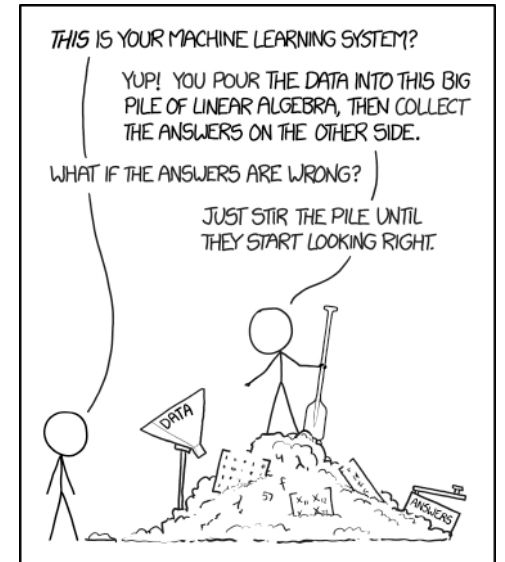
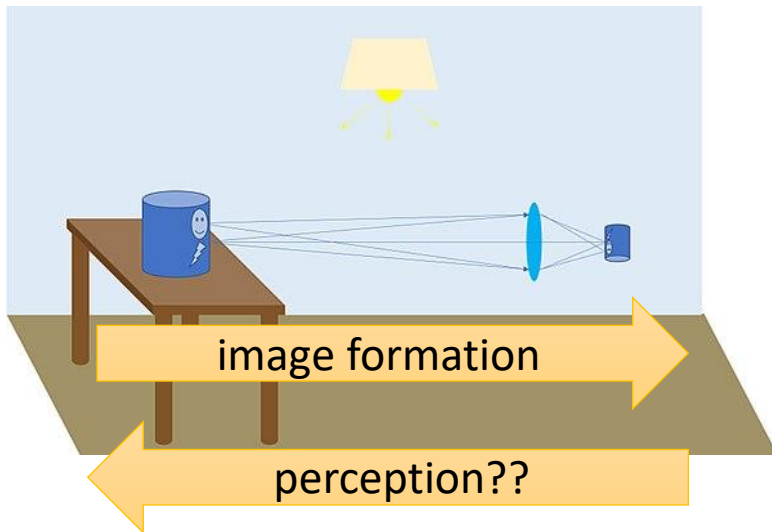
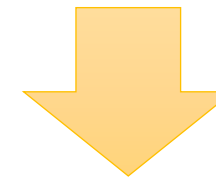
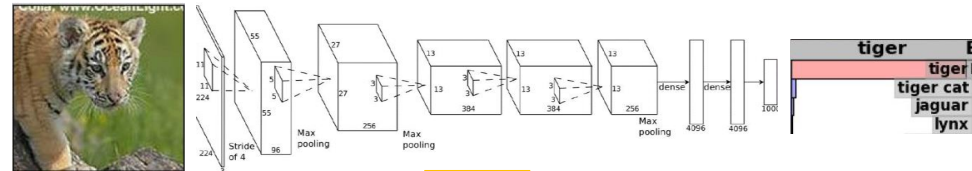
Sergey Levine

UC Berkeley

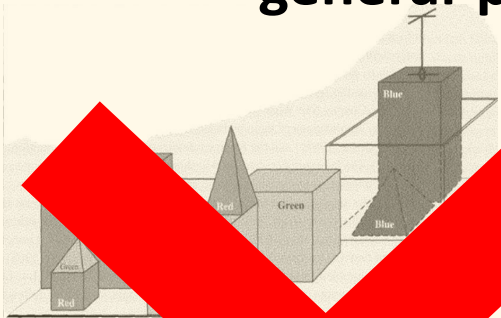




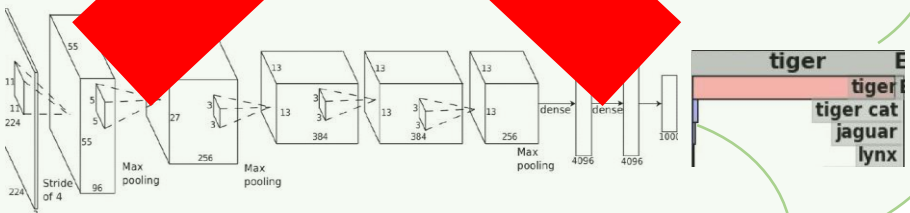
end-to-end training



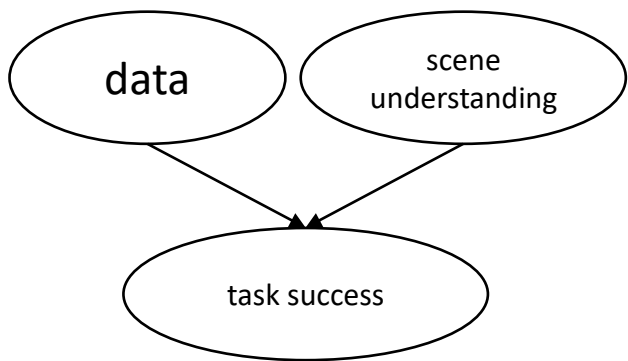
general-purpose?



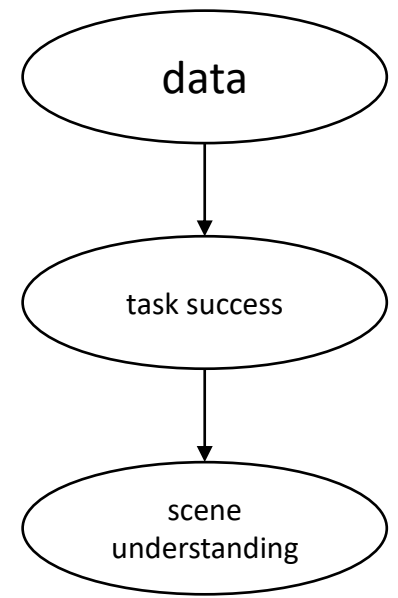
task-specific?



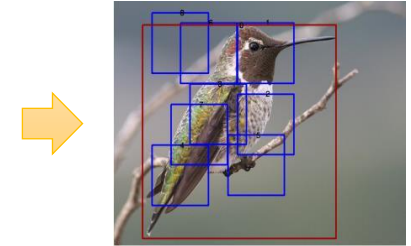
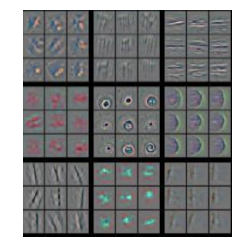
“classic” view



“epiphenomenon view”



end-to-end systems learn **general** concepts!



Understanding the world via end-to-end learning?

Our universe:

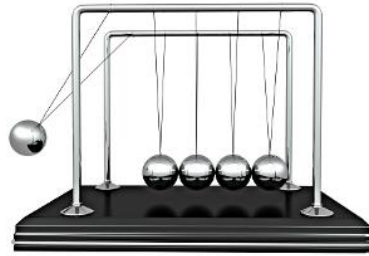
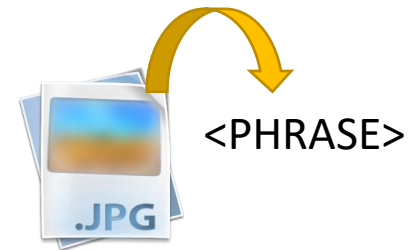
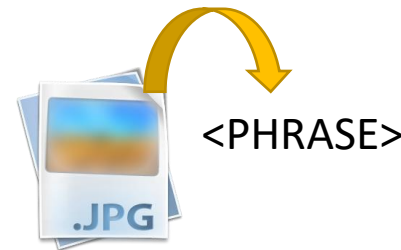
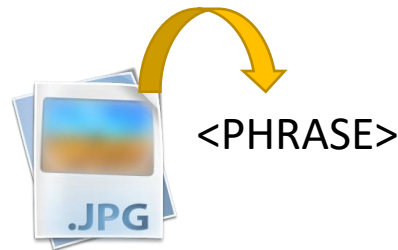
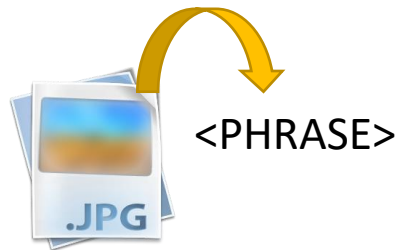


Image captioning
“universe”:



a group of people standing around a room with remotes
logprob: -9.17



a young boy is holding a baseball bat
logprob: -7.61



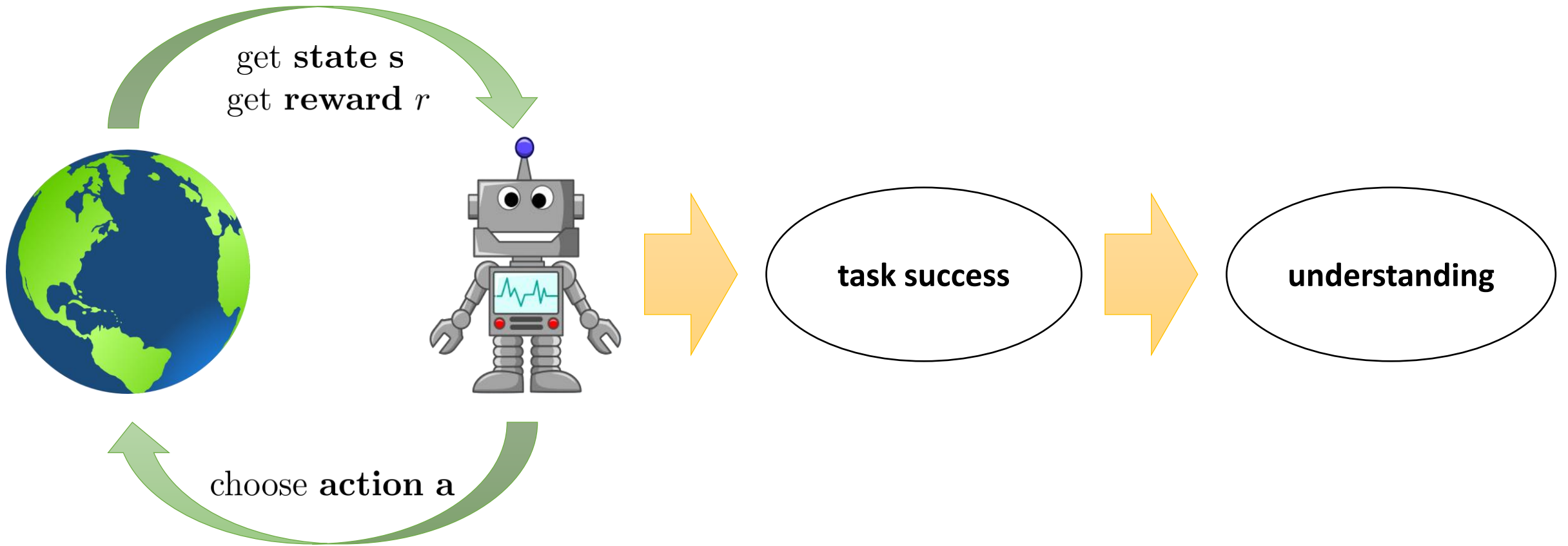
a toilet with a seat up in a bathroom
logprob: -13.44



a woman holding a teddy bear in front of a mirror
logprob: -9.65

Karpathy & Li, 2016

An embodied learning recipe for scene understanding?



Which end-to-end task should we use?

Model-free algorithms:
predict future *rewards*

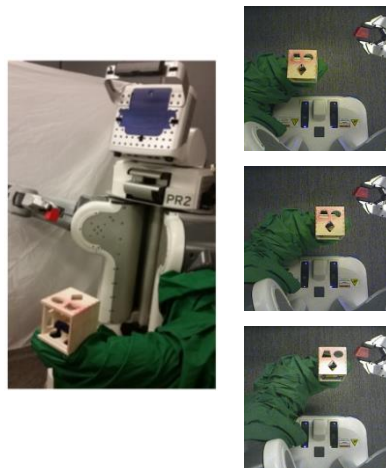
Model-based algorithms:
predict future *observations*

Which end-to-end task should we use?

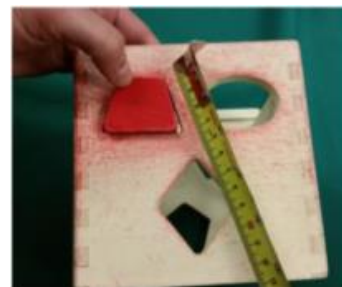
Model-free algorithms:
predict future *rewards*

Model-based algorithms:
predict future *observations*

End-to-end training



network architecture	test error (cm)
softmax + feature points (ours)	1.30 ± 0.73
softmax + fully connected layer	2.59 ± 1.19
fully connected layer	4.75 ± 2.29
max-pooling + fully connected	3.71 ± 1.73

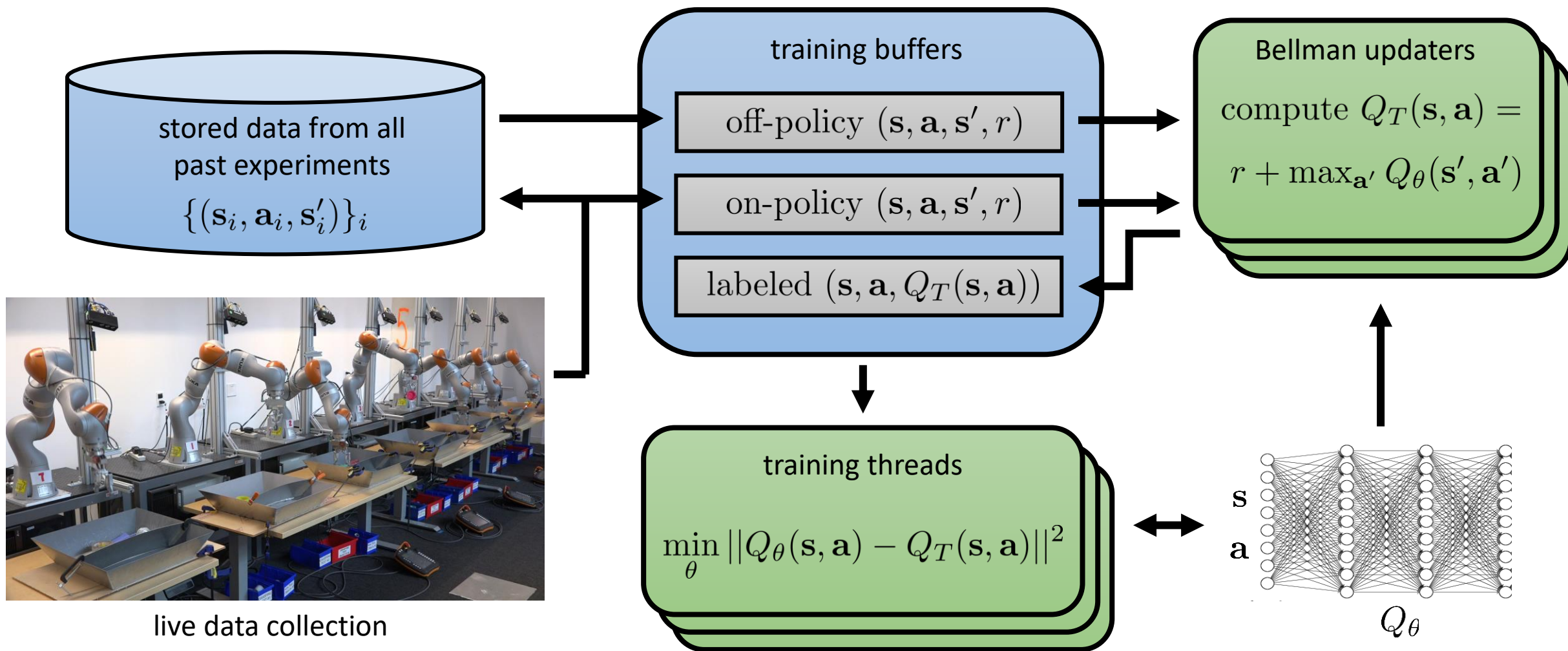


Meeussen et al. (WG)

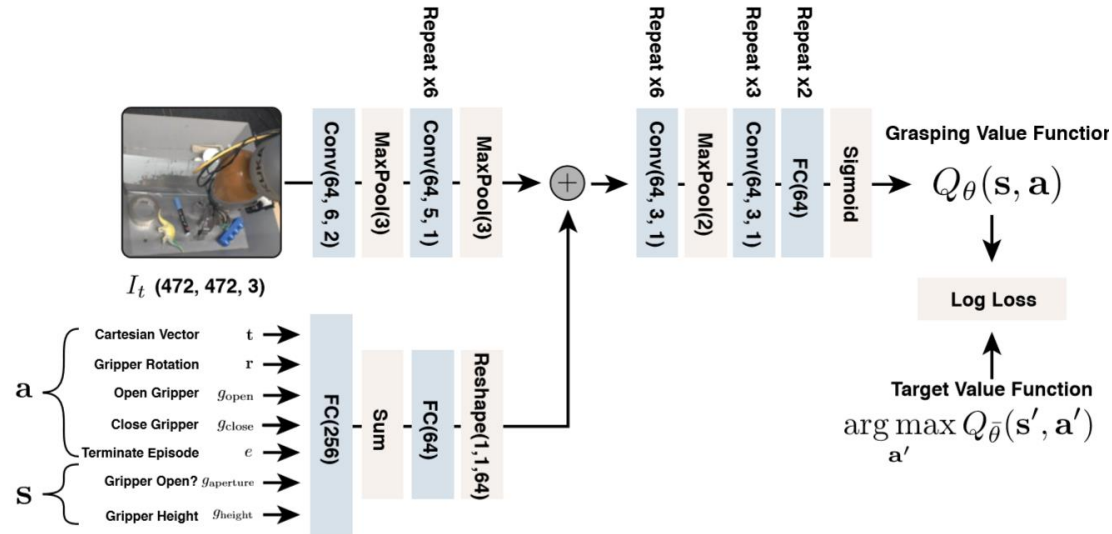
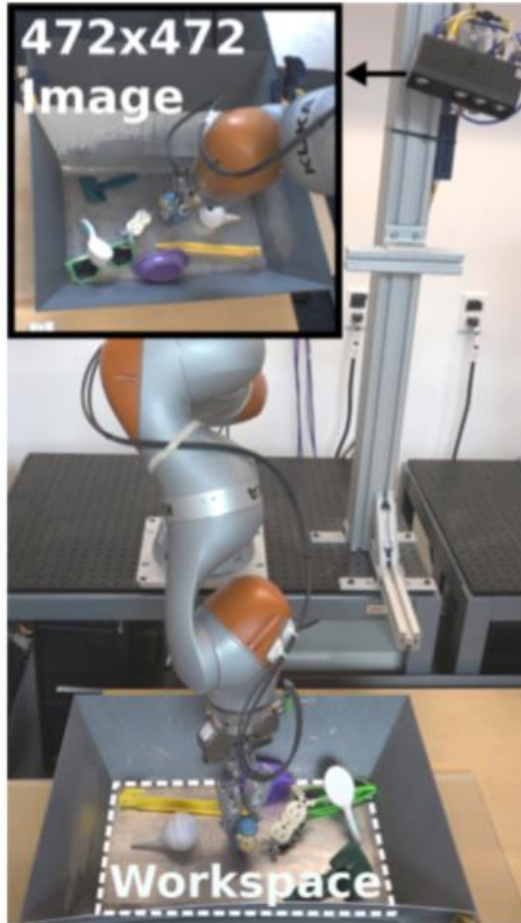


shape sorting cube	success rate
pose prediction	0%
pose features	70.4%
end-to-end training	96.3%

QT-Opt: robotic RL at scale



Grasping with QT-Opt



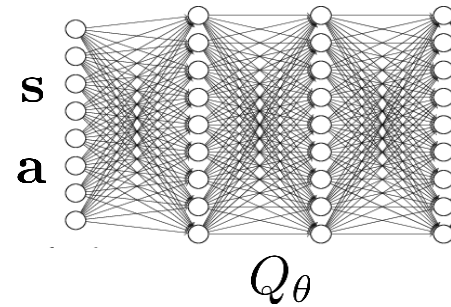
- About 1000 training objects
- About 600k training grasp attempts
- Q-function network with 1.2M parameters
- The only grasp-specific feature is the reward (1 if grasped)

Learning on the job

training



96% success rate



“on the job”

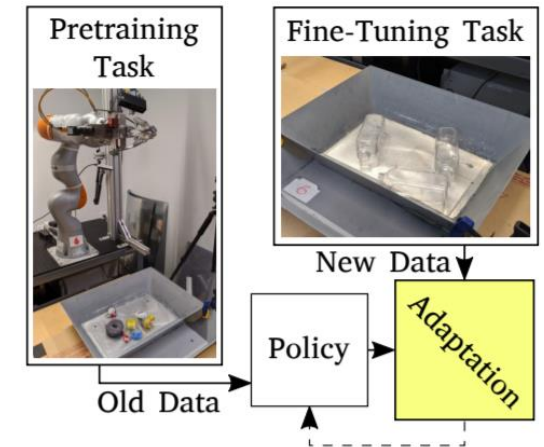


49% success rate

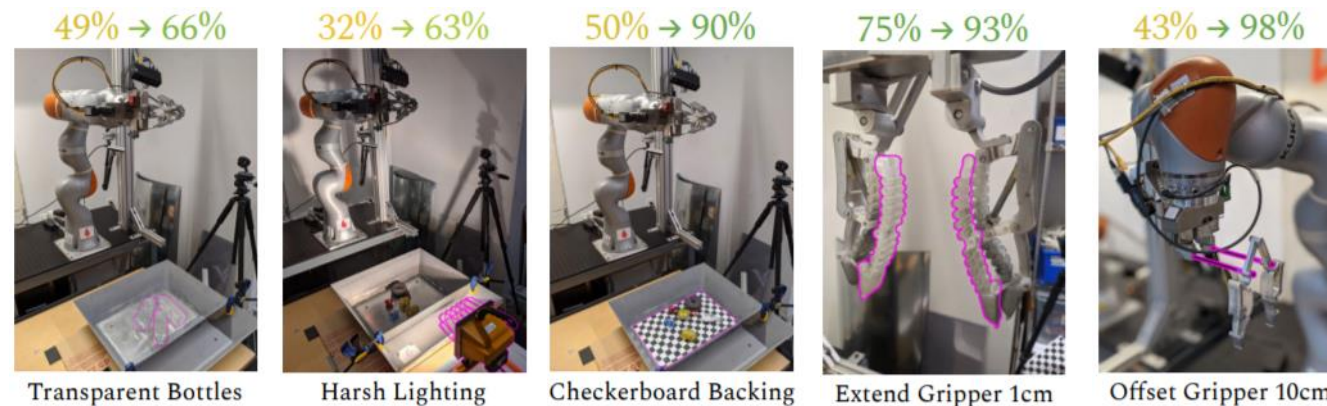
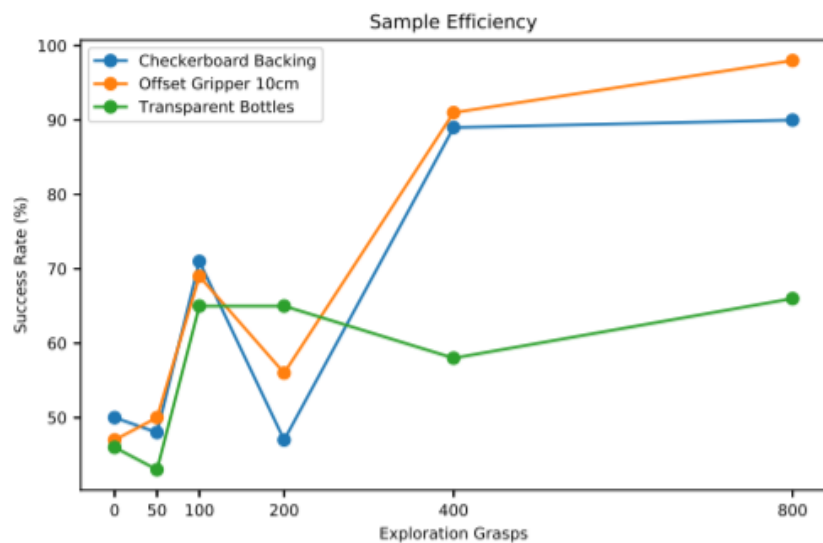
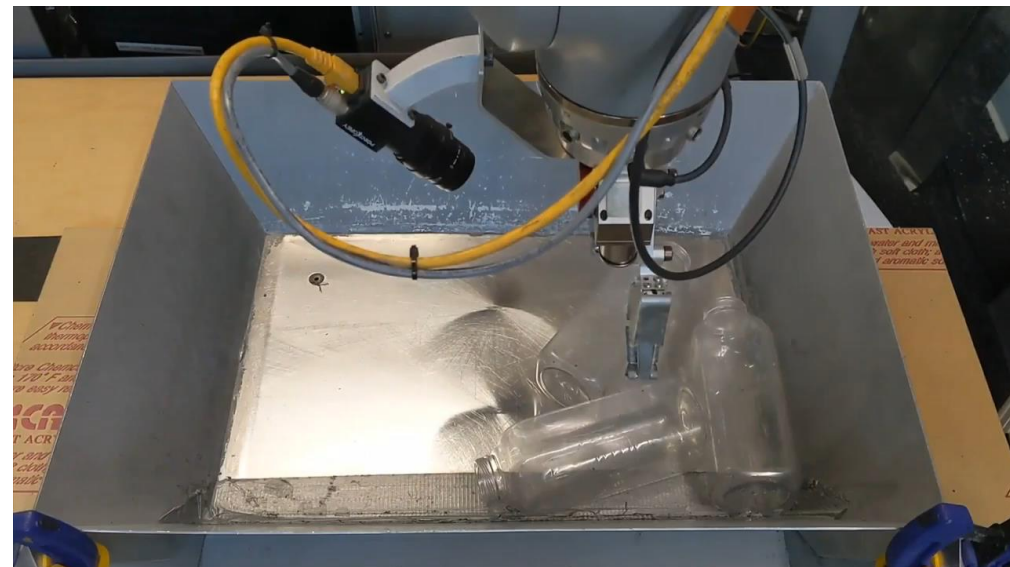
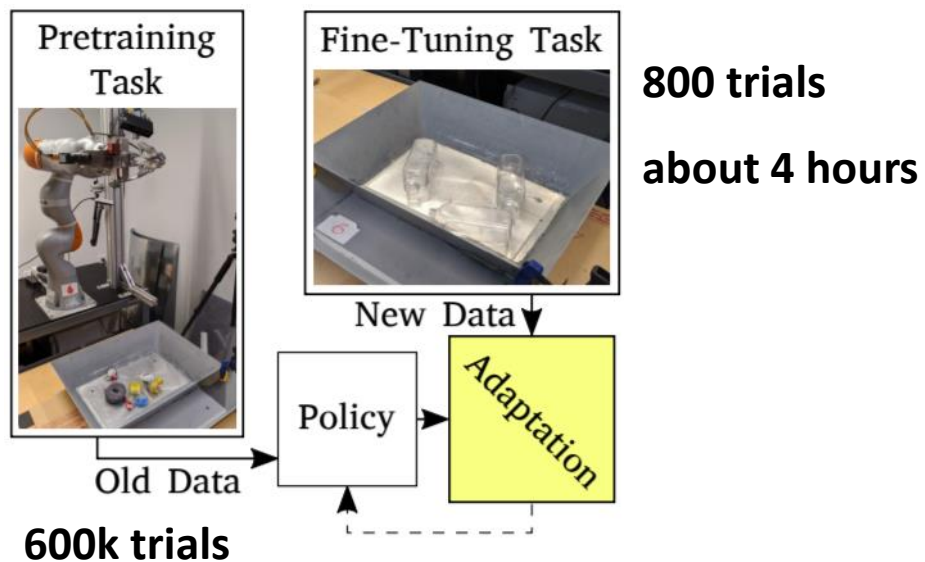


just keep training!

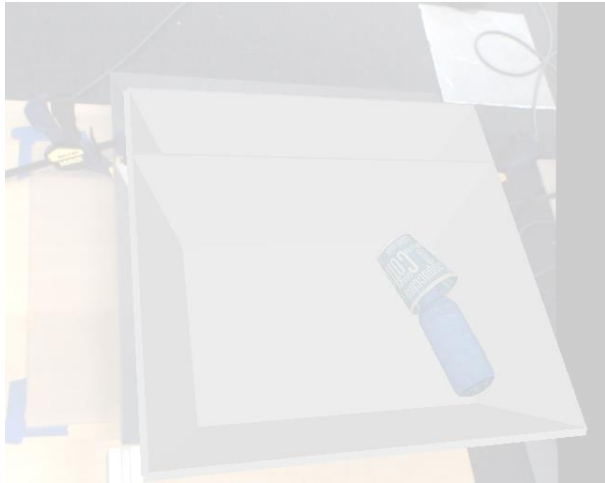
no human effort required



Learning on the job



Grasping provides supervision



pre-grasp scene

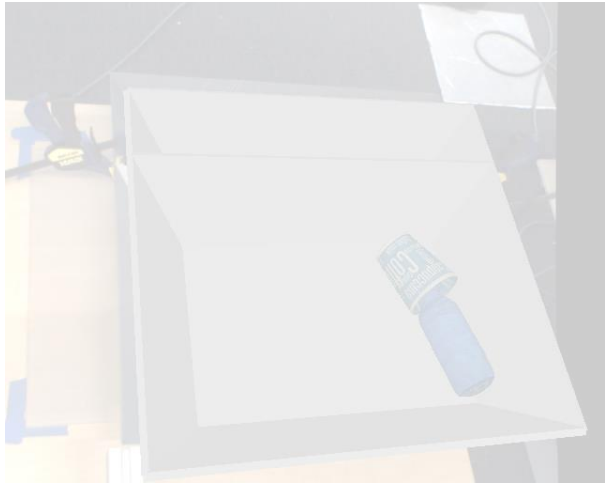


grasped object



post-grasp scene

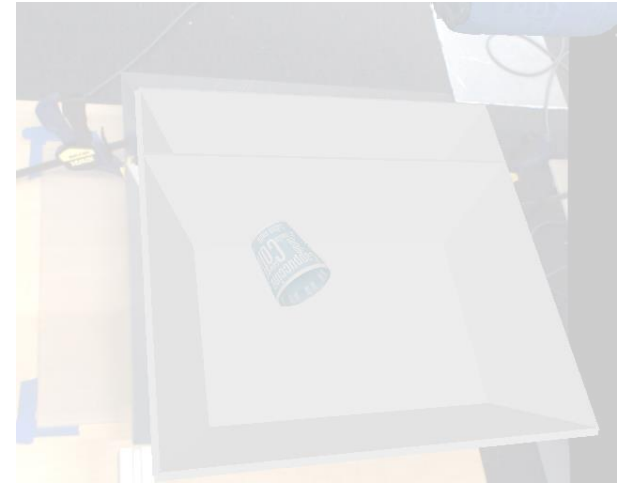
Grasping provides supervision



pre-grasp scene

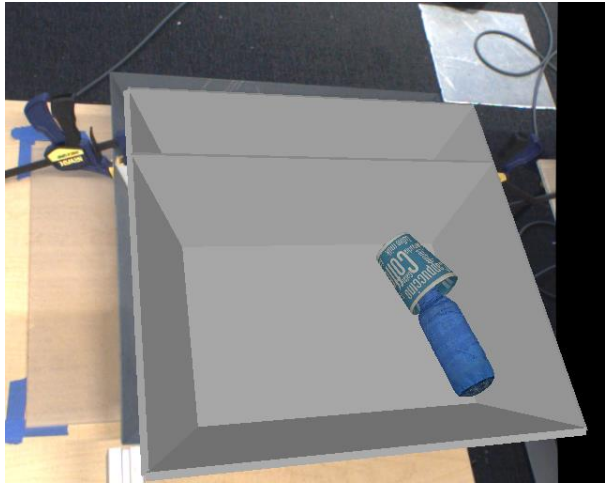


grasped object



post-grasp scene

Grasping provides supervision



pre-grasp scene

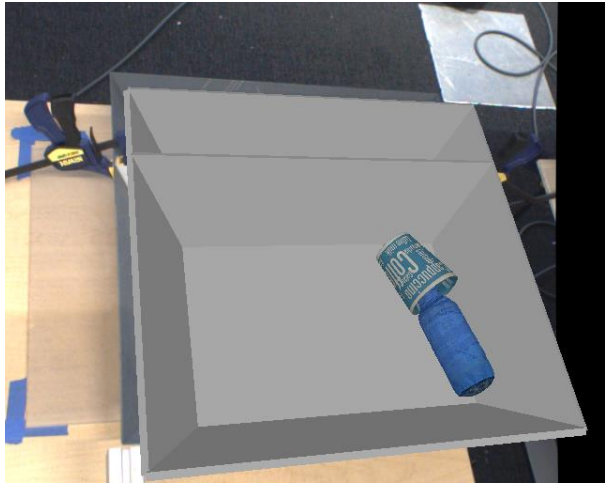


grasped object



post-grasp scene

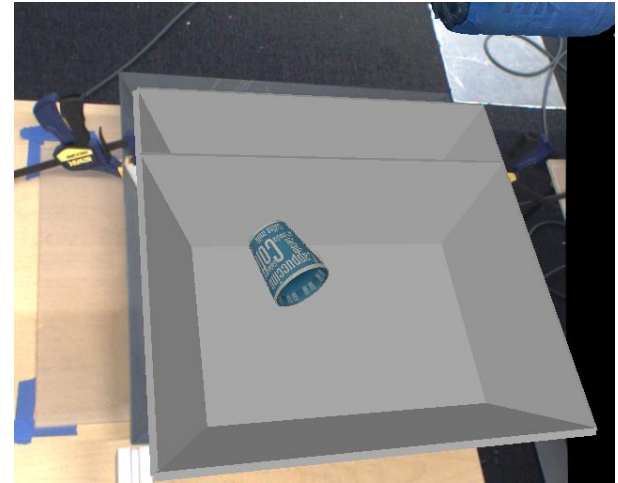
Grasping provides supervision



pre-grasp scene



grasped object



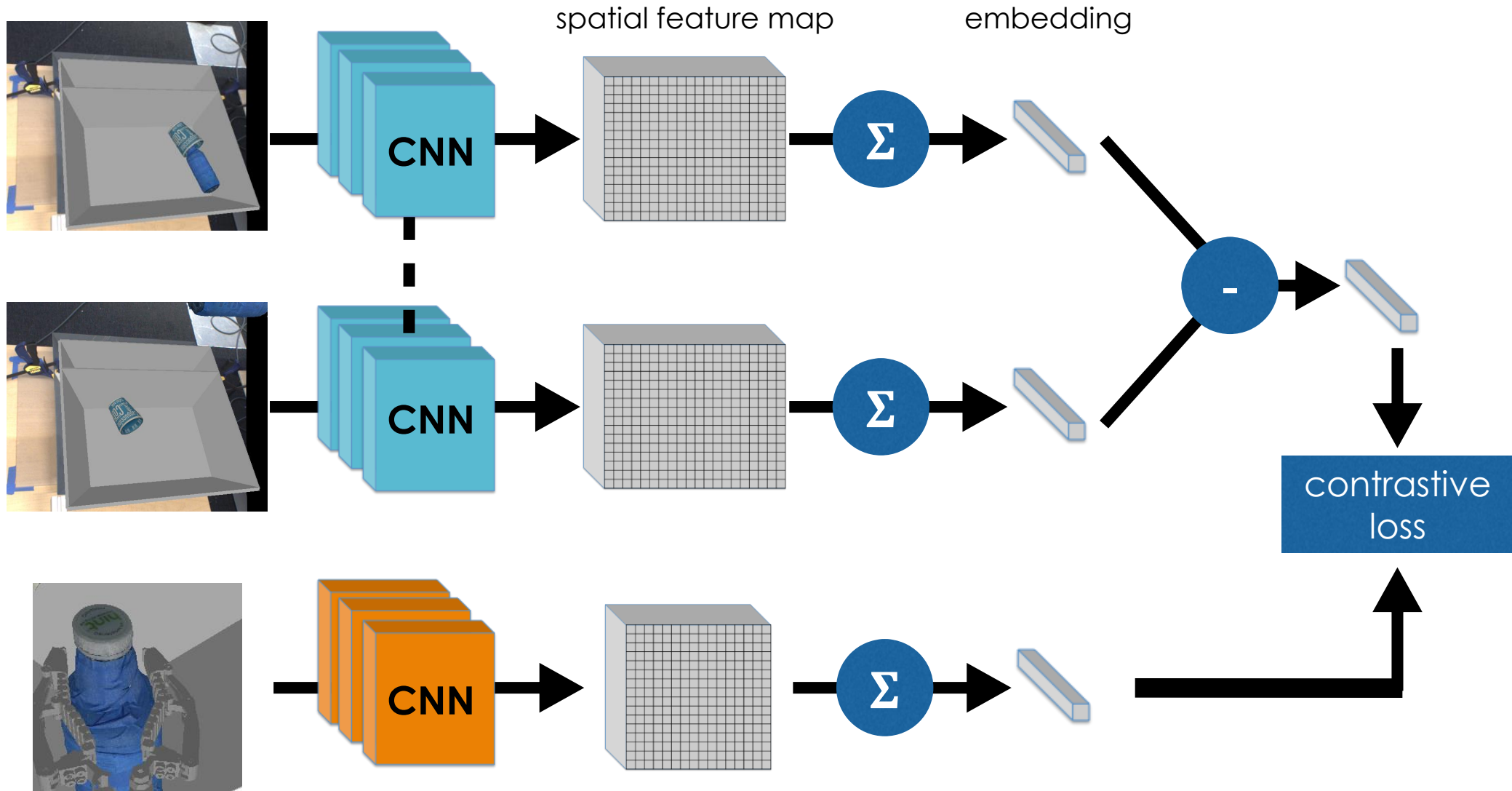
post-grasp scene

Representation learning from grasping

$$\varphi \left(\text{pre-grasp scene} \right) - \varphi \left(\text{post-grasp scene} \right) = \varphi \left(\text{grasped object} \right)$$

The diagram illustrates the relationship between scene representations before and after grasping. It shows three images: a pre-grasp scene with a blue can in a gray tray, a post-grasp scene with the can removed, and a close-up of the grasped object. The equation states that the difference in the learned representation φ between the pre-grasp and post-grasp scenes is equal to the representation of the grasped object.

Training with contrastive loss



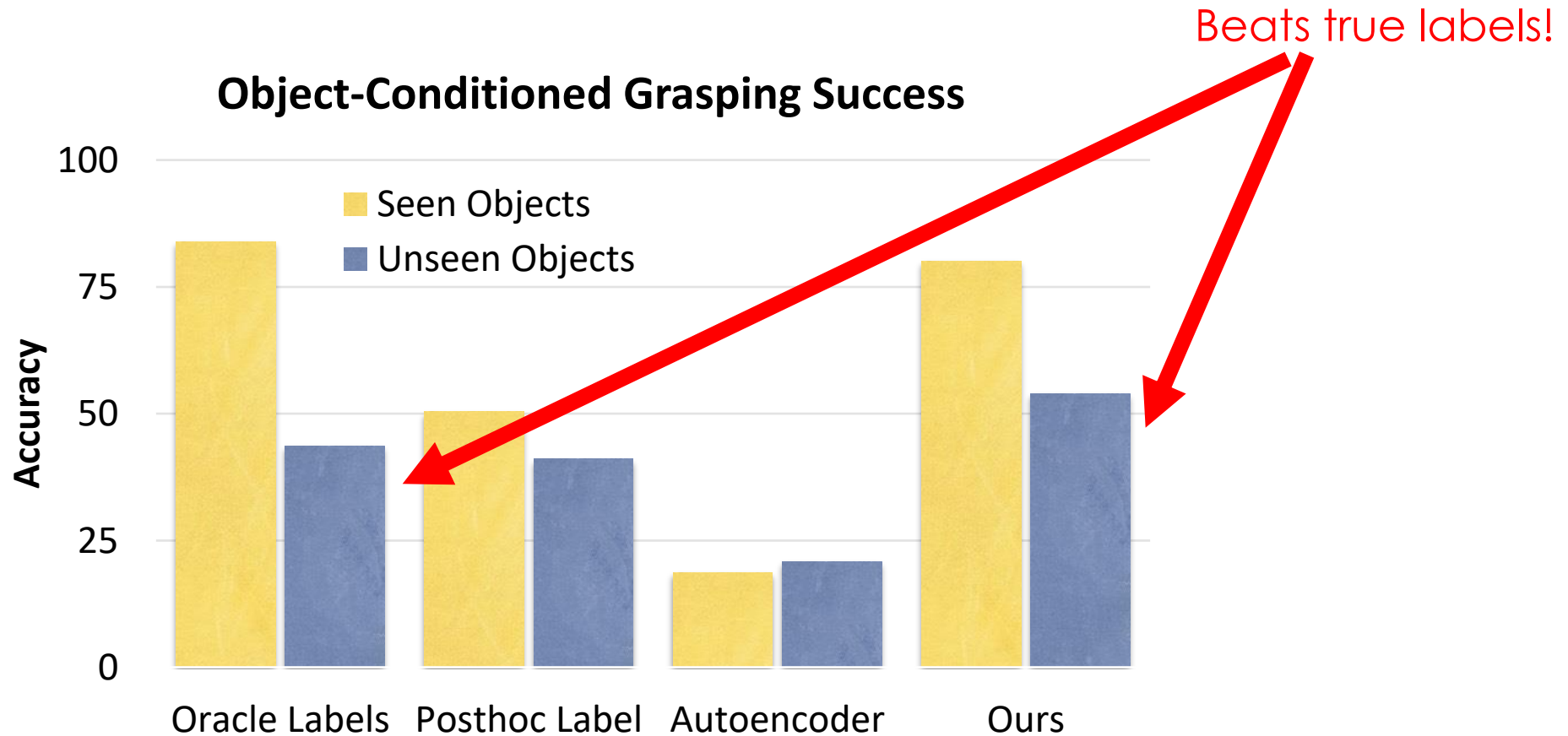
What can we do with the learned representation?

Object-specific grasping



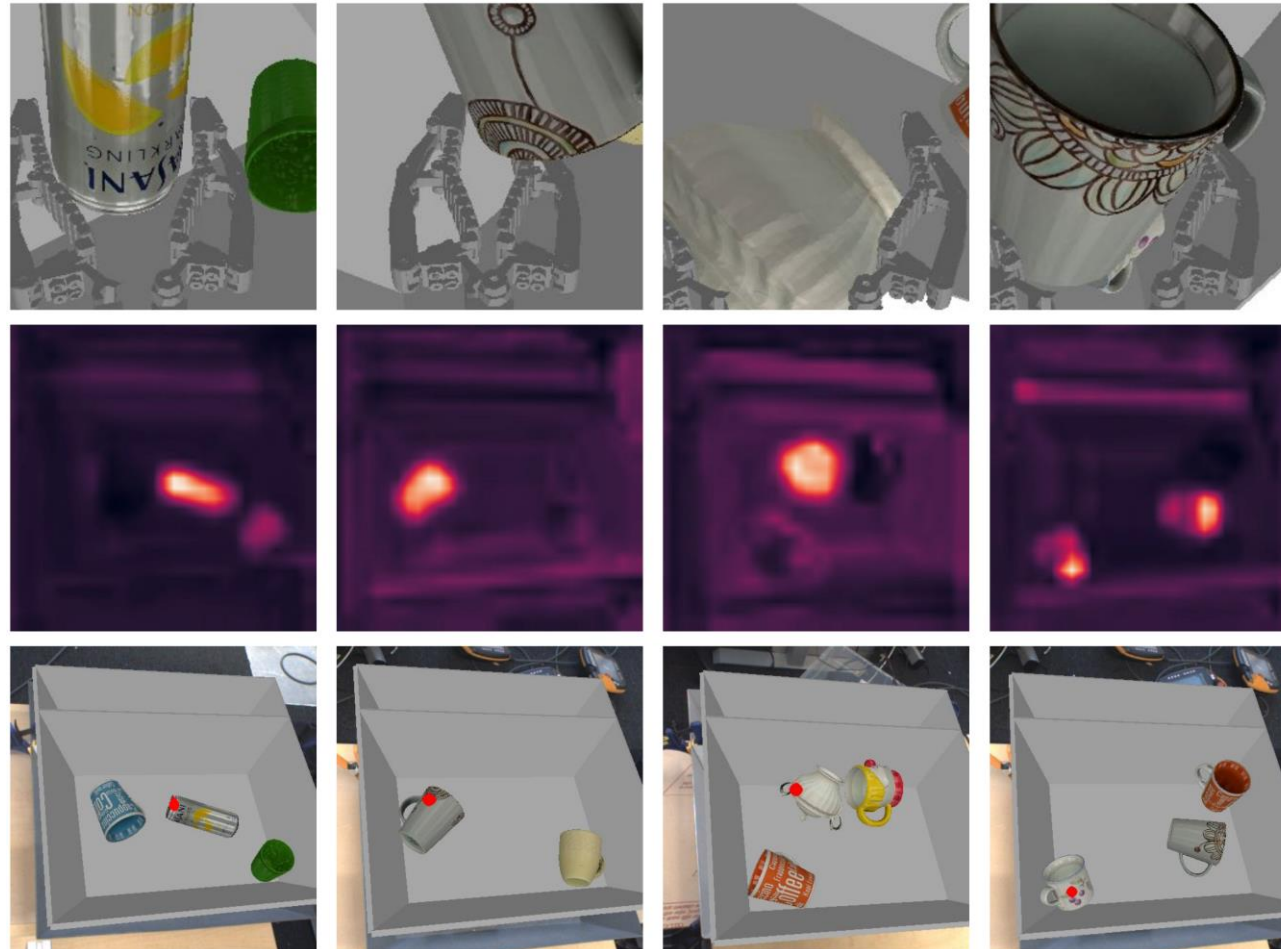
What can we do with the learned representation?

Object-specific grasping

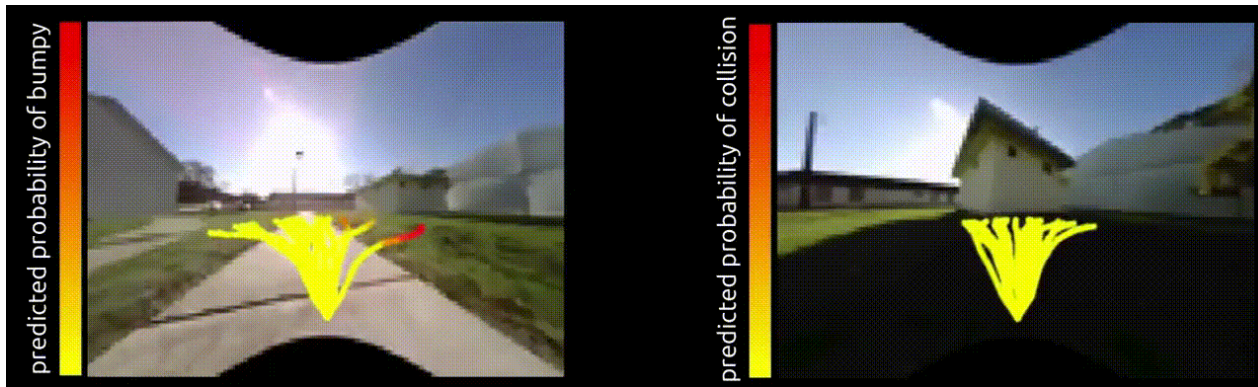
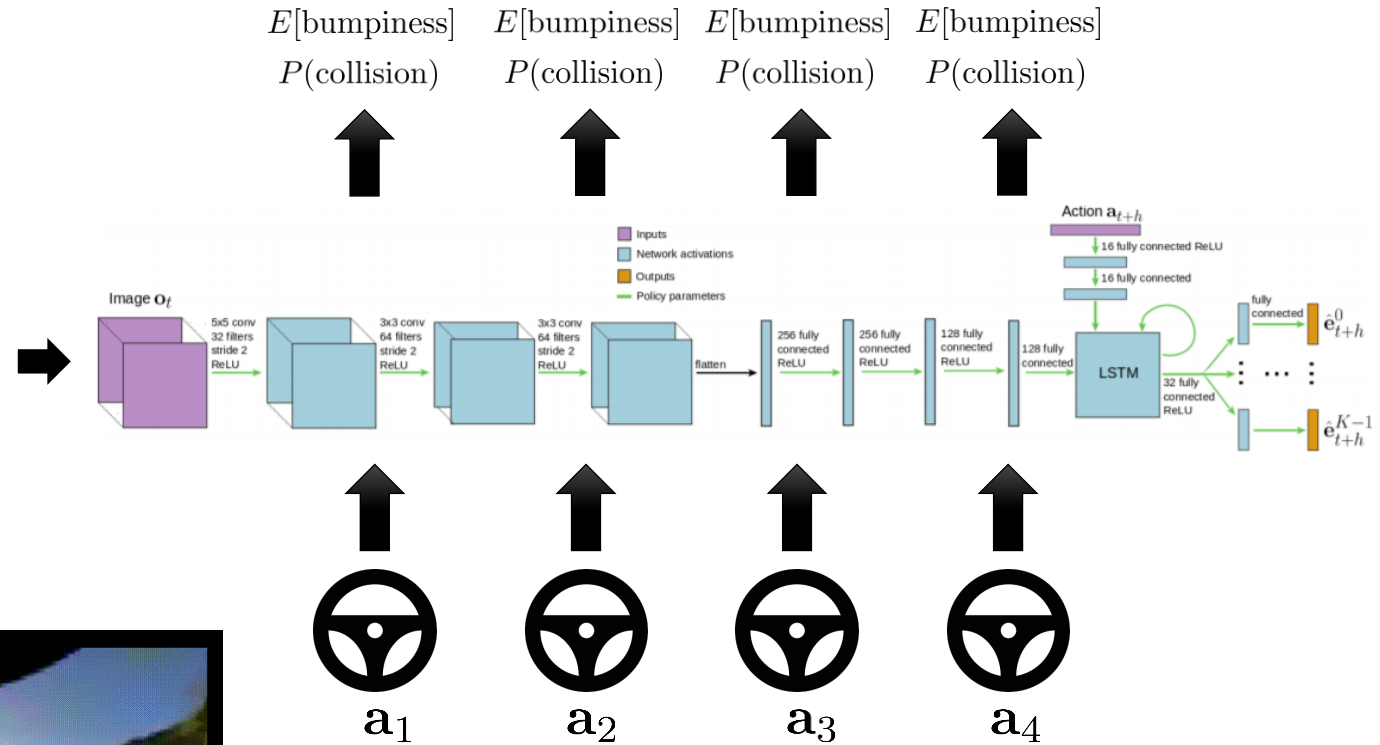


What can we do with the learned representation?

Fully self-supervised localization



Can we learn to understand *open-world* scenes?



Navigational affordances



baseline method

our method

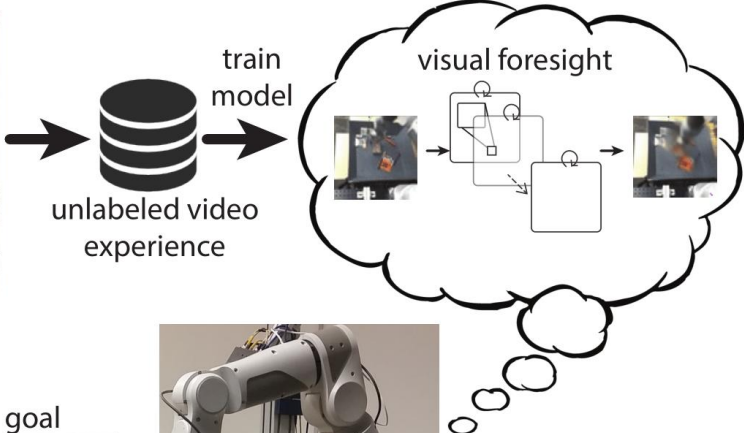
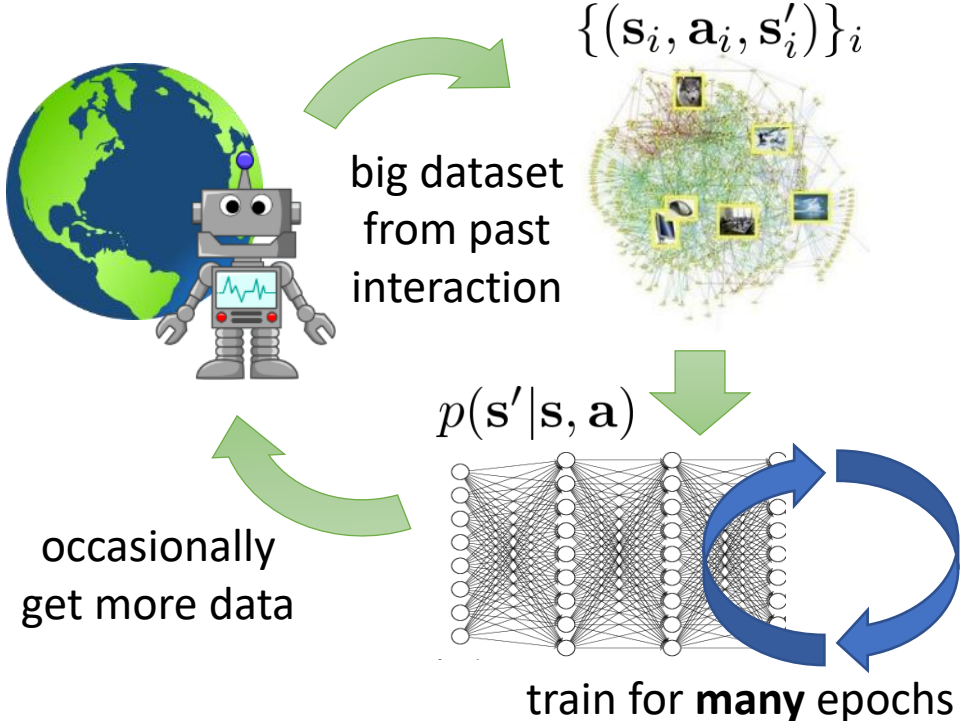
- Learns a kind of “navigational common sense” from experience
- Some obstacles (e.g., grass) are traversable
- Concrete paths are good for avoiding bumpiness

Which end-to-end task should we use?

Model-free algorithms:
predict future *rewards*

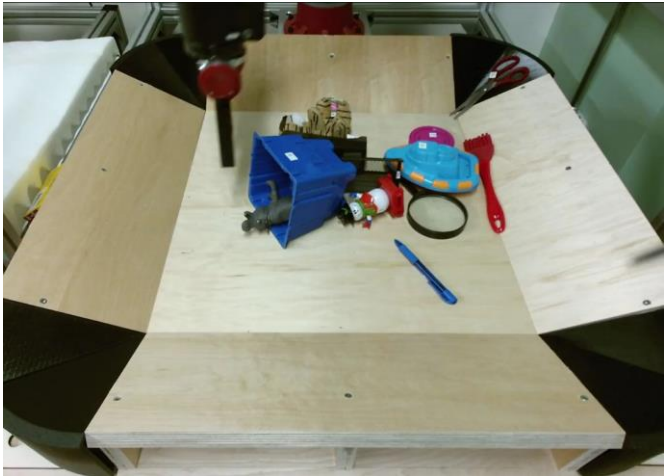
Model-based algorithms:
predict future *observations*

Learning to predict the future

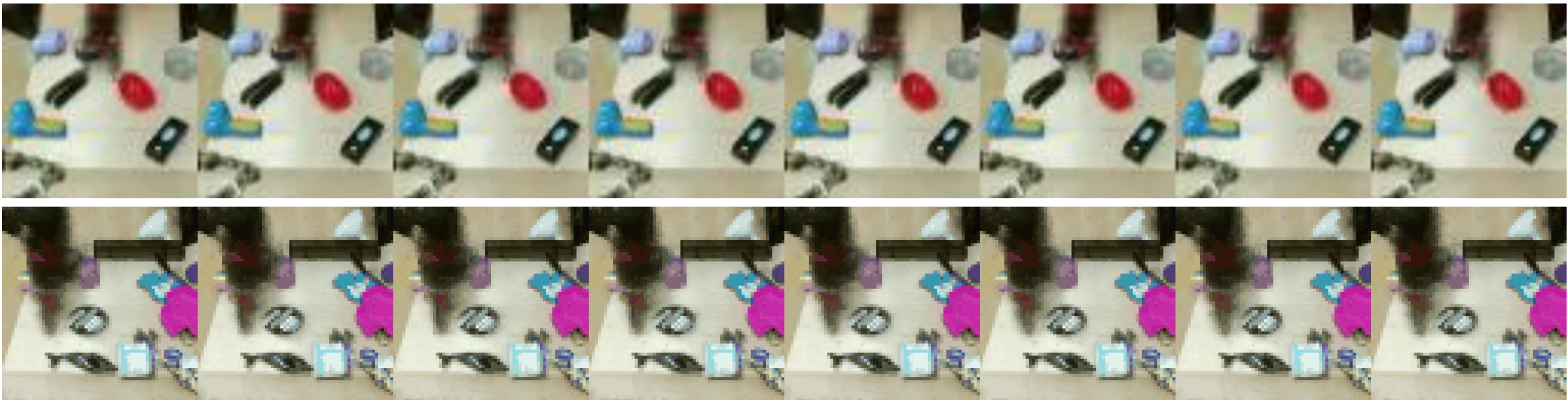
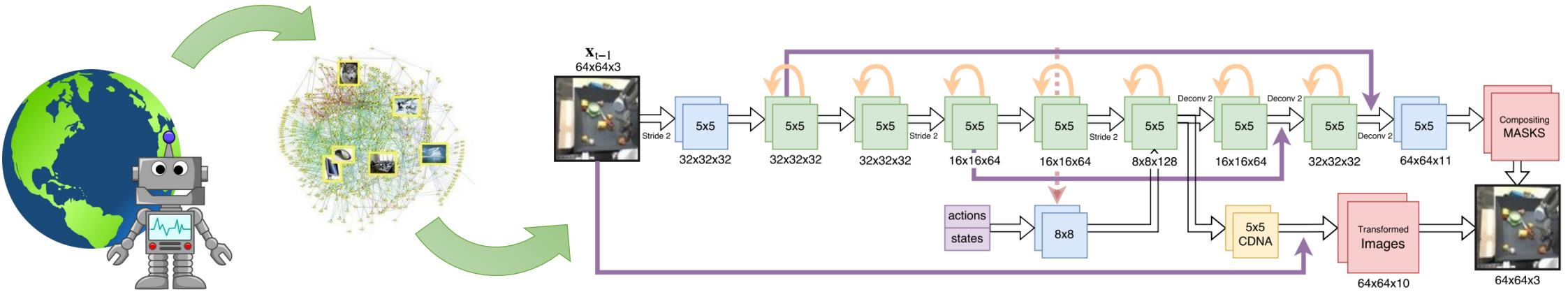


Finn, Levine. **Deep Visual Foresight for Planning Robot Motion.**
Ebert, Finn, Lee, Levine. **Self-Supervised Visual Planning with Temporal Skip Connections.**
Lee, Zhang, Ebert, Abbeel, Finn, Levine. **Stochastic Adversarial Video Prediction.**

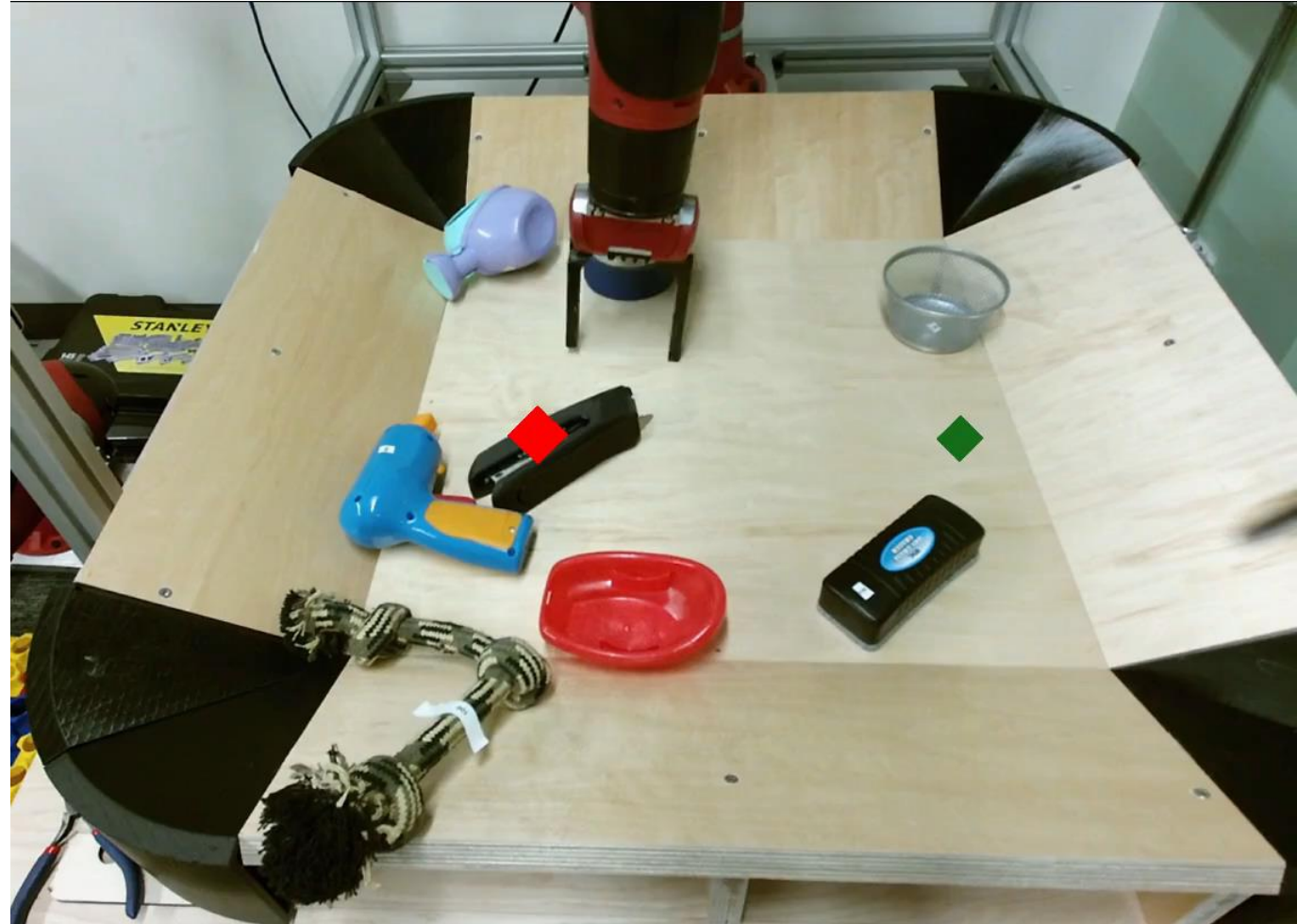
Collect data by playing with objects



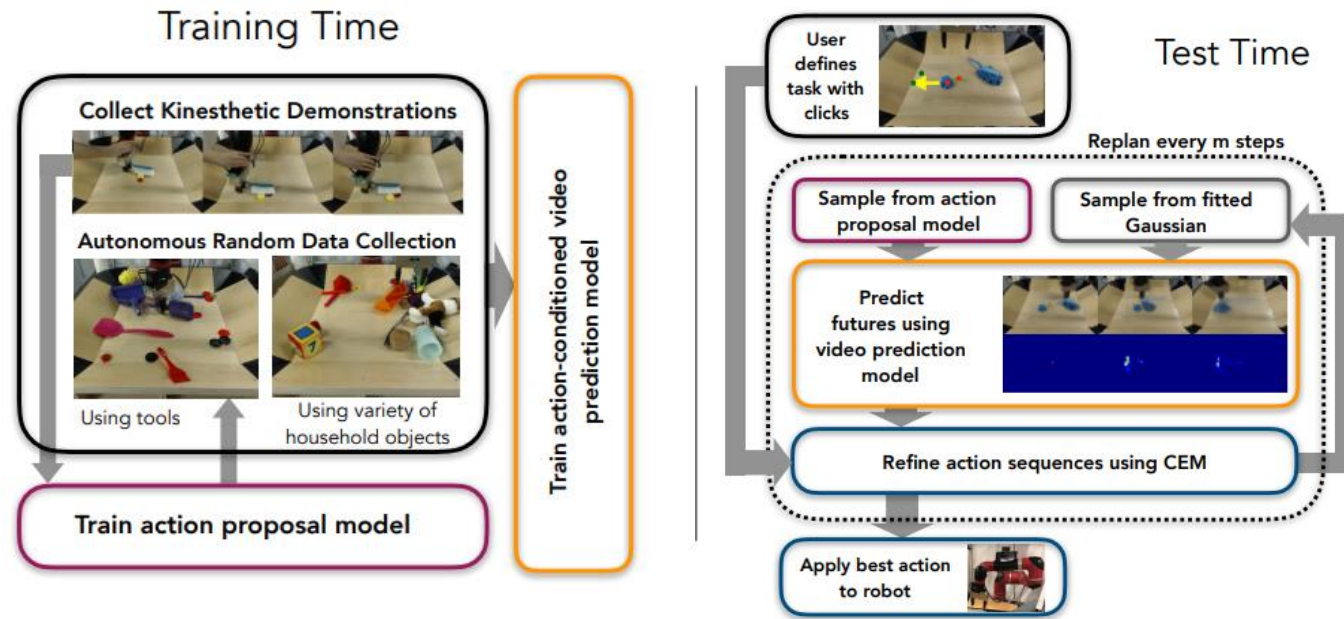
Learn to predict the future



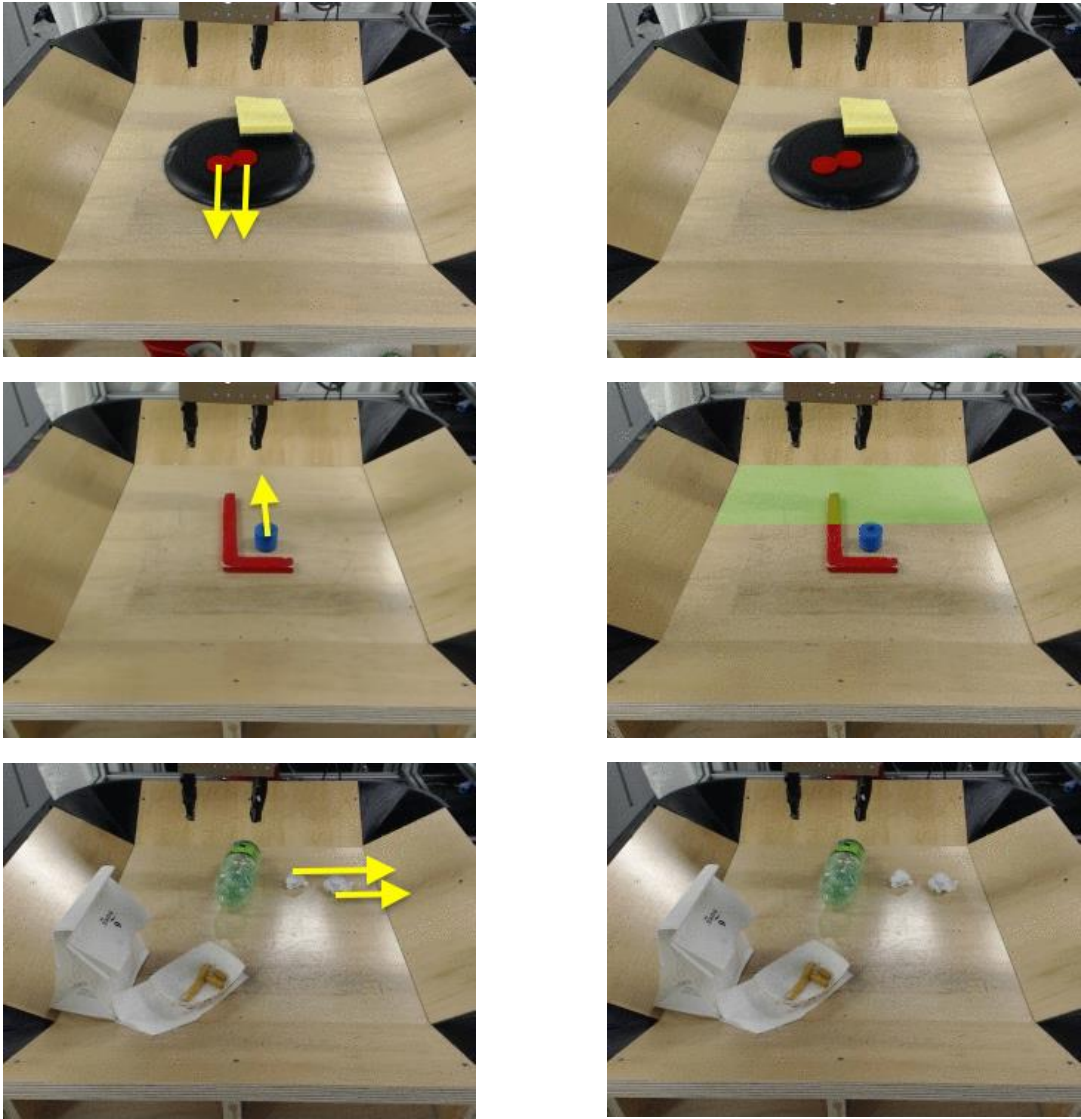
Example execution



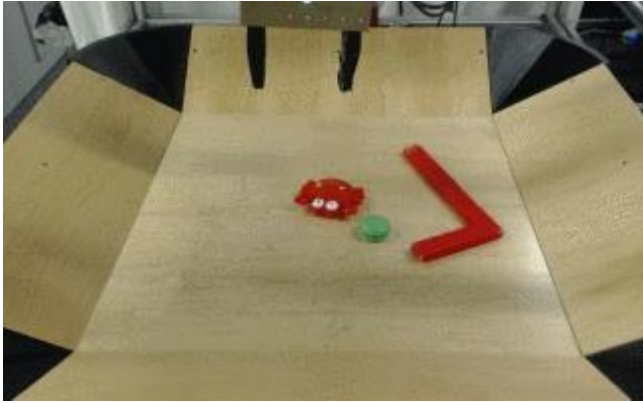
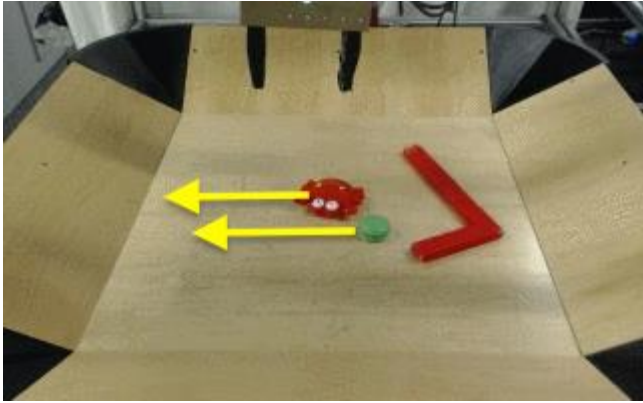
Model-based RL with improvised tool use



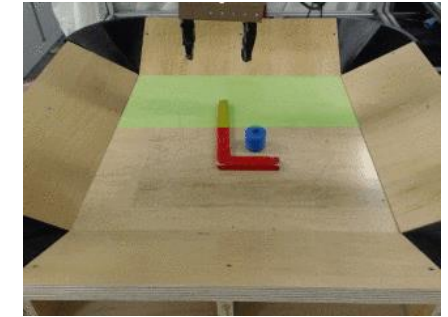
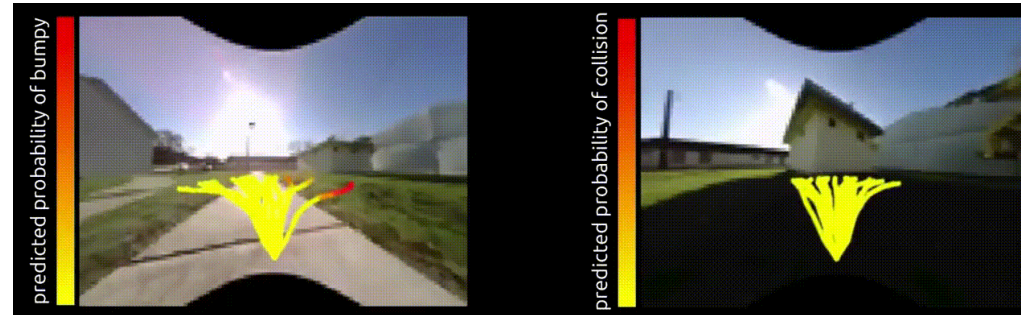
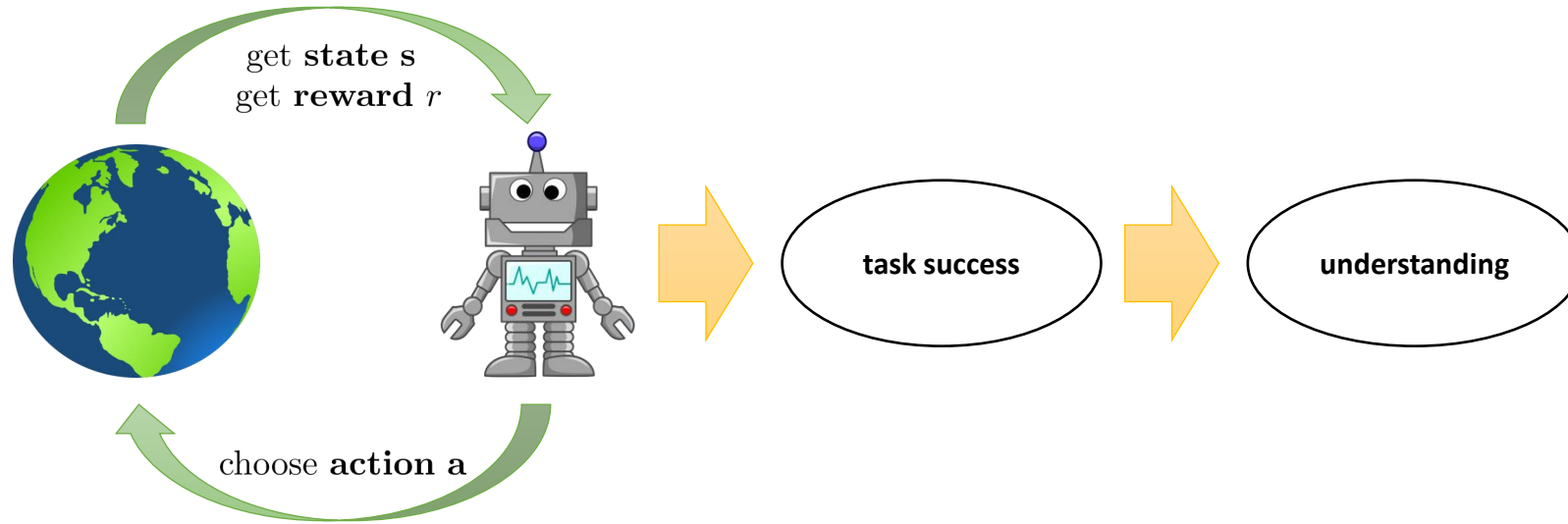
Model-based RL with improvised tool use



Model-based RL with improvised tool use



An embodied learning recipe for scene understanding?





RAIL
Robotic AI & Learning Lab

website: <http://rail.eecs.berkeley.edu>
source code: <http://rail.eecs.berkeley.edu/code.html>