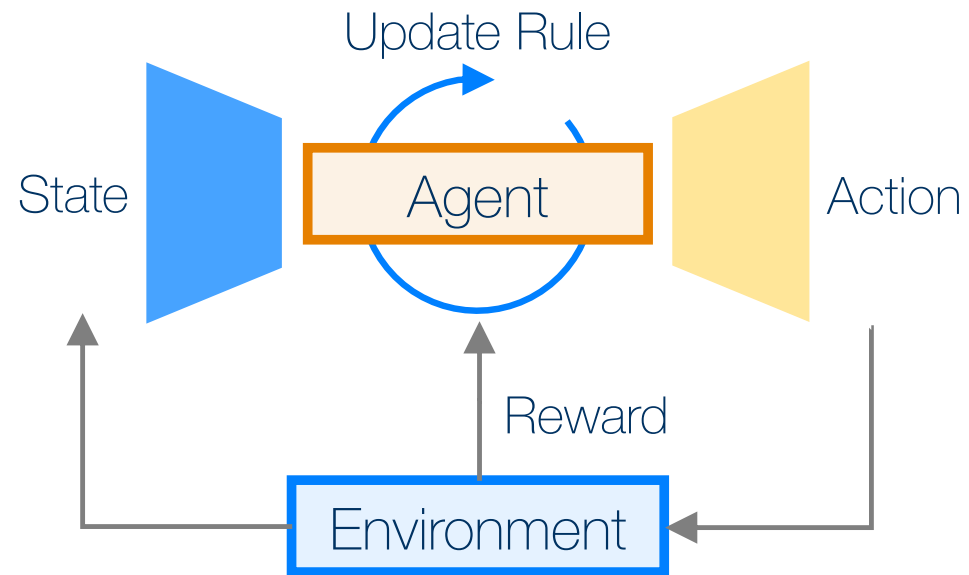


Towards Generalizable Autonomy

Structure in Reinforcement Learning for Robotics



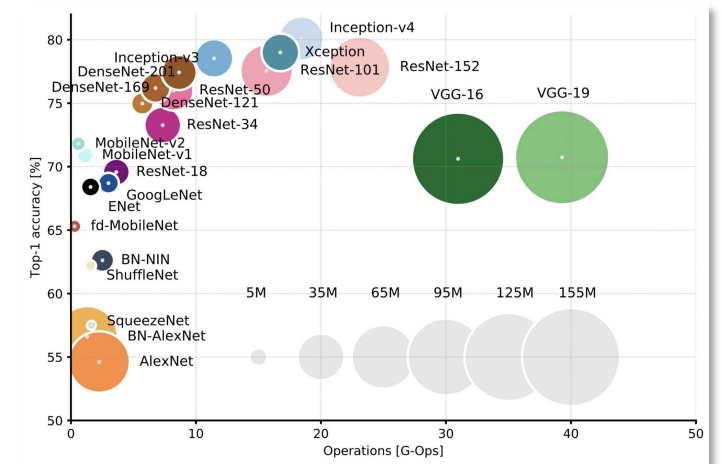
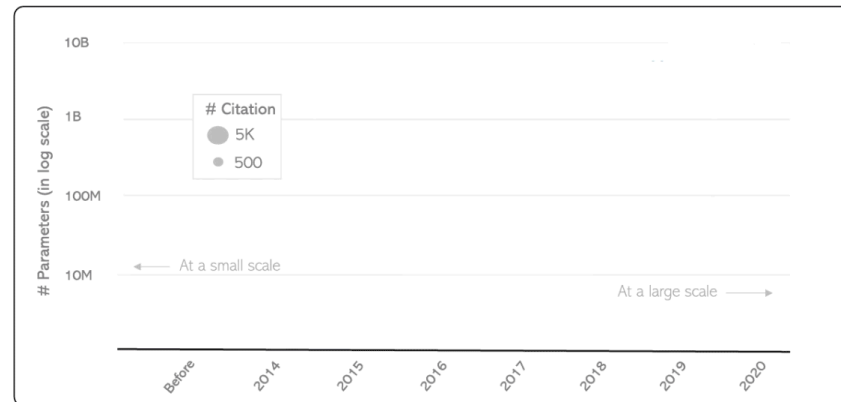
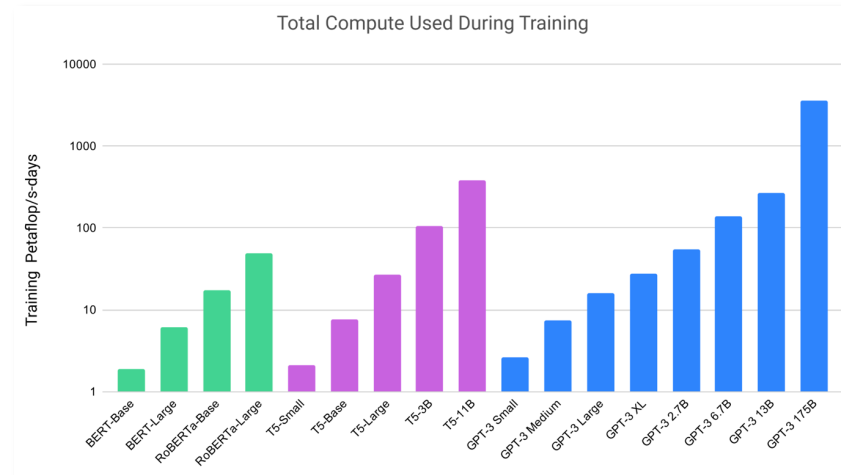
Animesh Garg

Generalizable Autonomy: Computer Vision & Language

Structured Models + Data + Compute → Performance



Open Images Dataset



Model	EM	F1
Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
IE-Net (ensemble) RICOH_SRCB_DML	90.939	93.214

Generalizable Autonomy: Computer Vision & Language

Ingredients of Modern Machine Learning & Applications



Large Structured Models

- Over-parameterized
- Structured Biases



IID Data & Datasets

- Concise problem Definition
- IID Data, easier to label



Distributed Deployment

- Large Scale Compute
- Distributed Deployment

Visual
Perception

Natural
Language

Passive Offline Decisions

Intelligent
Robotics

Embodied

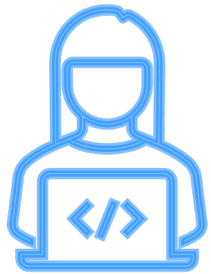
Generalizable Autonomy: Duality of Discovery & Bias



Domain
Expertise

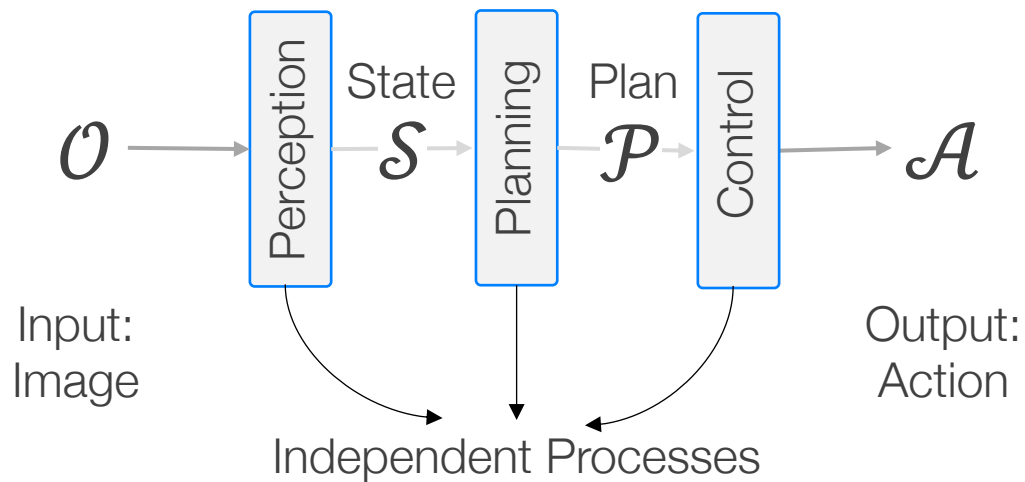
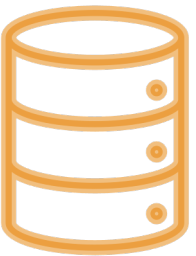


One problem,
One solution!



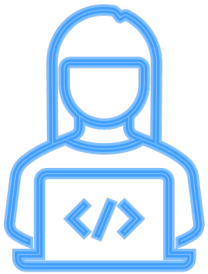
Structured
Environments

Data
Driven



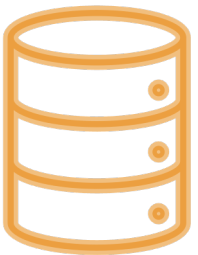
Generalizable Autonomy: Duality of Discovery & Bias

Domain
Expertise



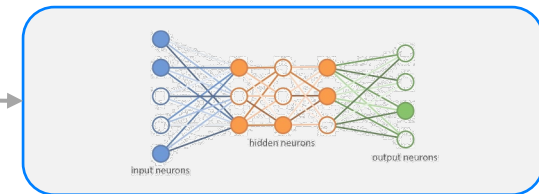
Just add
data...

Data
Driven



**The Unreasonable
Effectiveness of Data**
Alon Halevy, Peter Norvig, and Fernando Pereira, Google

\mathcal{O}



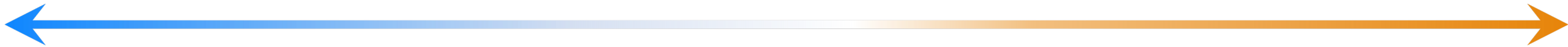
\mathcal{A}

Input:
Image

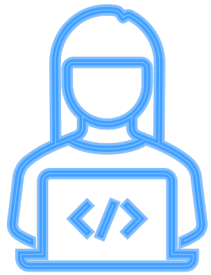
End-to-End Policy
 π

Output:
Action

Generalizable Autonomy: Duality of Discovery & Bias



Domain
Expertise



One problem,
One solution!

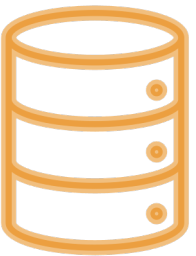
- X Need for experts
- X Limited applicability
- X Perf vs Flexibility



Just add
data...

- X Computational sustainability
- X Data accessibility
- X Out-of-distribution errors

Data
Driven



Neither achieves
Generality at Scale

Generalizable Autonomy: Duality of Discovery & Bias

Domain
Expertise

Data
Driven

...make the **inductive** leap necessary to classify instances beyond observed...

...other sources of information, or **biases** for choosing one generalization over the other...

No Generalization
without Structure!



Philosophy
David Hume
1739



Psychology
Nelson Goodman
1955



Machine Learning
Tom Mitchell
1980



Deep Learning
Bengio, Hinton,
LeCun 2020s

Generalizable Autonomy: Duality of Discovery & Bias



Domain
Expertise

Data
Driven

Generalizable Autonomy

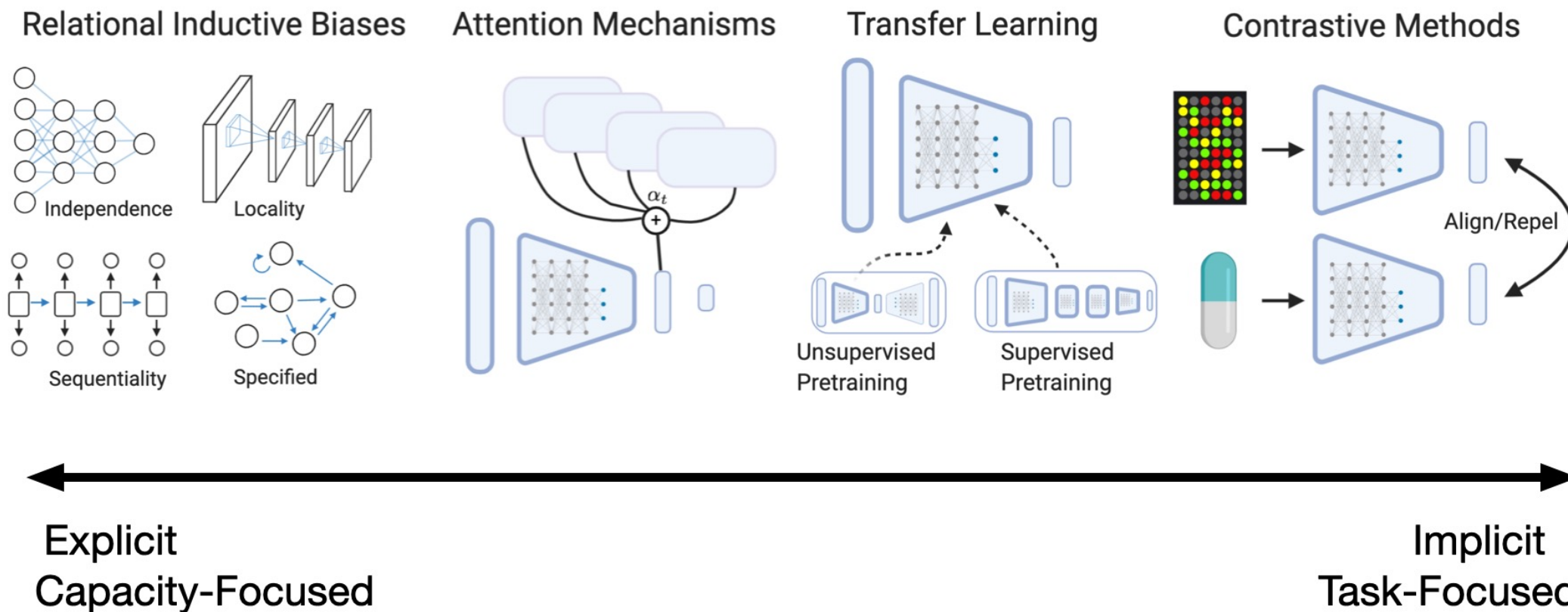
Structure + Data

- Domain knowledge,
- Inductive bias,
- Symmetries,
- Priors
- ...

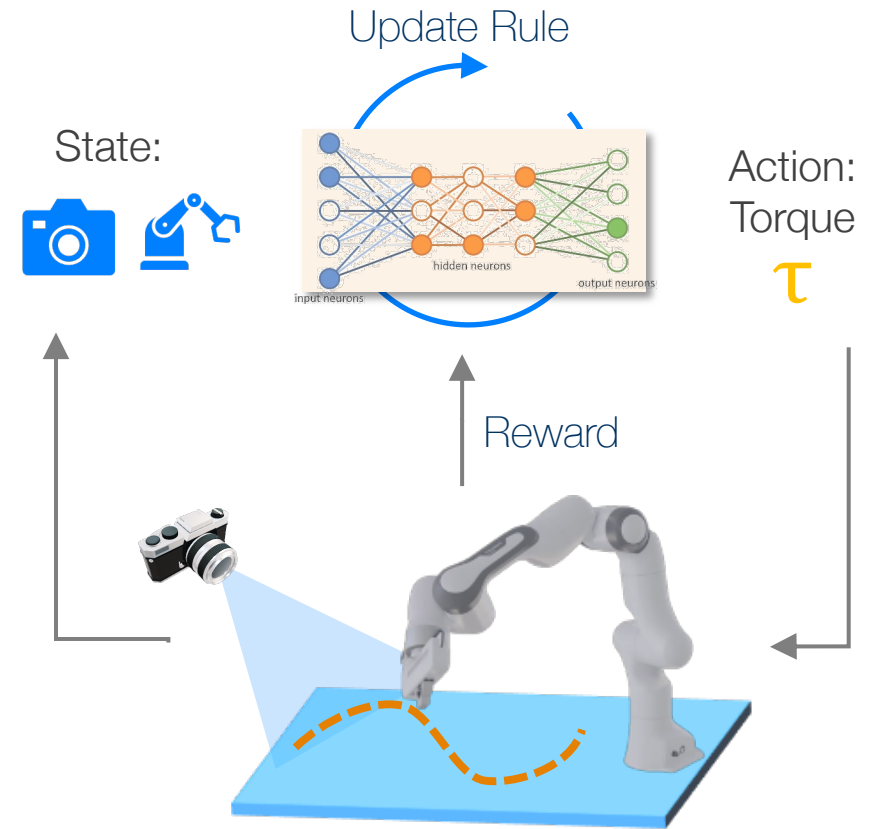
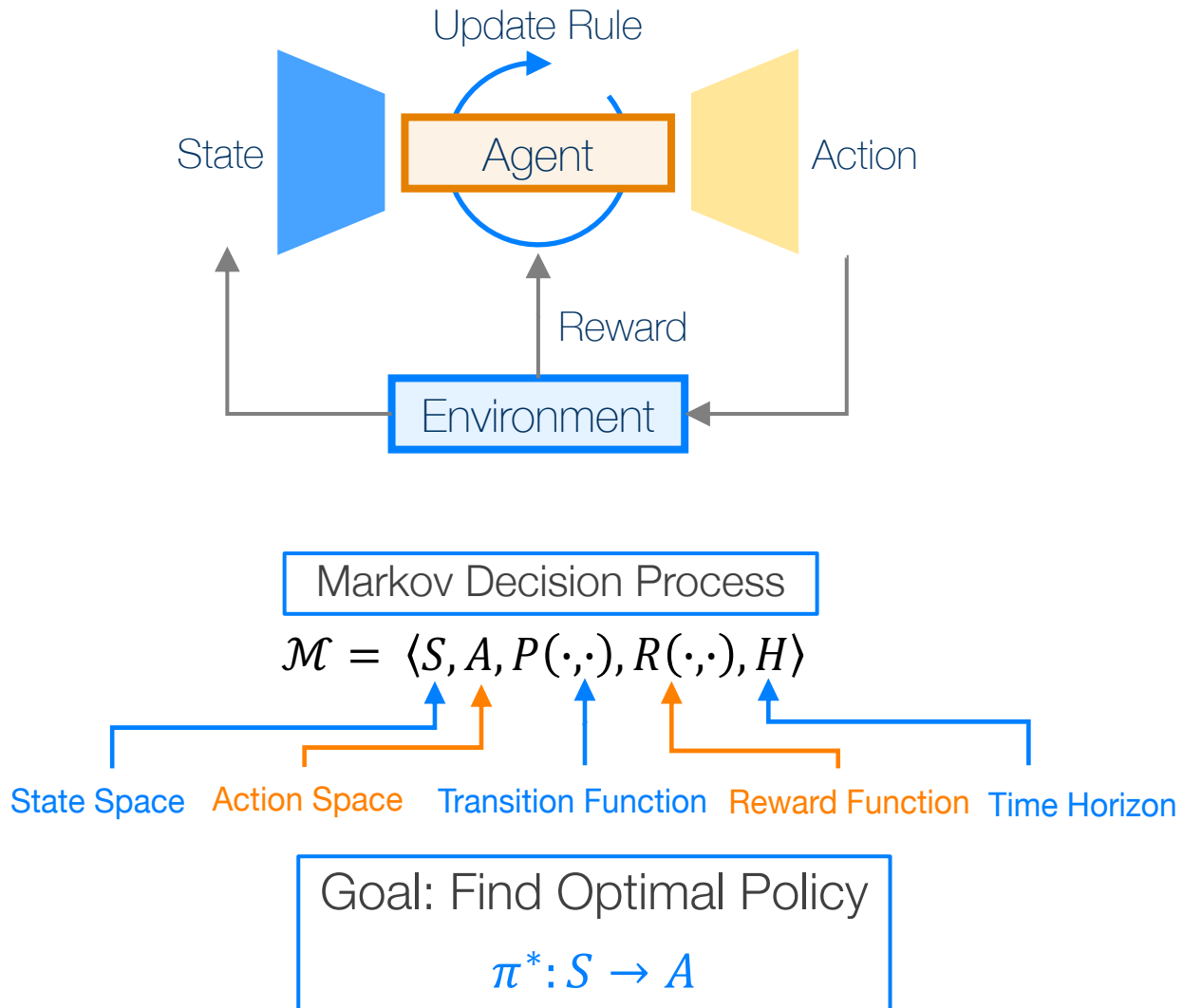
- Online & Offline,
- Simulation & Real,
- Labelled & self-supervised
- Human in the loop
- ...

Structured Representations: Vision & Language

Insight: Structure makes learning possible



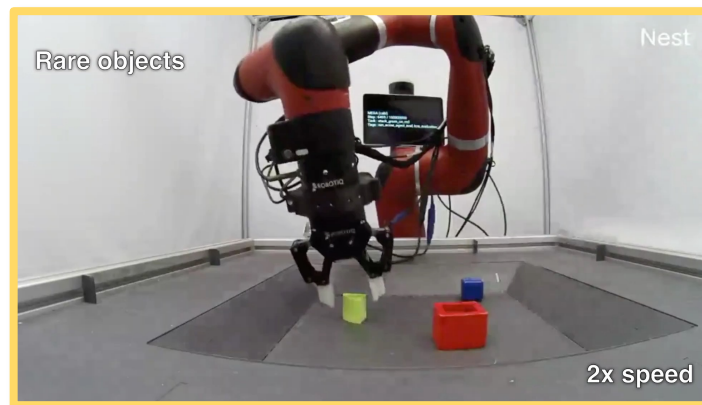
Structure in Reinforcement Learning



Structure in Skill Learning: or the lack of it

Slow and Narrow:

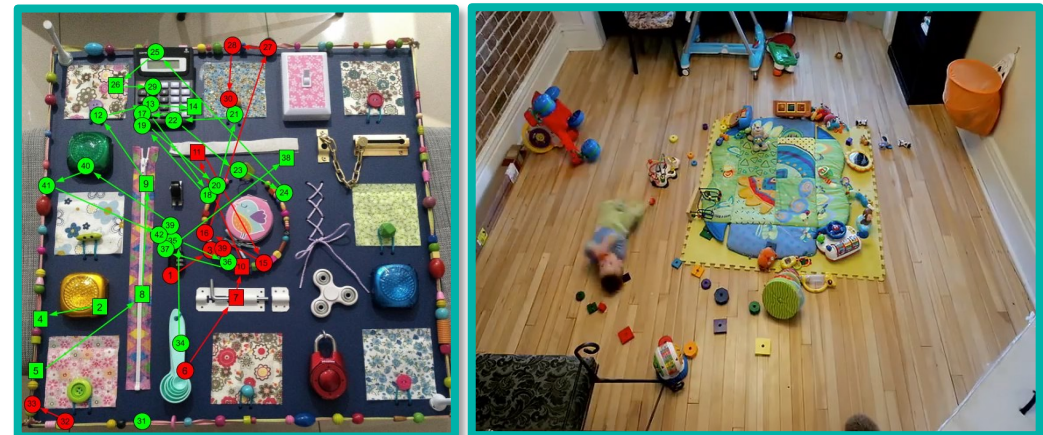
- Specific tasks (Grasp/Stack)
- Often Supervised and rigid!



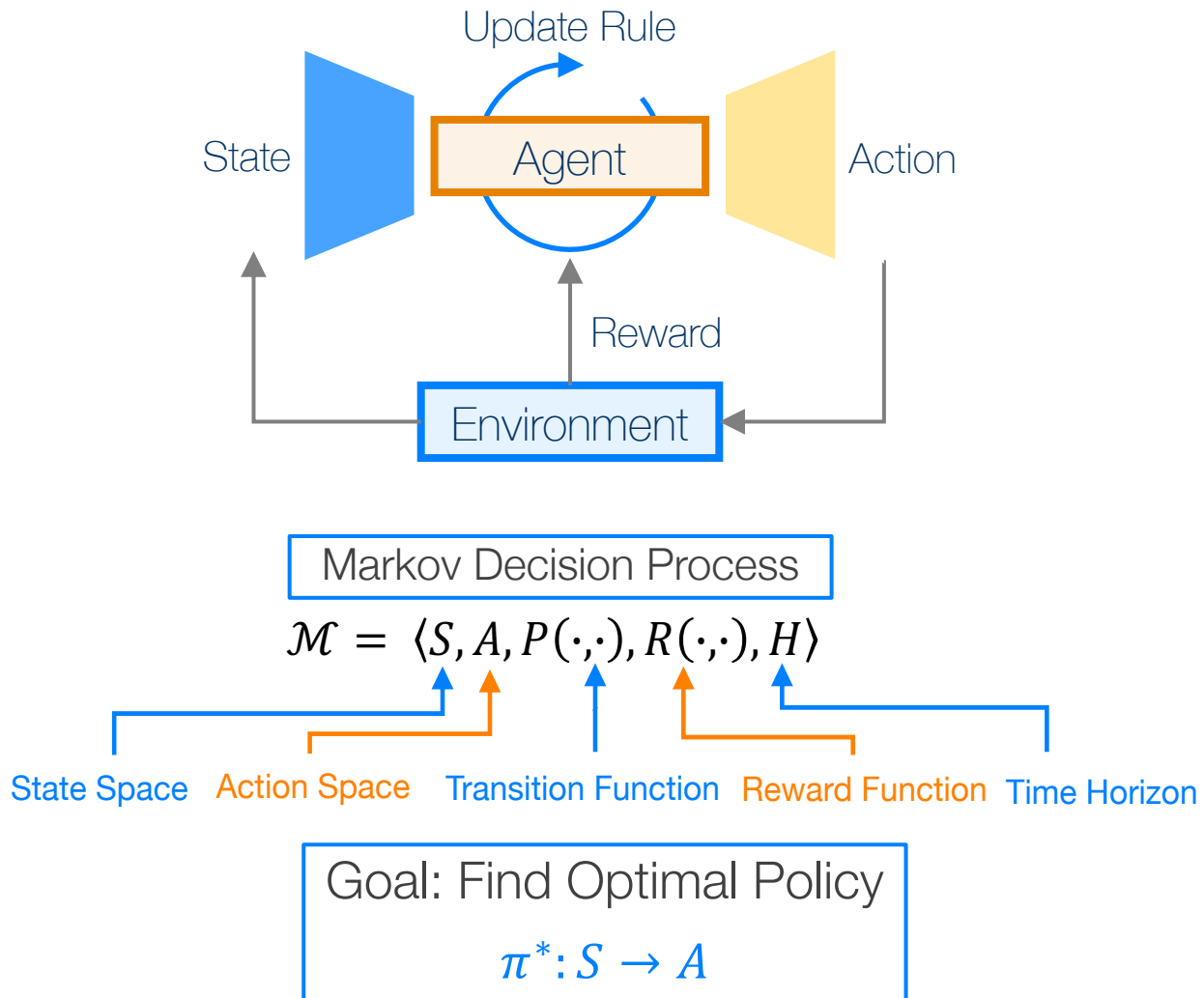
Learning Fast: Elephants Learning to use trunk



Learning Broad: Human infant learns to interact

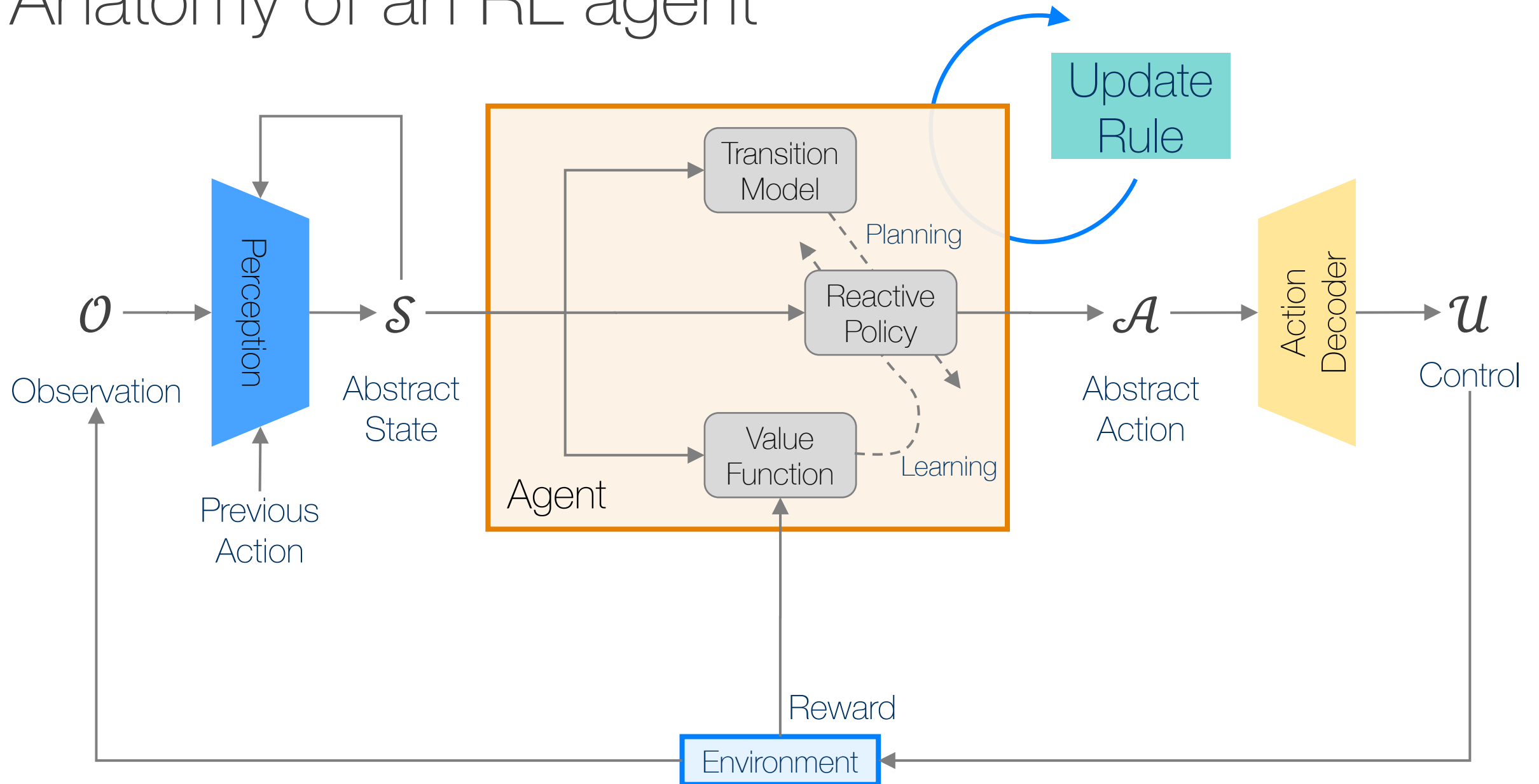


Structure for Reinforcement Learning

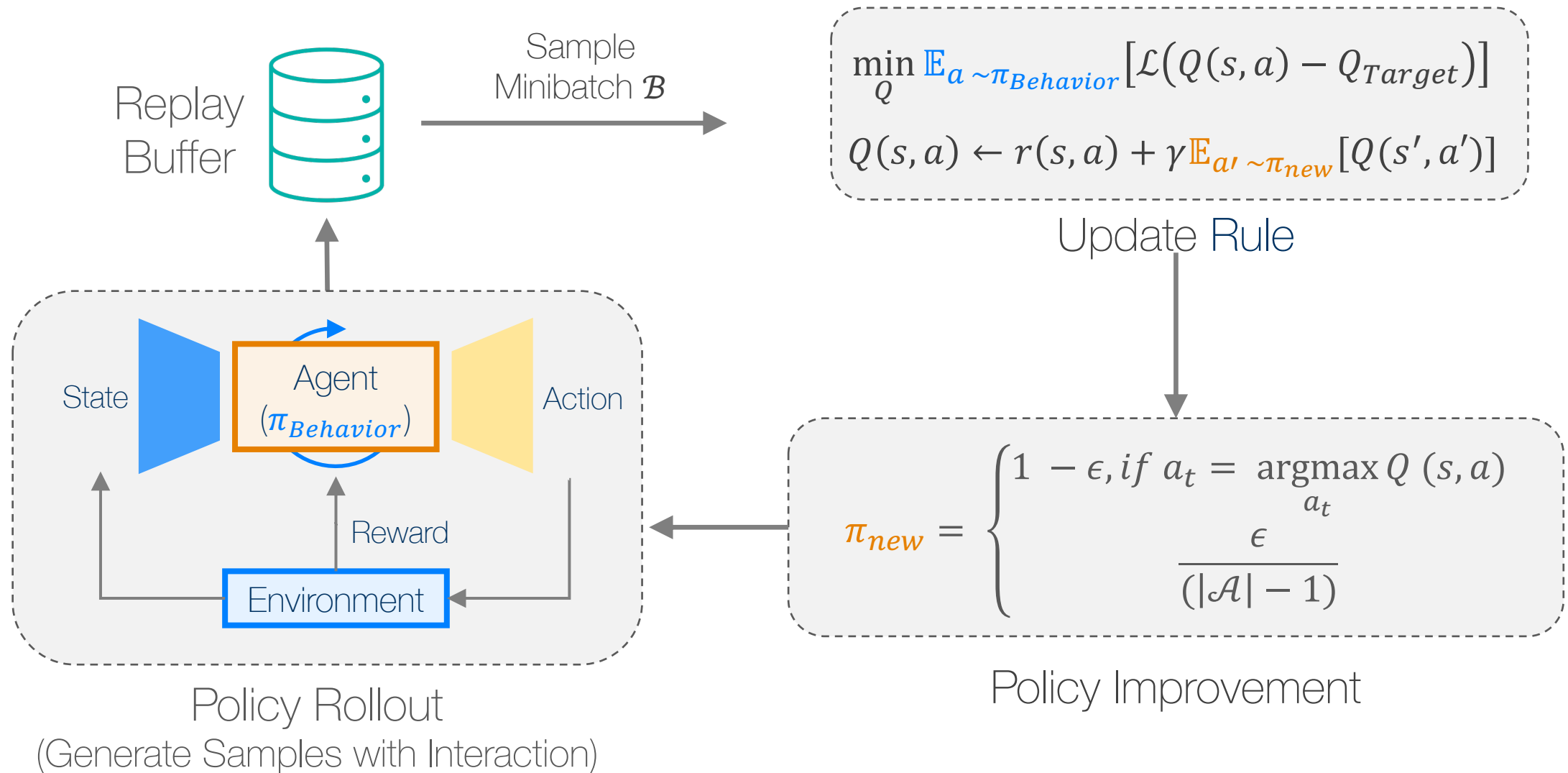


Which structured biases enable generalizable autonomy in decision-making?

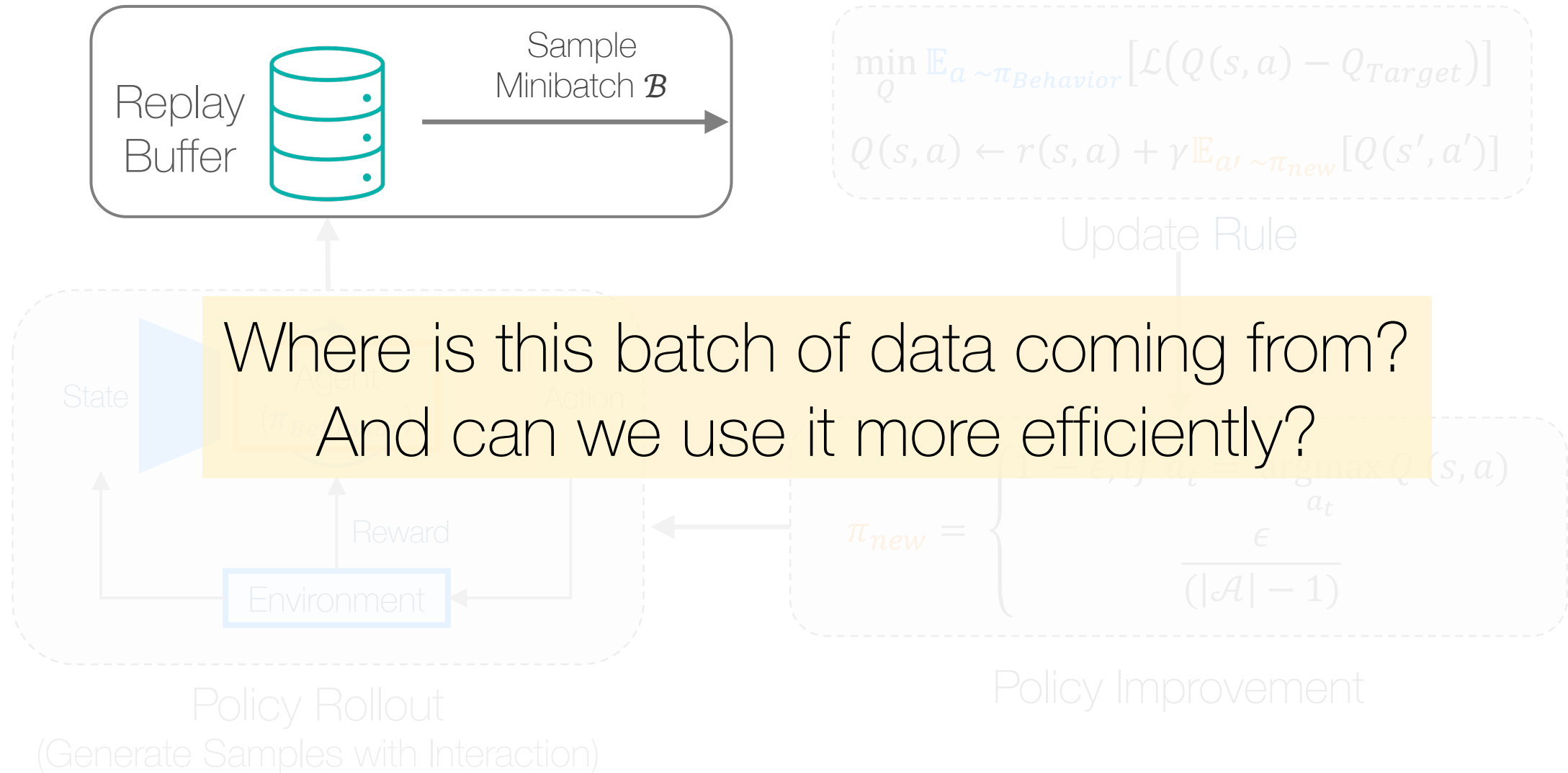
Anatomy of an RL agent



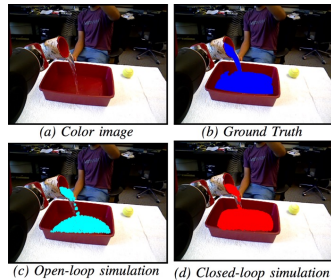
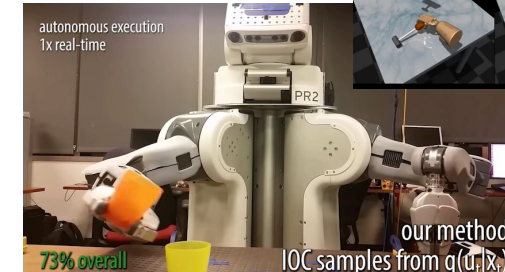
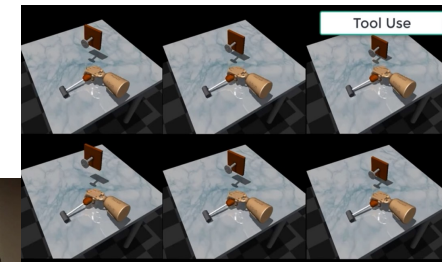
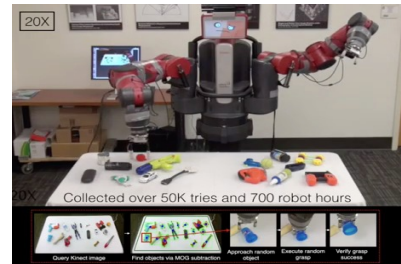
Structure for RL: Off-policy RL



Structure for RL: Off-policy RL



Data in Robotics



Manipulation

- | | |
|-------------------------|-----------------------|
| Mason & Salisbury 1985 | Li, Allen et al. 2015 |
| Srinivasa et al 2010 | Yahya et al, 2016 |
| Berenson 2013 | Schenck et al. 2017 |
| Odhner1 et al 2014 | Mar et al. 2017 |
| Chavan-Dafle et al 2014 | Laskey et al 2017 |
| Yamaguchi, et. al, 2015 | Quispe et al 2018 |
| ... | ... |

Grasping

- | | |
|-------------------------|---------------------|
| Mishra et al 1987 | Pinto & Gupta, 2016 |
| Ferrari & Canny, 1992 | Levine et al 2016 |
| Ciocarlie & Allen, 2009 | Mahler et al 2017 |
| Dogar & Srinivasa, 2011 | Jang et al 2017 |
| Rodriguez et al. 2012 | Viereck et al 2017 |
| Bohg et al 2014 | ... |

Imitation

- | | |
|------------------------|---------------------------|
| Abbeel et al, 2004 | Krishnan et al 2017 |
| Ratliff et al 2006, | Finn et al. 2017 |
| Ziebart et al, 2009 | Vecerik et al. 2017 |
| Argall et al, 2009, | Rajeswaran et al 2018 |
| Boularias et al., 2011 | Zhu et al 2018 |
| Montfort et al 2015, | Ravichandar et al 2020... |
| Wulfmeier et al 2015, | |

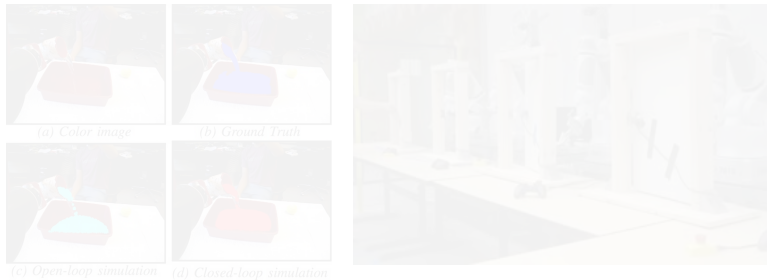
Data in Robotics

Robot Data is Expensive!

X Short-Horizon skills
X Skill Specific learning

X Platform dependent data
X Scaling to other skills ?

X Small datasets (minutes)
X Low diversity ?



Manipulation

Mason & Salisbury 1985
Srinivasa et al 2010
Berenson 2013
Odhner et al 2014
Chavan-Dafle et al 2014
Yamaguchi, et. al, 2015
...

Li, Allen et al. 2015
Yahya et al, 2016
Schenck et al. 2017
Mar et al. 2017
Laskey et al 2017
Quispe et al 2018
...



Grasping

Mishra et al 1987
Ferrari & Canny, 1992
Ciocarlie & Allen, 2009
Dogar & Srinivasa, 2011
Rodriguez et al. 2012
Bohg et al 2014

Pinto & Gupta, 2016
Levine et al 2016
Mahler et al 2017
Jang et al 2017
Viereck et al 2017
...



Imitation

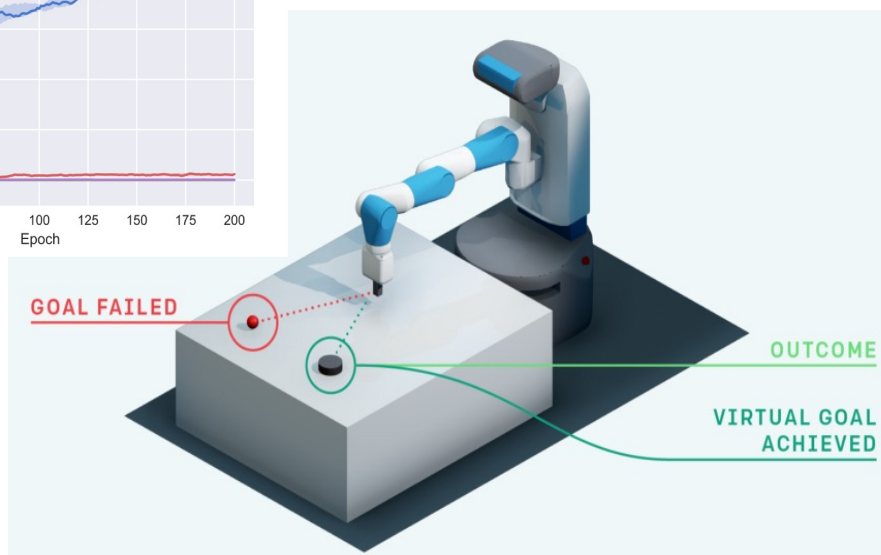
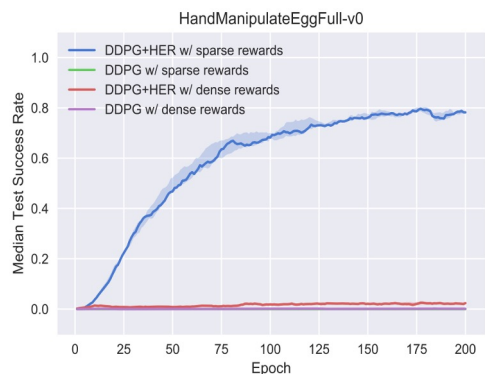
Abbeel et al, 2004
Ratliff et al 2006,
Ziebart et al, 2009
Argall et al, 2009,
Boularias et al., 2011
Montfort et al 2015,
Wulfmeier et al 2015,

Krishnan et al 2017
Finn et al. 2017
Vecerik et al. 2017
Rajeswaran et al 2018
Zhu et al 2018
Ravichandar et al 2020...

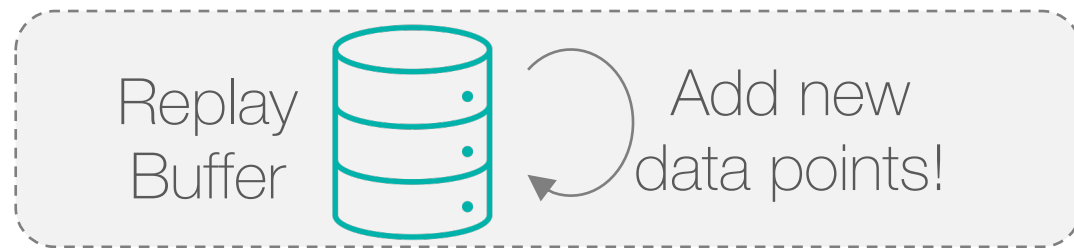
Data Augmentation in RL

How to do this Algorithmically

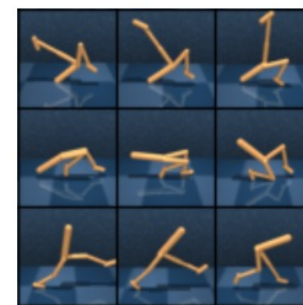
- Substantial performance boosts!



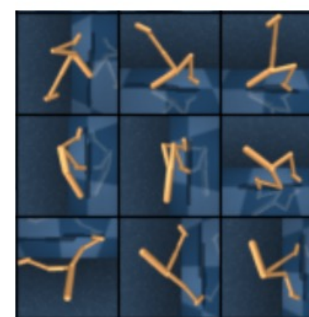
Goal relabeling (e.g, HER)]



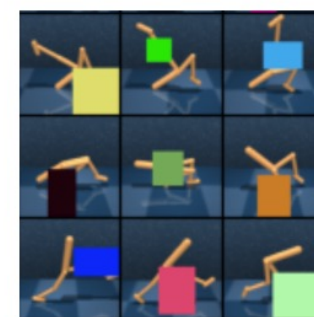
Input



Rotate



Cutout-color



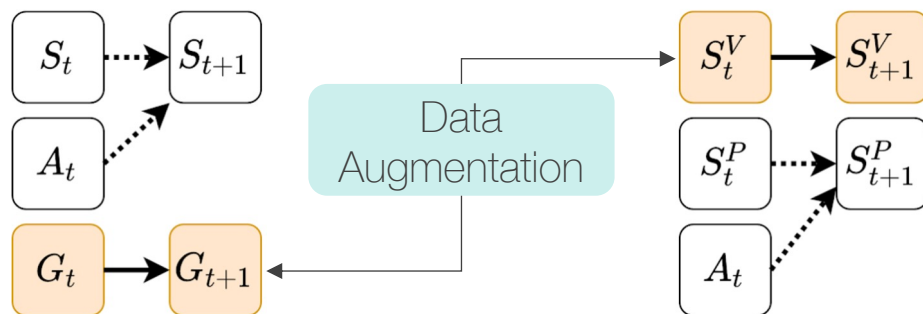
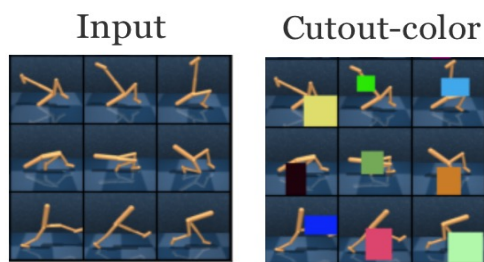
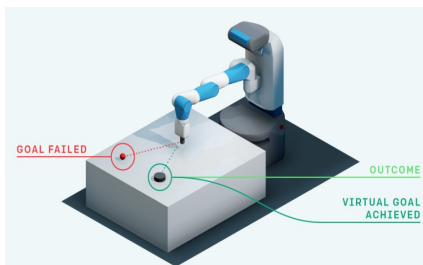
Random conv



Visual Input relabeling (e.g, RAD)]


Data Augmentation in RL

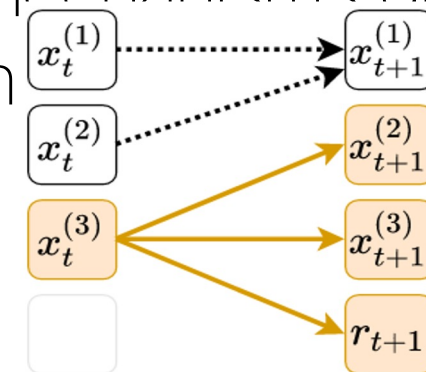
Unified View



Goal is independent of State/Action Dynamics

Visual characteristics (e.g., crop) are independent of physical dynamics

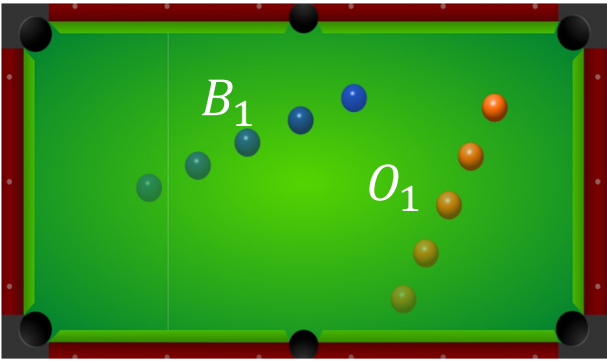
 **Insight** (counterfactual) data!
 Exploit the **Conditional Independence** of the causal mechanisms guiding transition



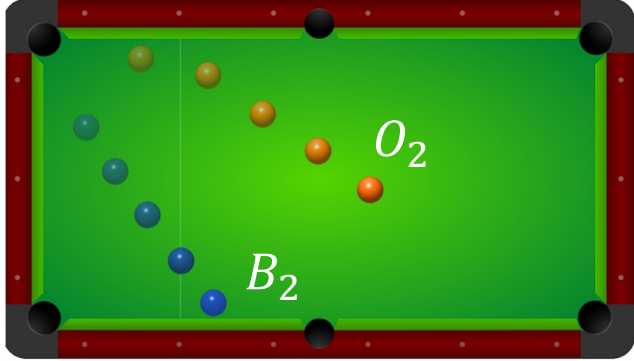
Given two independent mechanisms, Relabel one (conditional independence!)

Data Augmentation in RL

Do more with the same data

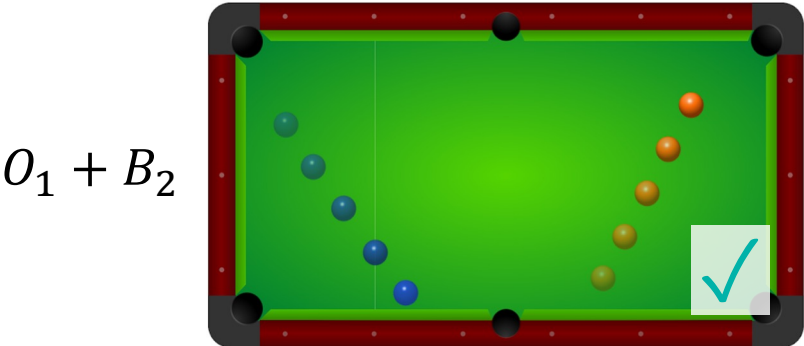


Scenario 1

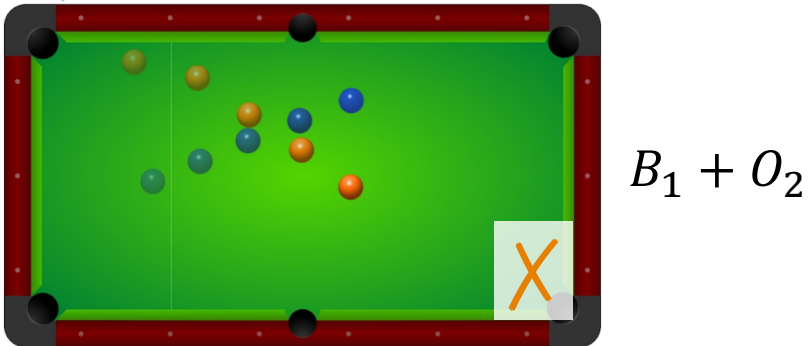


Scenario 2

Which of the following is possible (only based on observed data)



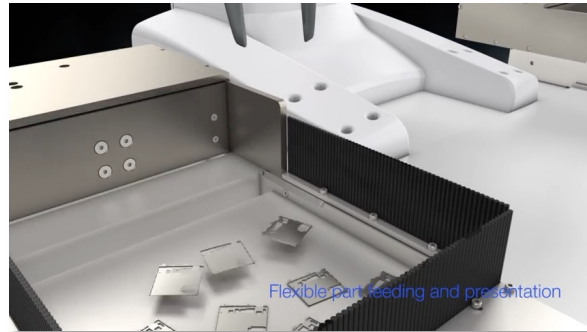
Independent
Compositional Generalization



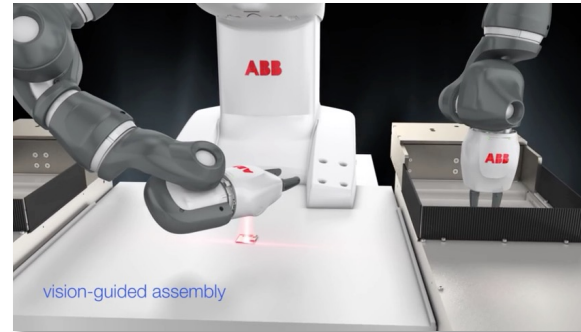
Not-Independent (!)
Hence need evidence of possibility

Data Augmentation in RL

Do more with the same data



Left Arm Pick and Place



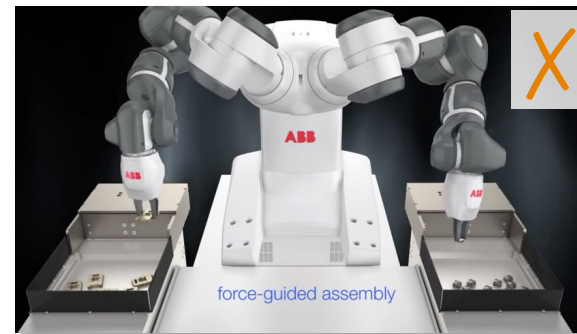
Right Arm Pick and Place



Which of the following is possible (only based on observed data)



Independent
Compositional Generalization



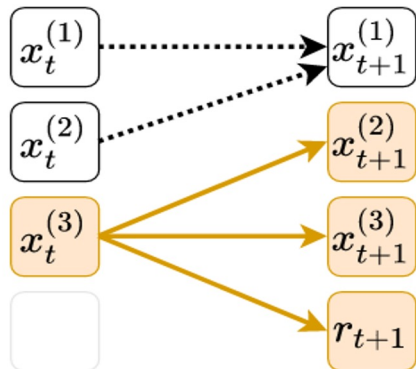
Not-Independent (!)
Hence need evidence of possibility



Counterfactual Data Augmentation

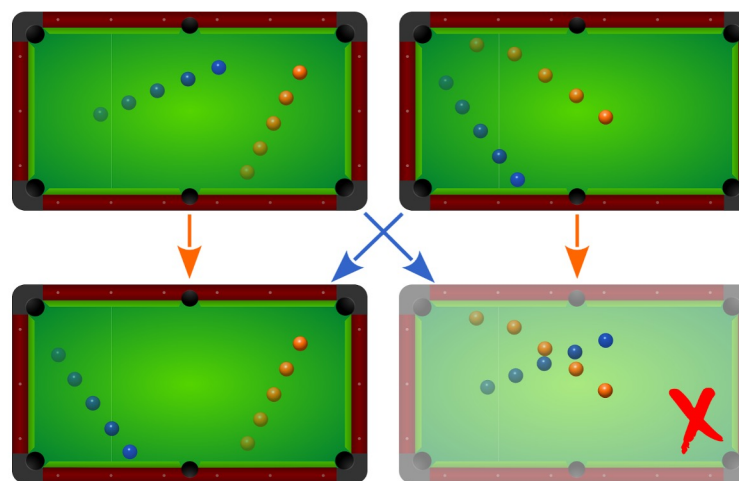
Counterfactual reasoning to generate new, causally valid (counterfactual) data!

Generic CoDA



Given two independent mechanisms, Relabel one (conditional independence!)

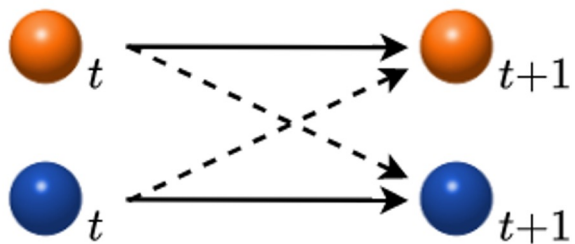
- ✓ Model-Free relabelling
- ✗ But Causal Independence is not Global



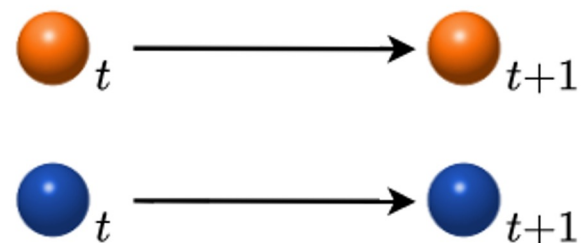
- For the most part, entities behave independently, and we can use CoDA
- But entities are not always independent, so this can also produce nonsense

Counterfactual Data Augmentation

Local Causal Model



Global Model



Local Model

$$\mathcal{M}_t = \langle V_t, U_t, \mathcal{F} \rangle \xrightarrow{\text{Condition on } (s_t, a_t) \in \mathcal{L}} \mathcal{M}_t^{\mathcal{L}} = \langle V_t^{\mathcal{L}}, U_t^{\mathcal{L}}, \mathcal{F}^{\mathcal{L}} \rangle$$

Structural Causal Model (SCM) that marginalizes across all possible transitions

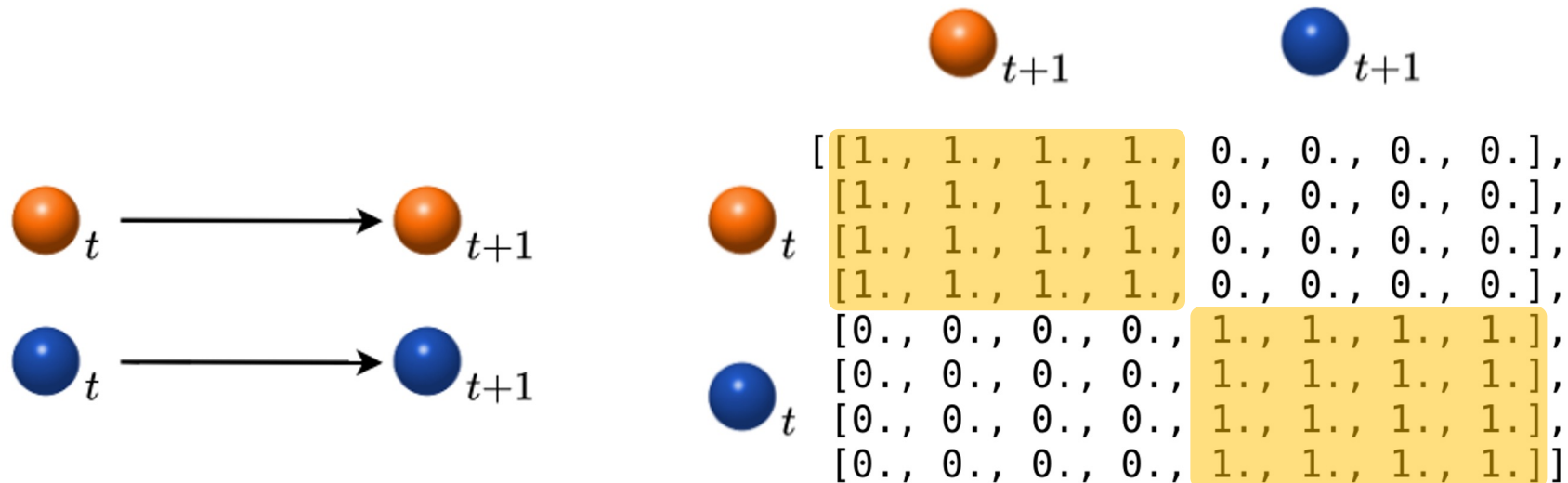
Local Causal Model (LCM) that behaves like the global SCM in local subspace \mathcal{L}

Where do local models come from?

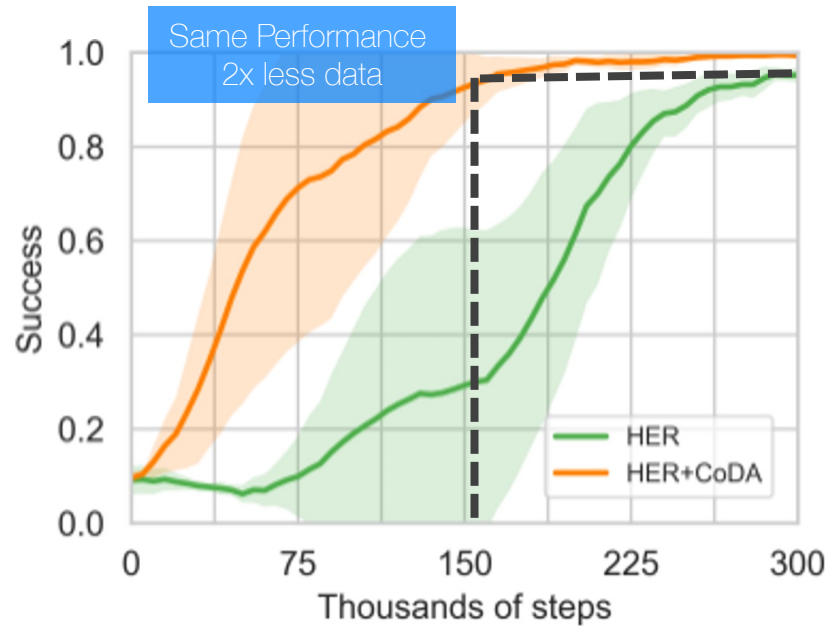
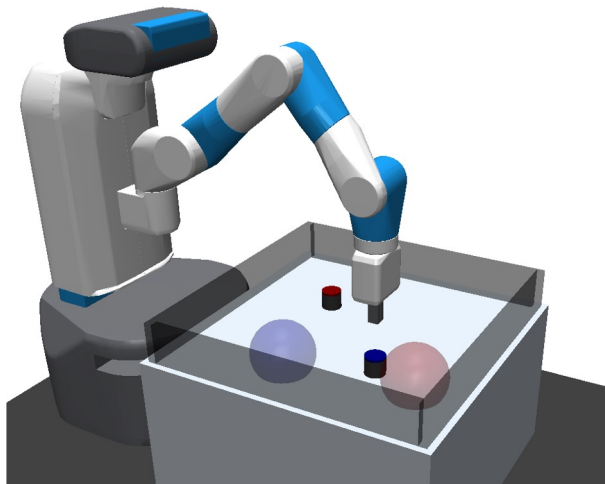
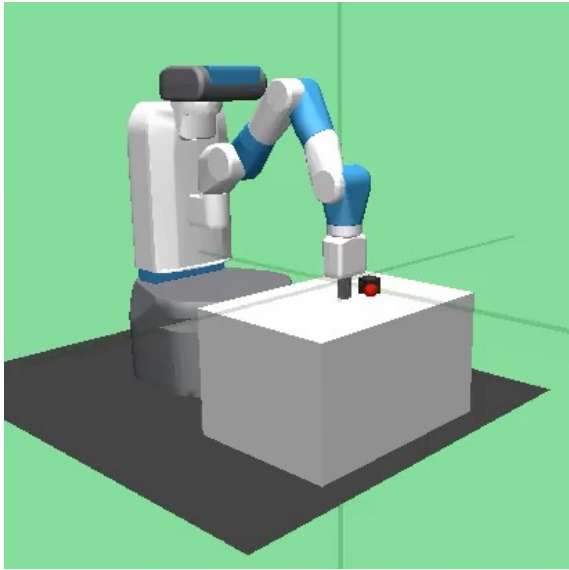
Counterfactual Data Augmentation

Learning Local Causal Model

- Input: 2 balls, each with 4 features: $[x, y, \dot{x}, \dot{y}]$
 - \bullet_t $[[1.23, -0.73, 1.31, 1.07],$
 - \bullet_t $[-0.6, 2.51, -1.51, -0.89]]$
- Output: Adjacency matrix M of the causal graph (between \mathbf{x}_t and \mathbf{x}_{t+1})
- (intuition) M : the input-output Jacobian is non-zero



CoDA: Goal-Conditioned (Online) RL

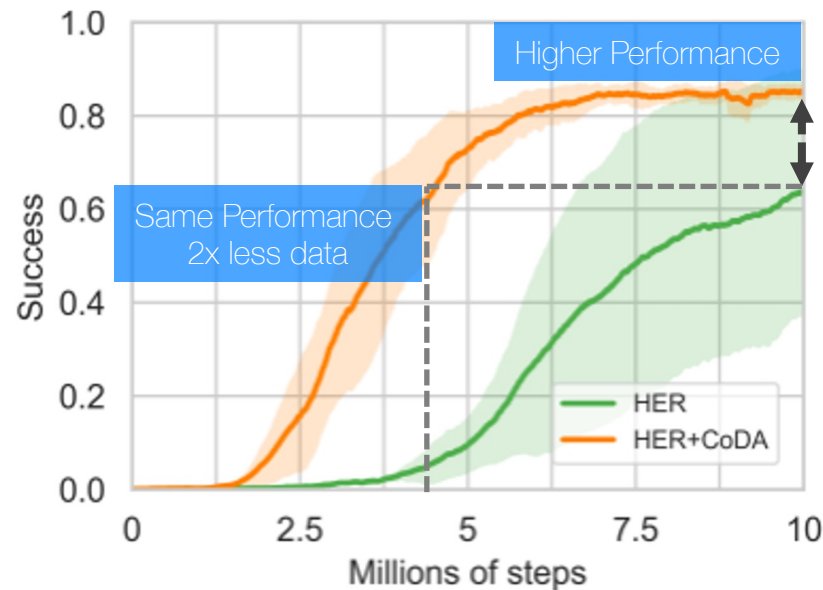


Fetch-Push-v1

state space: [Robot and 1 object]

CoDA Heuristic for independence:

Obj \perp Gripper | if $|x_g - x_o| \geq 10 \text{ cm}$



Fetch-Slide2

state space: [Robot and 2 objects]

Harder task (30x more samples)!

Structure for RL: Off-policy RL



$$\min_Q \mathbb{E}_{a \sim \pi_{\text{Behavior}}} [\mathcal{L}(Q(s, a) - Q_{\text{Target}})]$$
$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_{\text{new}}} [Q(s', a')]$$

Update Rule

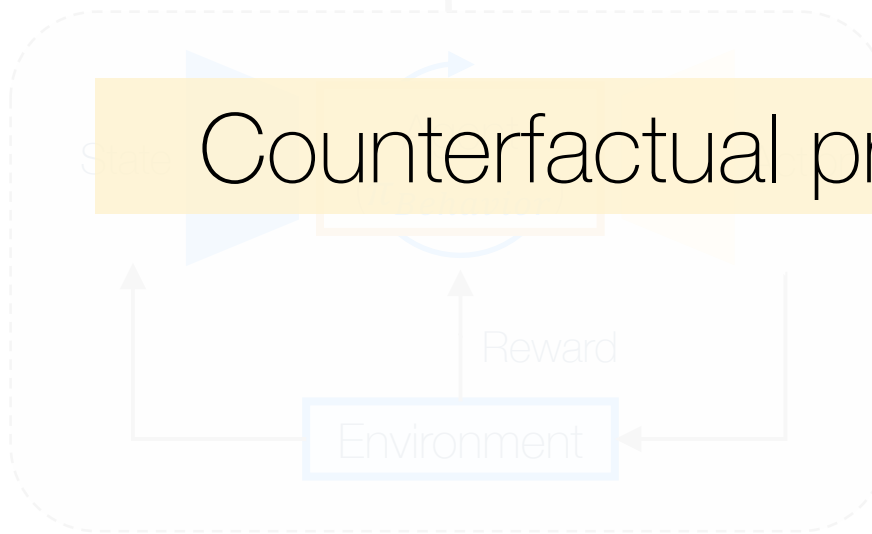
Counterfactual predictions help Data Efficiency

$$\pi_{\text{new}} = \begin{cases} 1 - \epsilon, & \text{if } a_t = \operatorname{argmax}_{a_t} Q(s, a) \\ \frac{\epsilon}{(|\mathcal{A}| - 1)} & \end{cases}$$

Policy Improvement

Policy Rollout

(Generate Samples with Interaction)

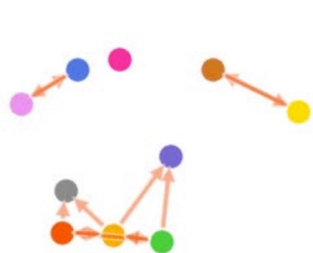


Structure for RL: Off-policy RL



Counterfactual Data Augmentation helps
How to learn this structure automatically?

Predicted graph



Predicted keypoint movements



Ground truth keypoint movements



$\pi_t = \arg\max_a Q(s_t, a)$
Discover Causal Dynamics Structure from Visual Data

Improvement

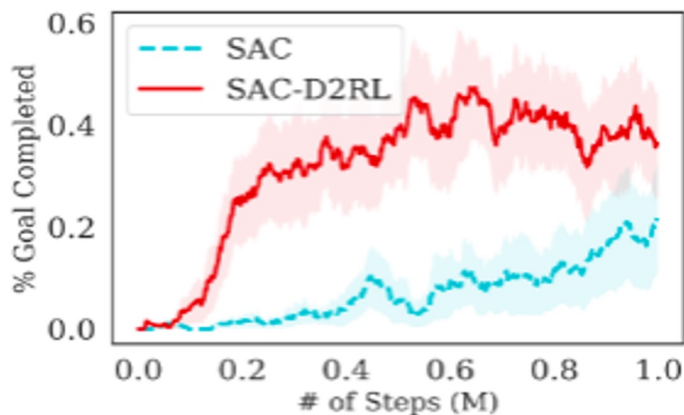
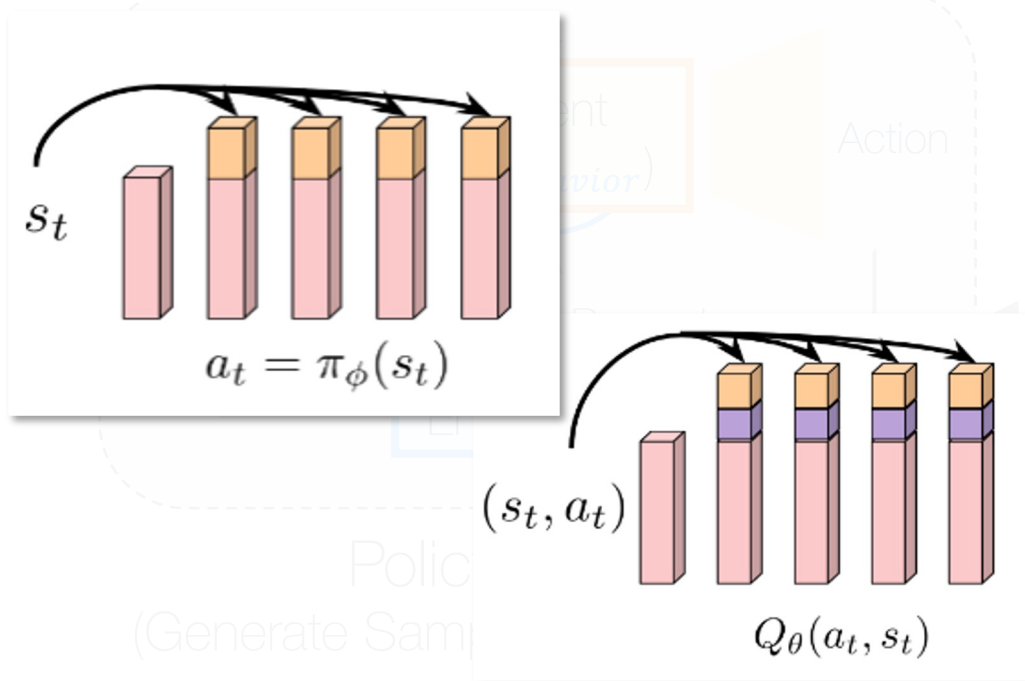
(Generate samples with interaction)

Structure for RL: Off-policy RL

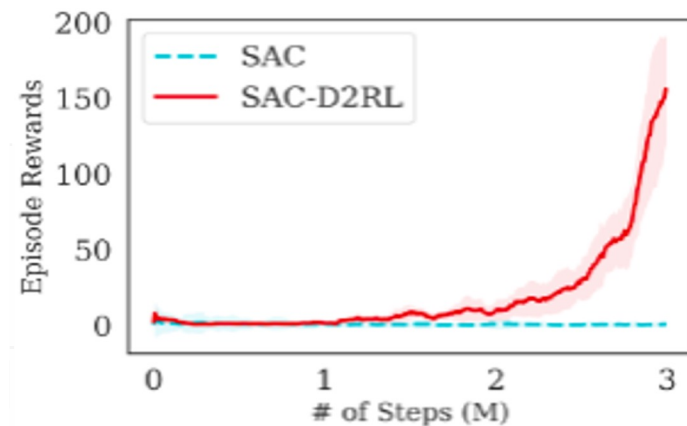
Does the choice of architecture matter?

$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_{new}} [Q(s', a')]$$

Using Dense connections in Policy/Value improves sample efficiency



(b) Fetch Slide SAC



(d) Jaco Reach SAC

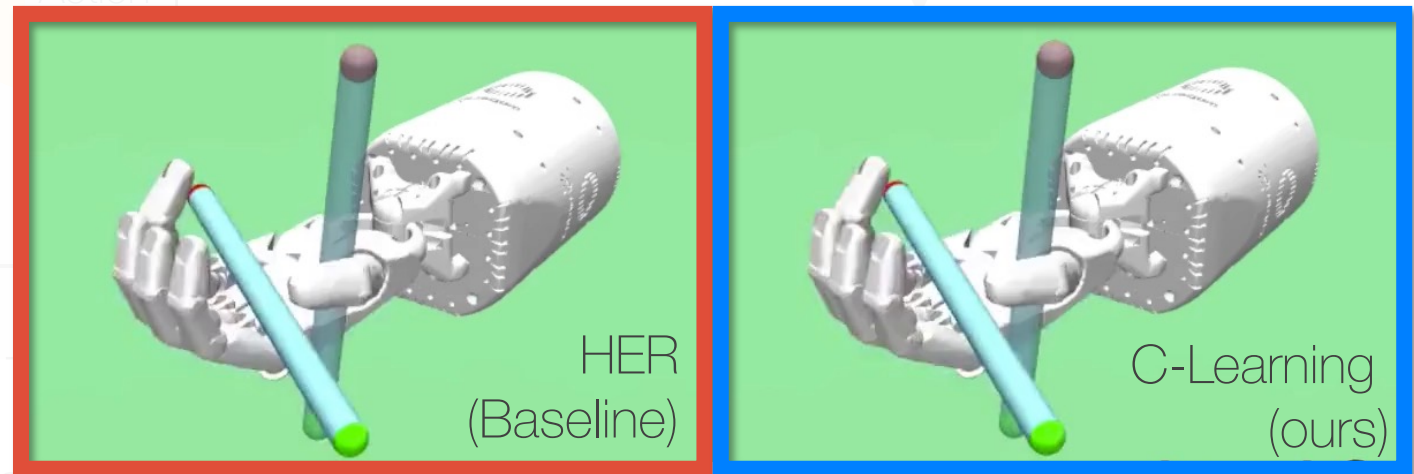
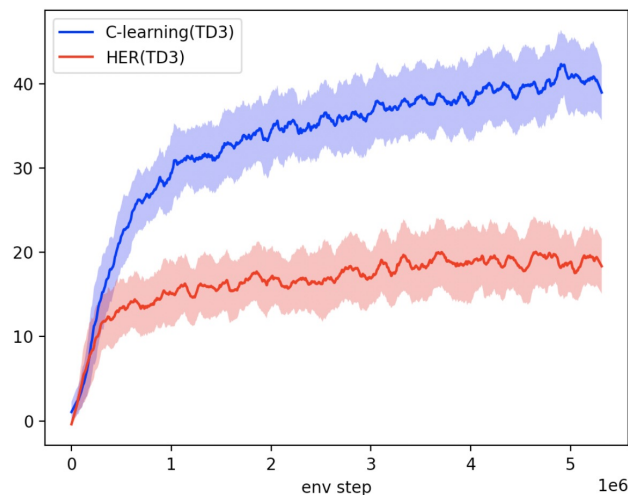
Structure for RL: Off-policy RL

Can we use better Utility Functions?

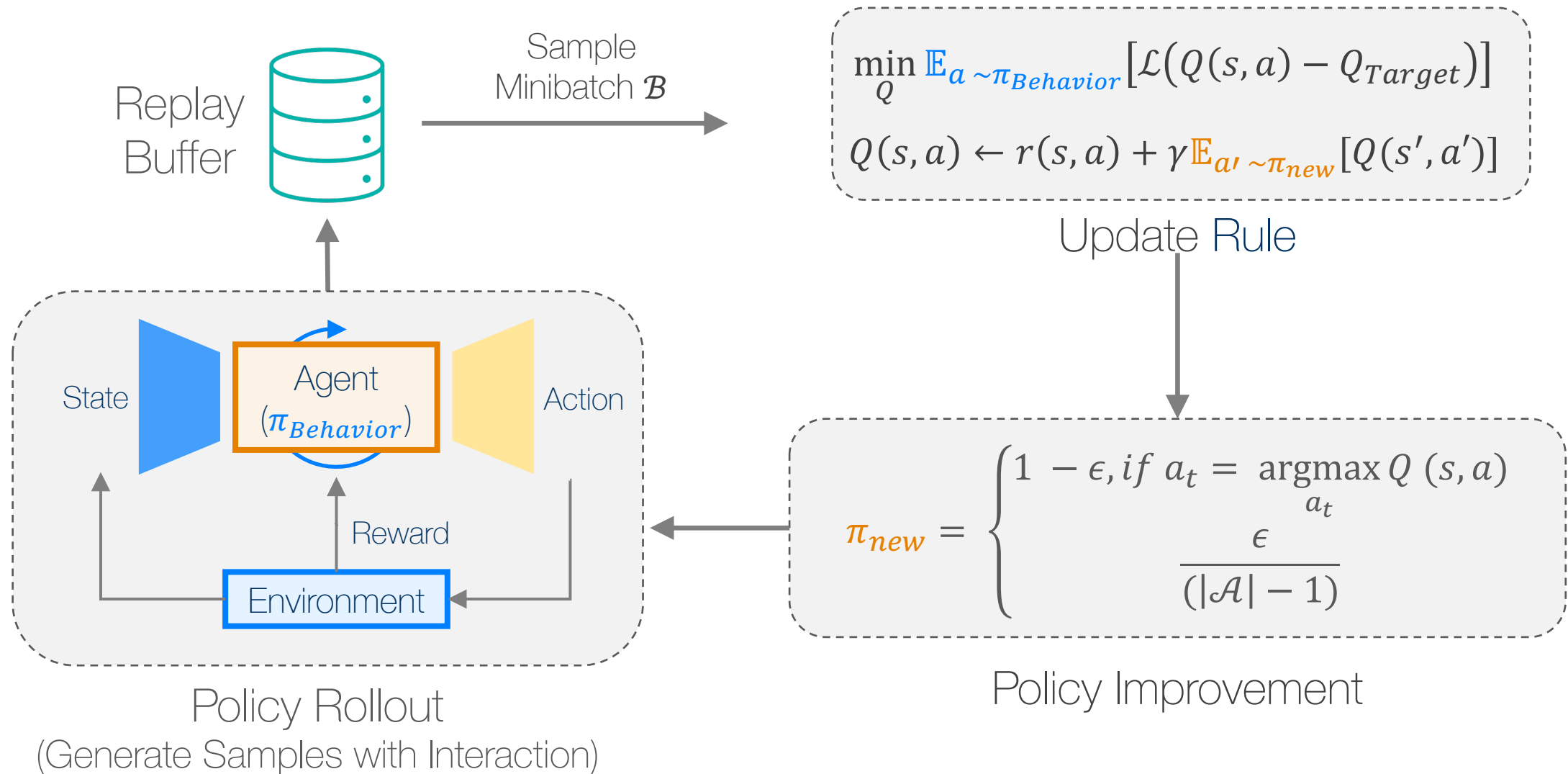
$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_{new}} [Q(s', a')]$$

Learning Cumulative Accessibility $\mathcal{C}(s, a, h)$ is better than $Q(s, a)$

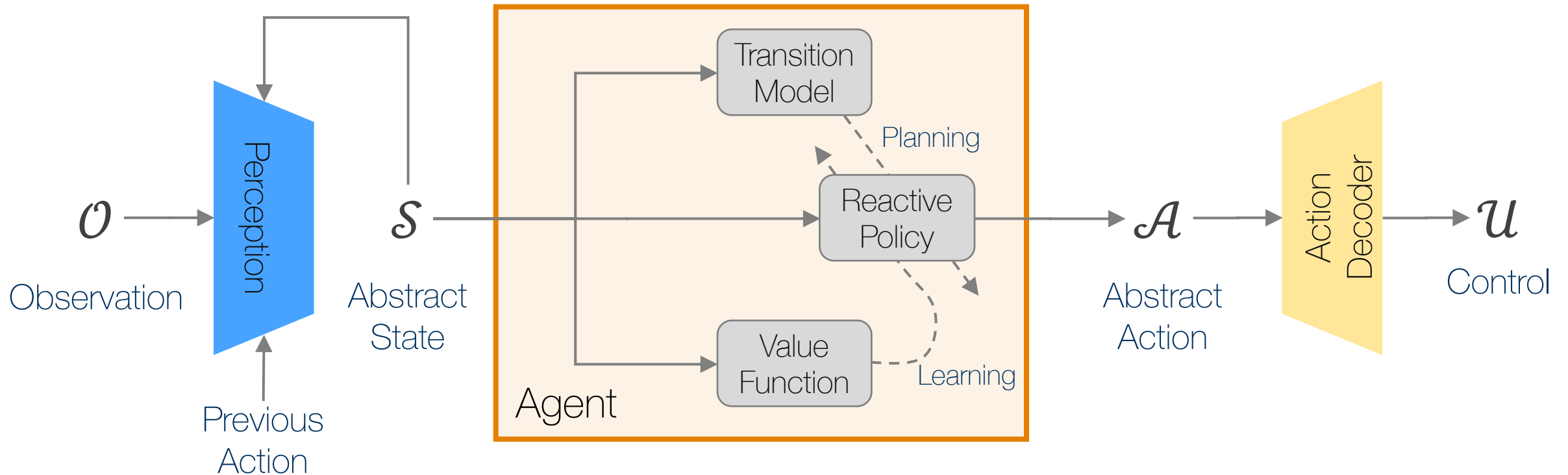
Can represent multimodal, multi-goal, horizon-aware solutions as well as reachability



Structure for RL: Off-policy RL



Structure for Reinforcement Learning



Update Rule

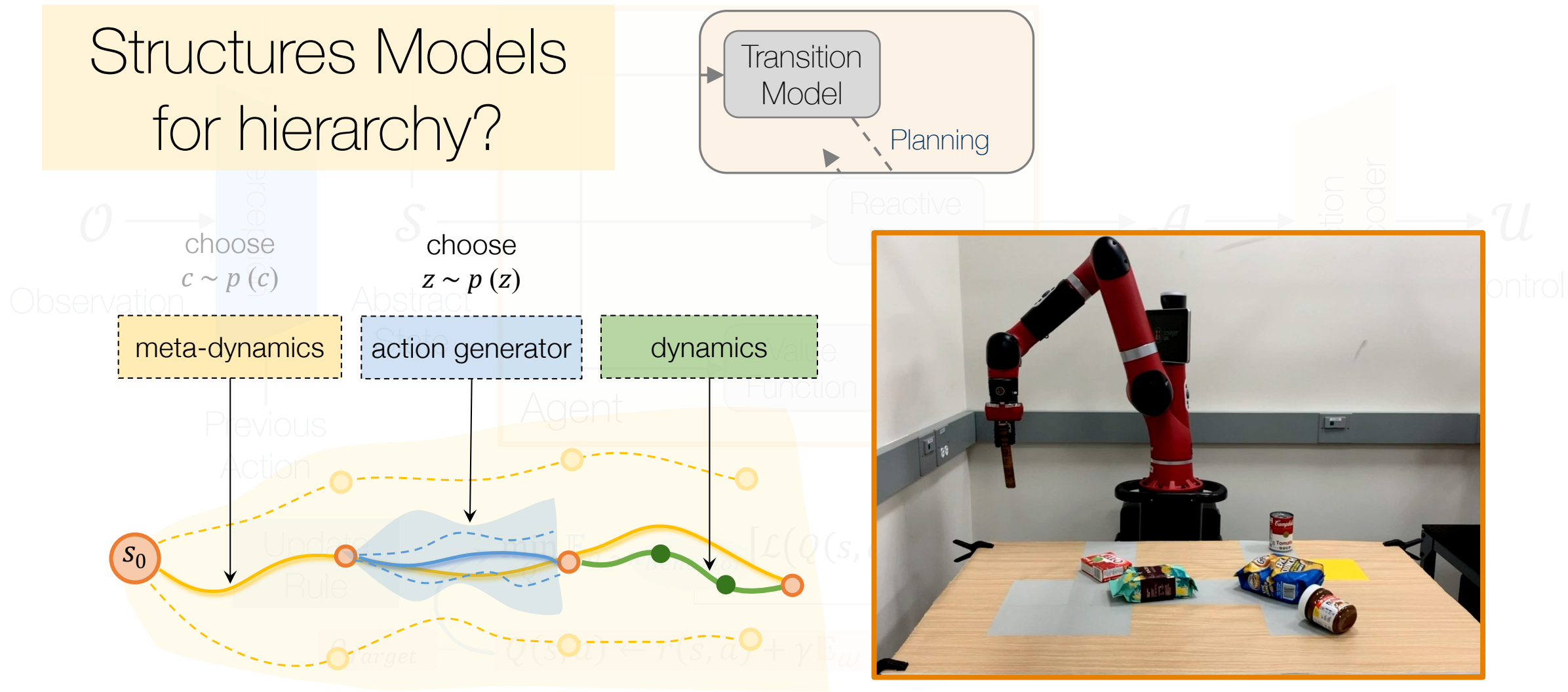
$$\min_Q \mathbb{E}_{a \sim \pi_{Behavior}} [\mathcal{L}(Q(s, a) - Q_{Target})]$$

$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_{new}} [Q(s', a')]$$

Not the same Policies

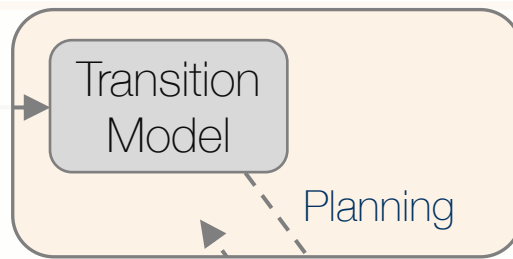
Structure for Reinforcement Learning

Structures Models
for hierarchy?

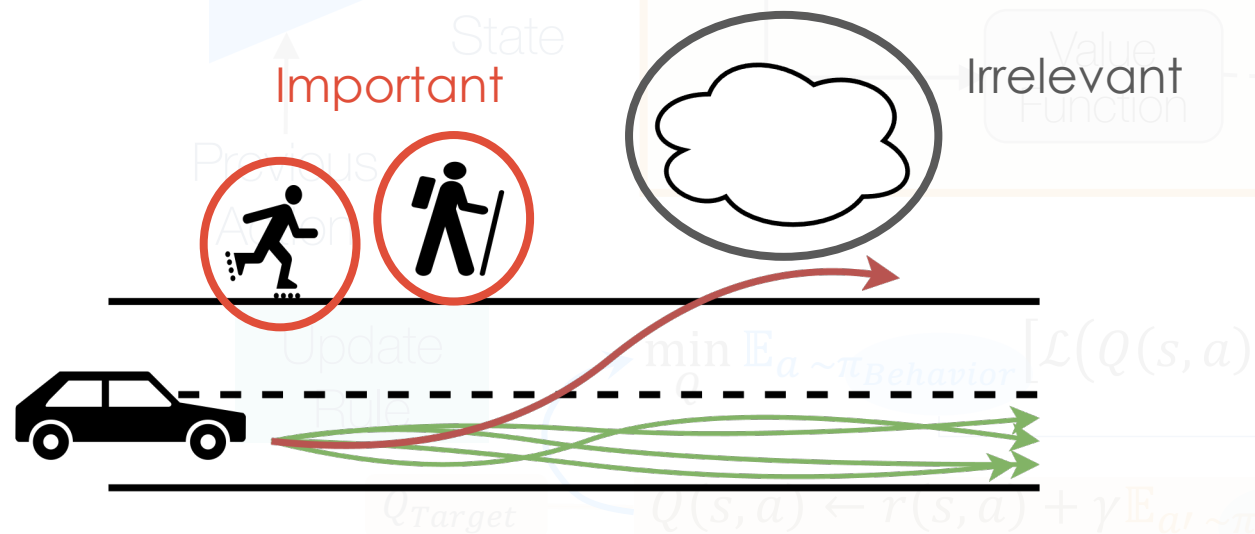


Structure for Reinforcement Learning

Dynamics Prediction,
Correct objective?

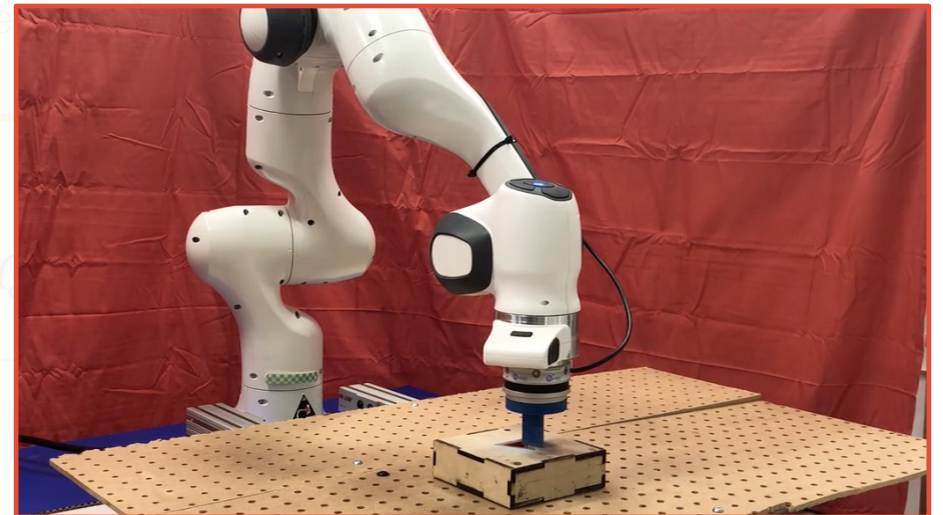
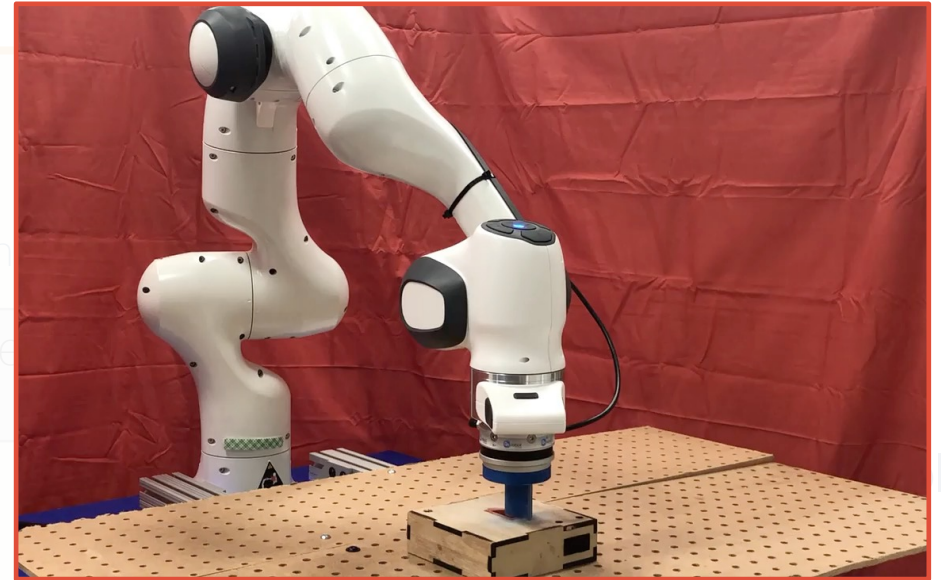
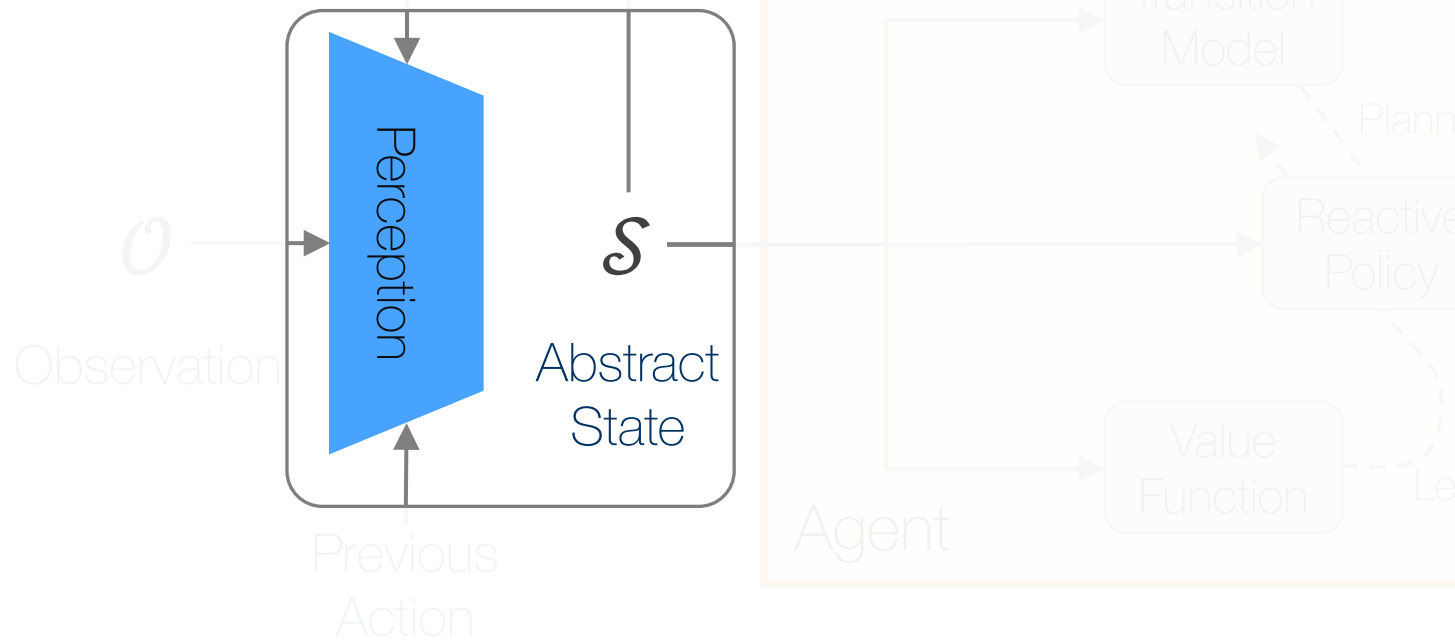


Task-Aware Objective helps efficient and targeted Model-Learning



Structure for Reinforcement Learning

How can better state representations capture multimodal data?



Generalizable Multi-modal State Representations

- Learn a joint Visuo-Tactile representation for Peg Transfer
- Representation transfers to new task, while Policy doesn't

Q_{target}

$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_{new}} Q(s, a')$$

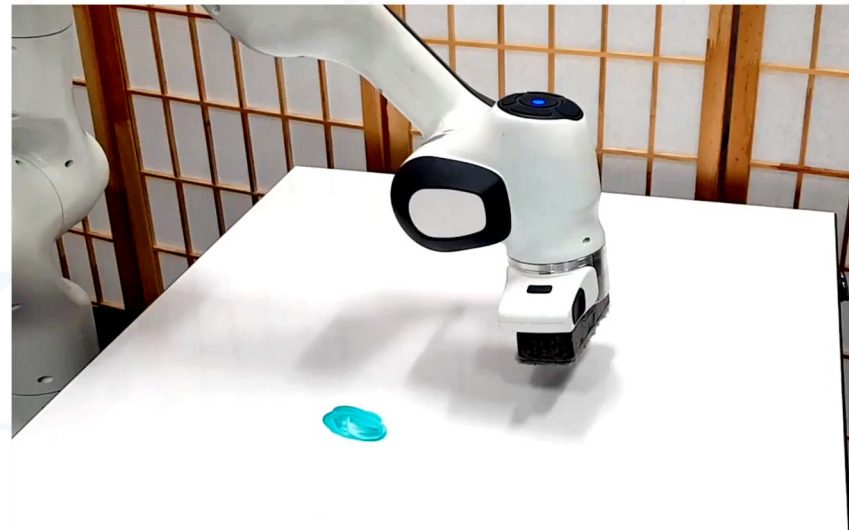
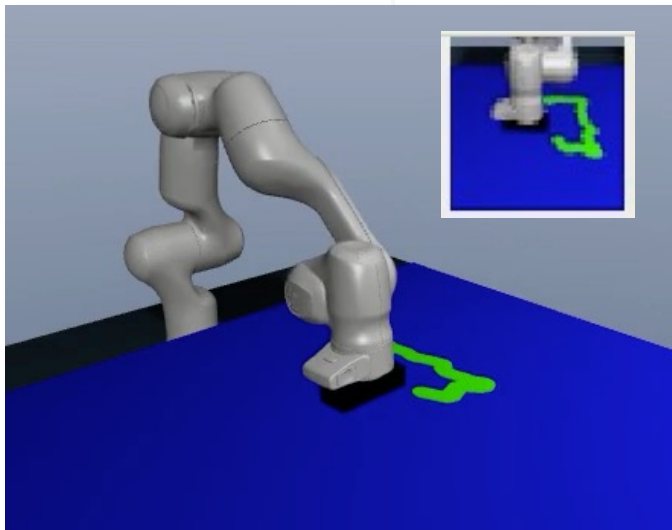
Structure for Reinforcement Learning

How can better action representations result in generalization?

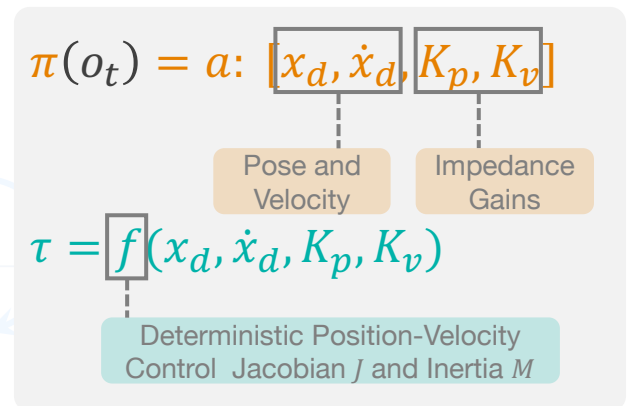
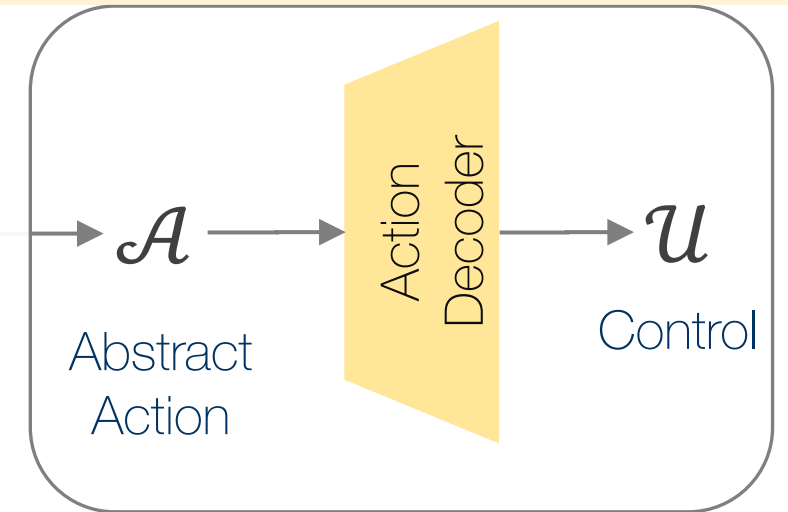
Variable Impedance Controller in End-Effector Space

- Efficiency in RL
- Ease of Sim2Real Transfer

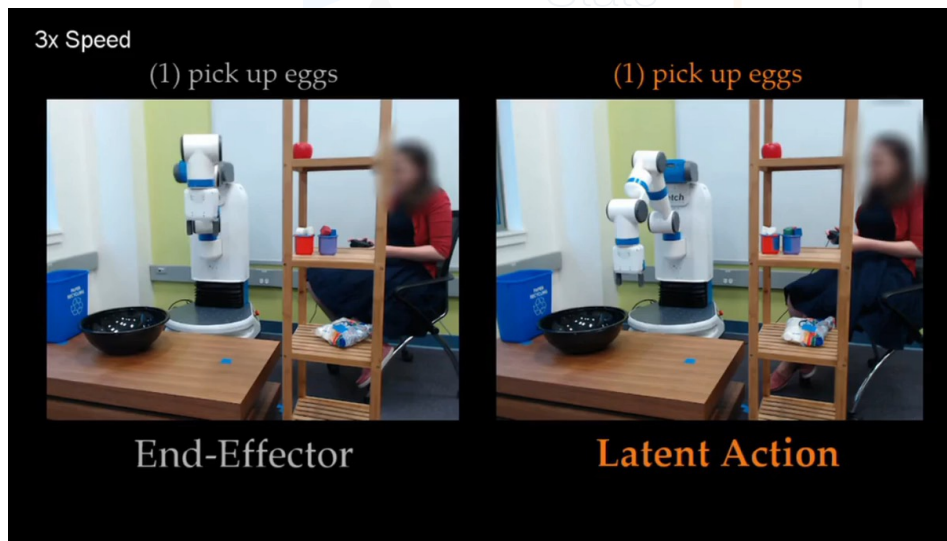
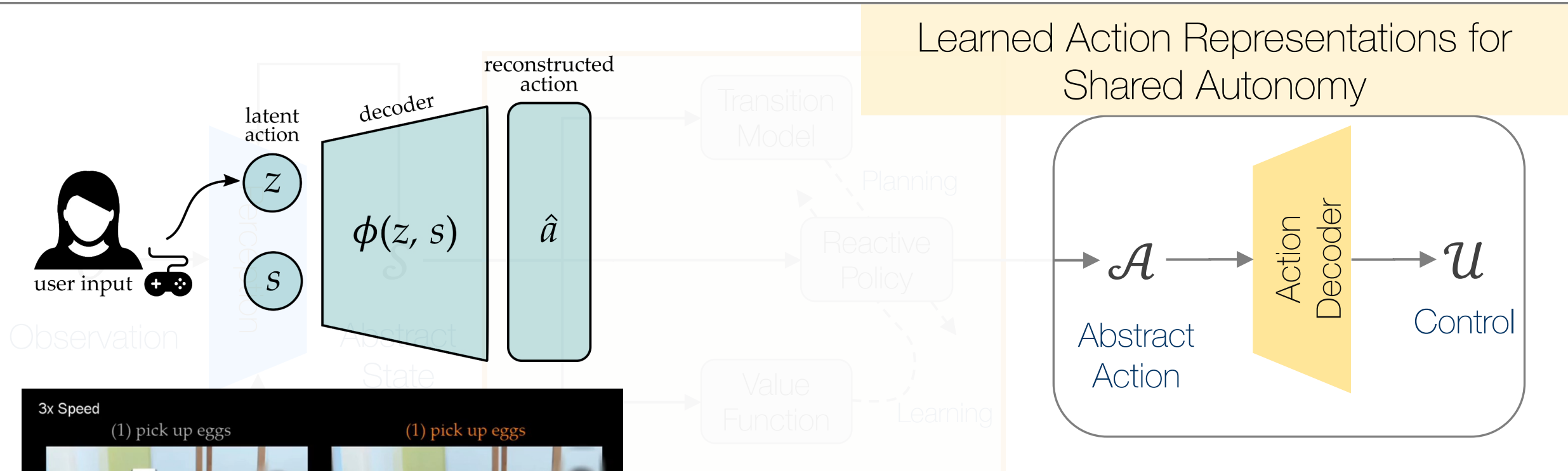
Observation $\tau = f_{Sim}(\pi(o_t))$ State



$\tau = f_{Real}(\pi(o_t))$



Structure for Reinforcement Learning



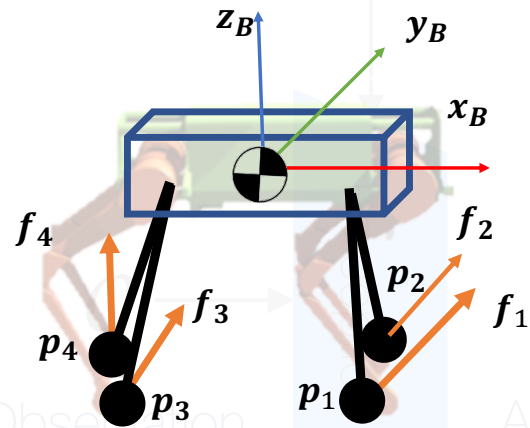
Learned Action Space

Easier to control **high-dimensional robots** by embedding the robot's actions into a **low-dimensional latent space**

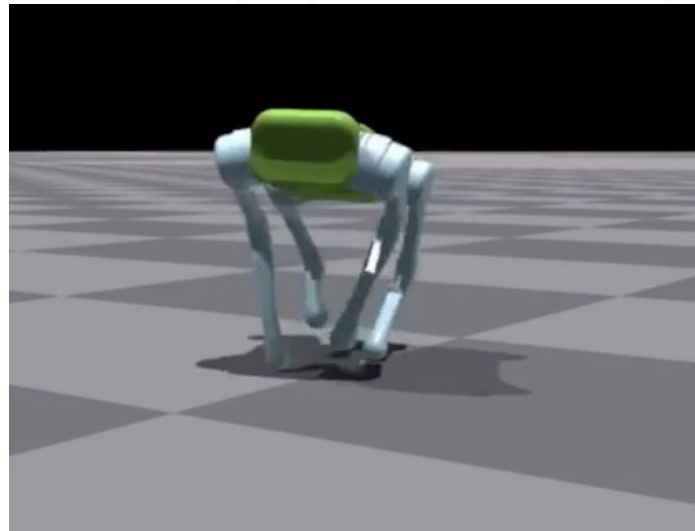
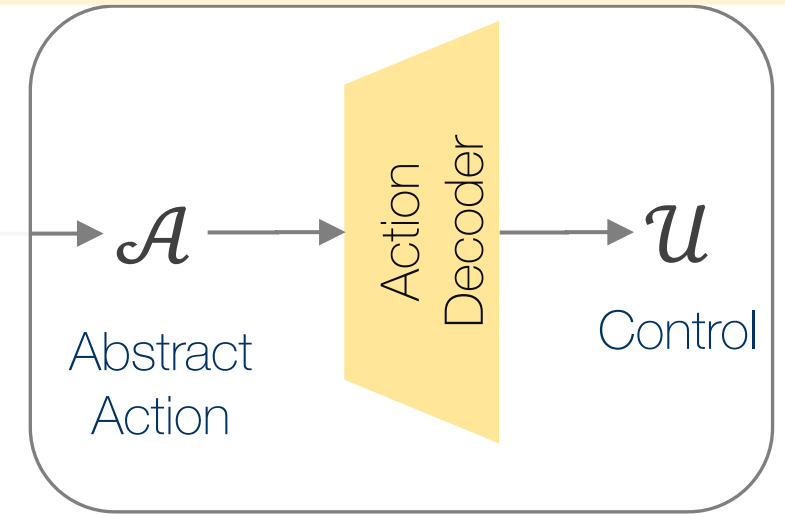
Structure for Reinforcement Learning

Centroidal Task Space

- Easier for learning
- RL + Optimal Control

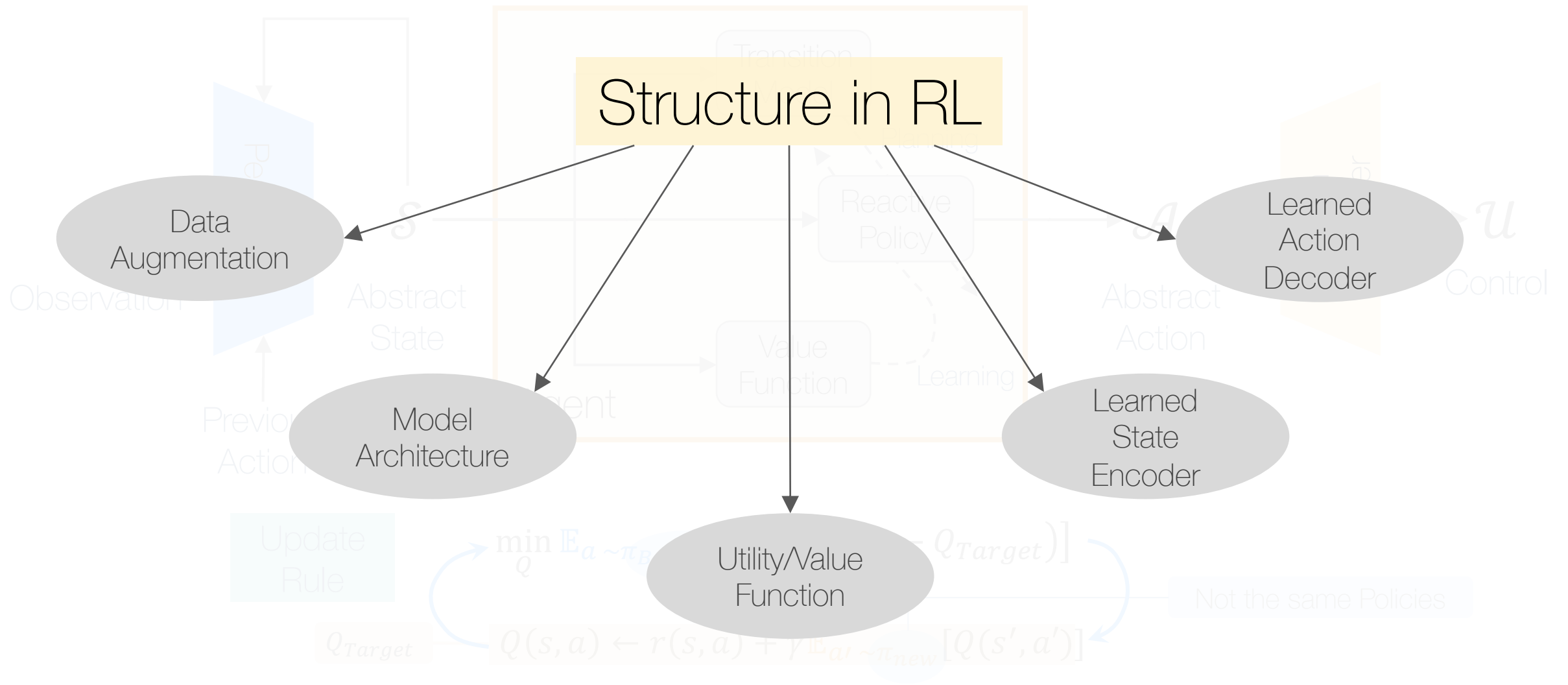


Action representations for Legged Locomotion?

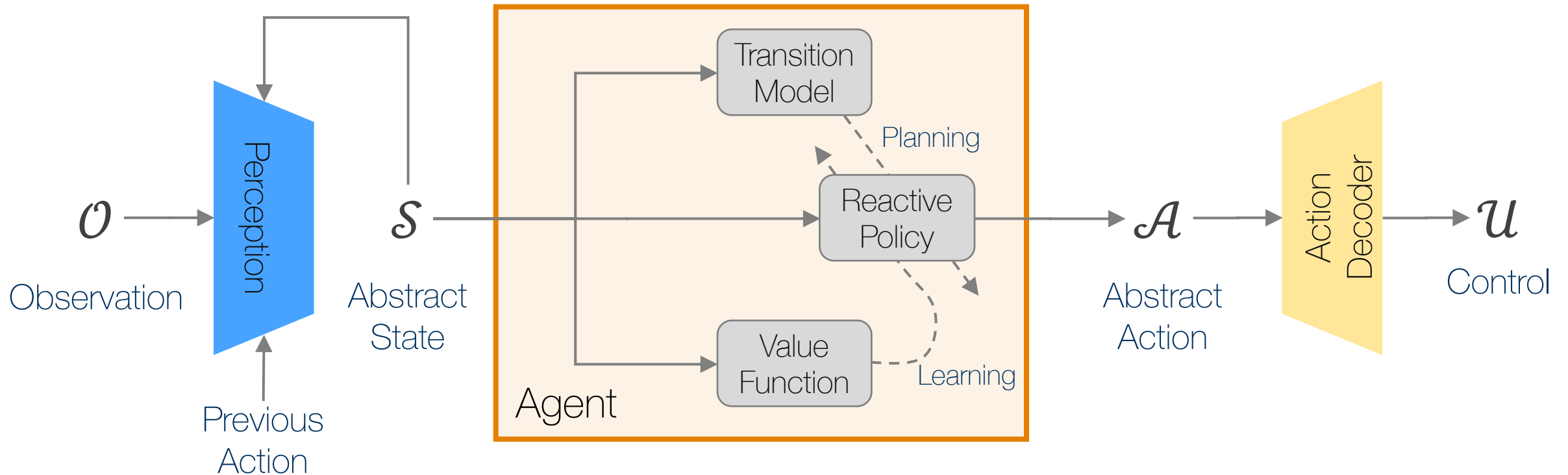


Sim2Real

Structure for Reinforcement Learning



Structure for Reinforcement Learning



Update Rule

$$\min_Q \mathbb{E}_{a \sim \pi_{Behavior}} [\mathcal{L}(Q(s, a) - Q_{Target})]$$

$$Q(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{a' \sim \pi_{new}} [Q(s', a')]$$

Not the same Policies

Structure in Compositional Planning

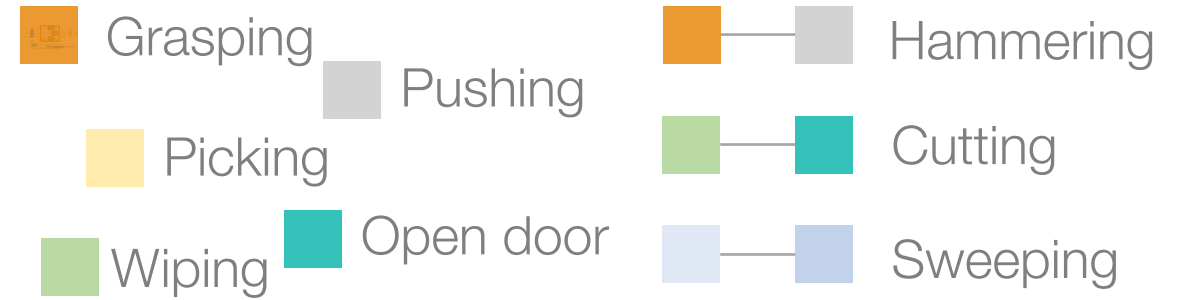


Visuo-Motor Skills

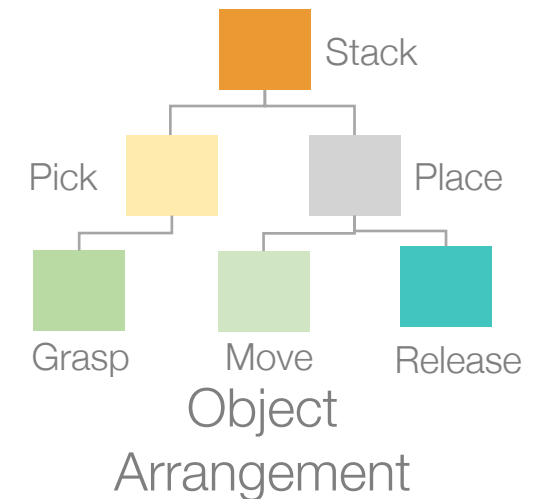
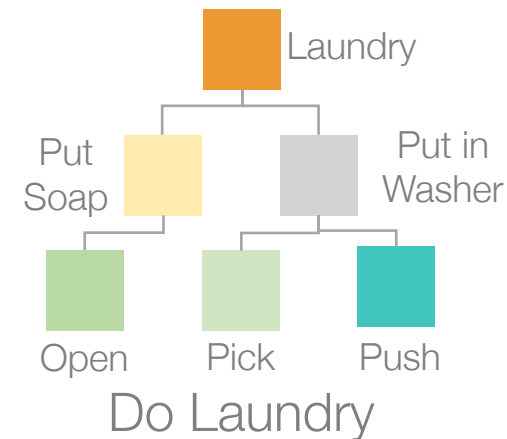


Compositional Planning

Visuo-Motor Skills



Compositional Planning



Structure in Compositional Planning

Imitation: But at which level? What should I copy?



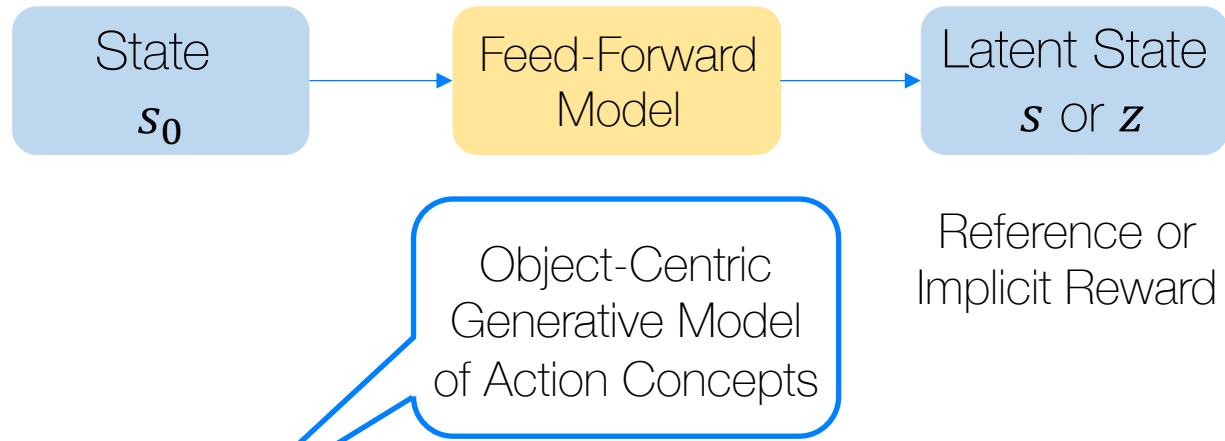
Imitative Babbling
Movement Skills

Dexterity
Skill Sequencing

Causality
Semantic Purpose

Task Specification
{Language, Video,
Kinesthetic}

Structure in Compositional Planning



Causal Generative Model

- Learn to predict the “effect” of “action”
- Compositional & Counterfactual
- Multi-step Semantic consistency
- Pre-trainable over large problem settings



What does it mean to “open” a “door”?

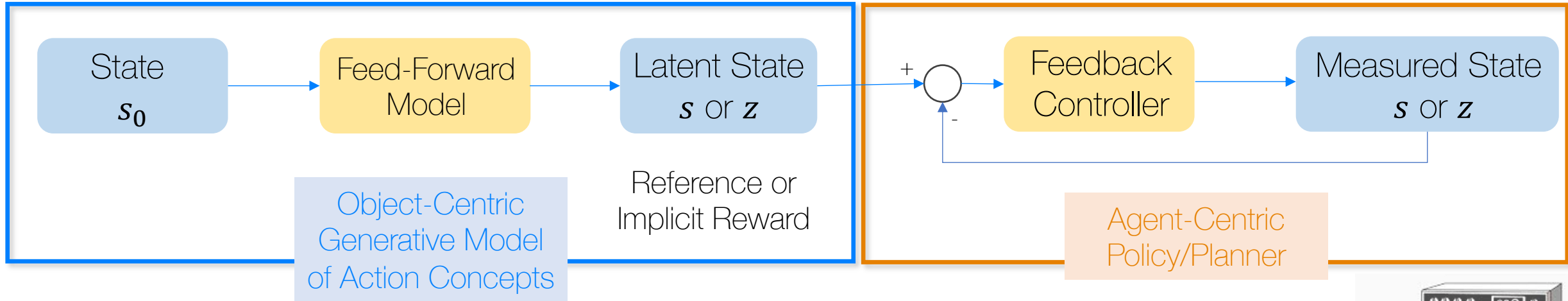
“open” a “jar”?

“open” <🏠>

- “Take” “Jug”
- “Open” “Fridge”
- “Put” “Jug” in “Fridge”



Structure in Compositional Planning



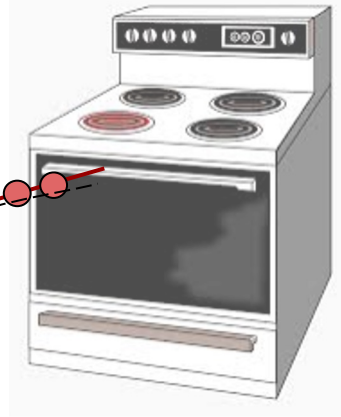
Input

- "Take" "Jug"
- "Open" "Fridge"
- "Put" "Jug" in "Fridge"

Goal Generation

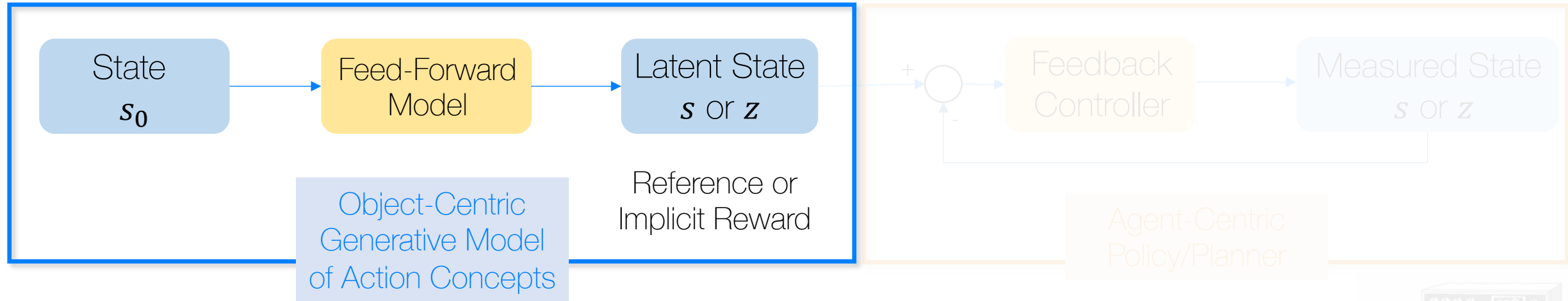


Goal-conditioned Reactive controller



Solvable online for different agents

Structure in Compositional Planning



Input

- "Take" "Jug"
- "Open" "Fridge"
- "Put" "Jug" in "Fridge"

Goal Generation



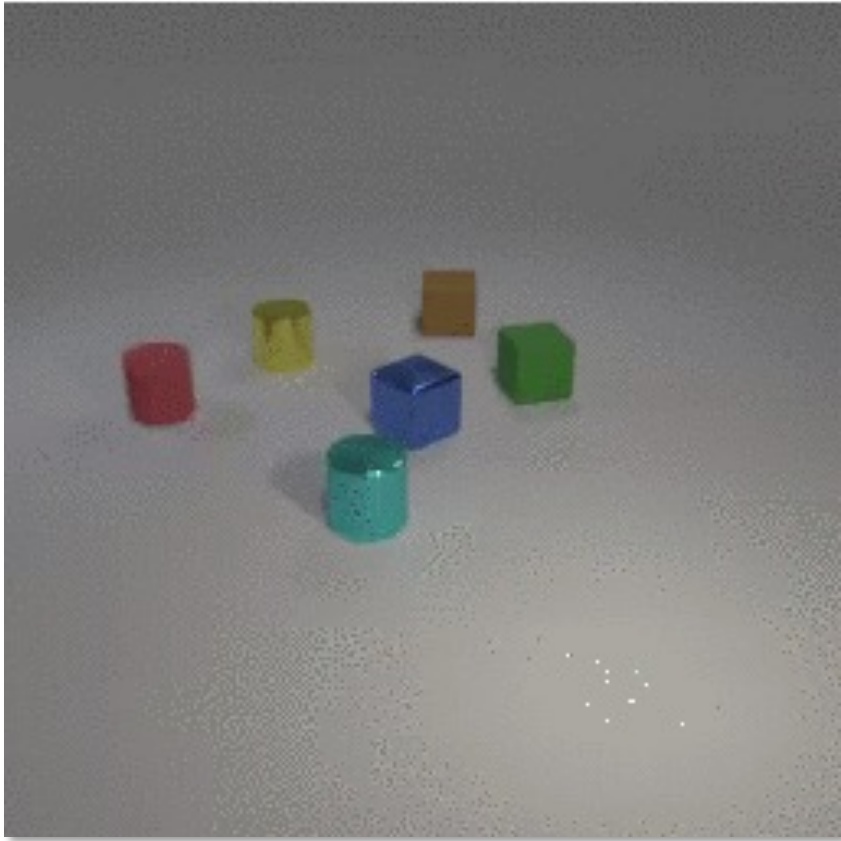
Goal-conditioned Reactive controller



Solvable online for different agents

Object-Centric Causal Generative Model

Semantic + Action-Conditional



Semantic action-conditional video prediction

Self-Supervised Modular Object Representation

Long-term Semantically Consistent Predication

No bounding box or object level supervision.

Prompt: Sequential Language Instruction

Object-Centric Causal Generative Model

Modular Action Concepts

Input: $t=1$

- "Take" "Jug"
- "Open" "Fridge"
- "Put" "Jug" in "Fridge"



- "Pick up" "Green Bowl"
- "Place in" "large pink bowl"



Ground Truth
Instruction

Ground
Truth

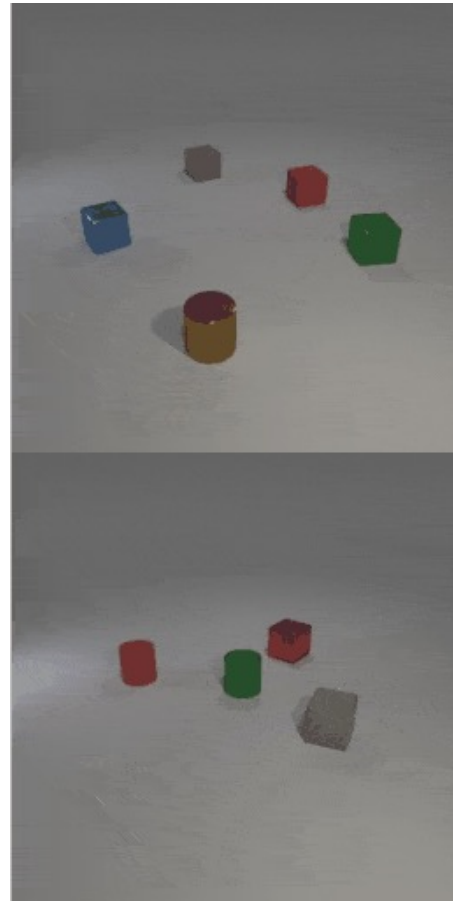
MAC
Prediction

Object-Centric Causal Generative Model

Systematic Generalization: Out of Distribution

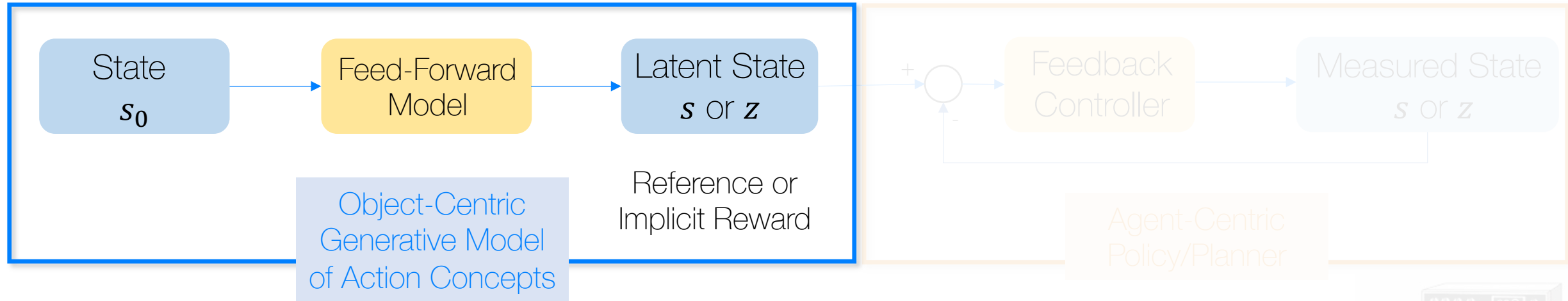


All **Red** cubes are removed
from training data



Testing:
Concurrent actions
Training:
Single action

Structure in Compositional Planning



Input

- "Take" "Jug"
- "Open" "Fridge"
- "Put" "Jug" in "Fridge"

Goal Generation

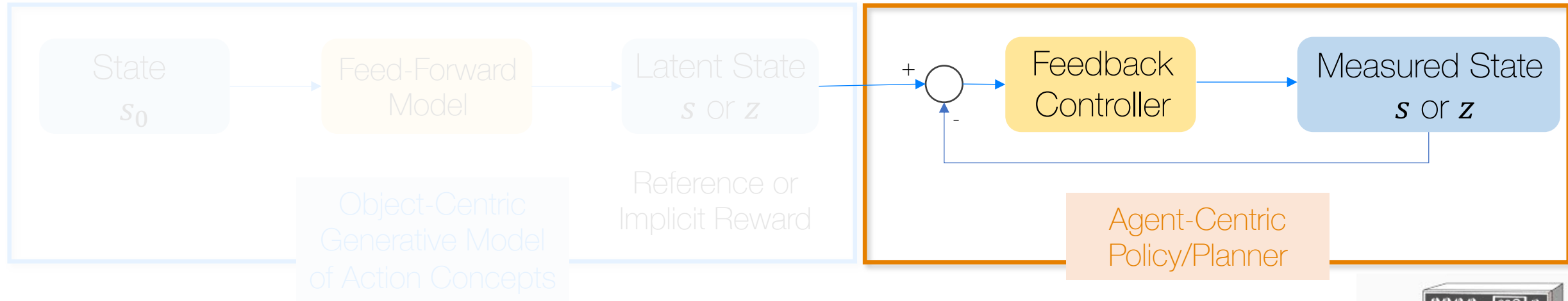


Goal-conditioned
Reactive controller



Solvable online for
different agents

Structure in Compositional Planning



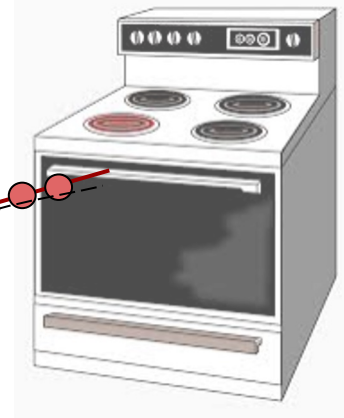
Input

- "Take" "Jug"
- "Open" "Fridge"
- "Put" "Jug" in "Fridge"

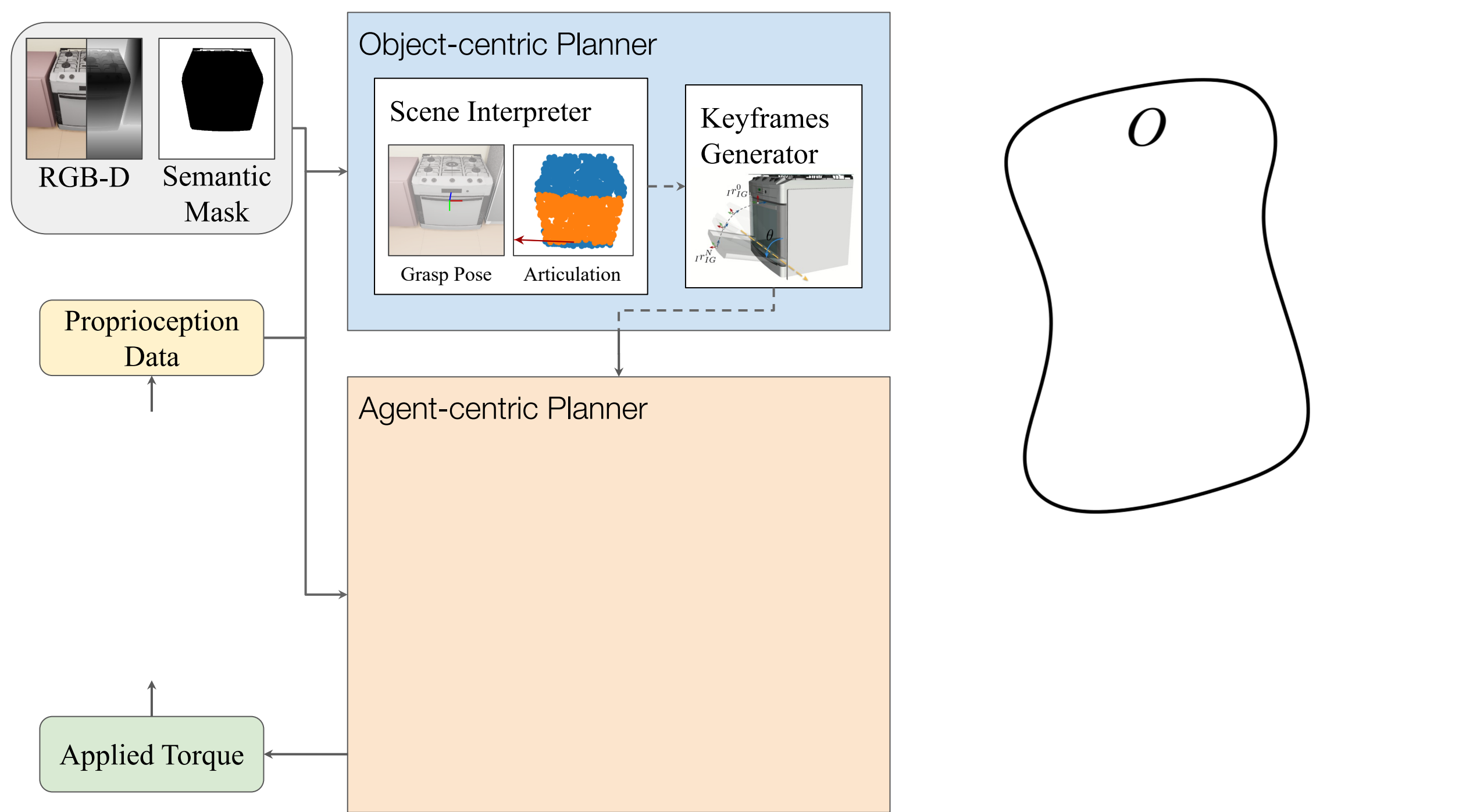
Goal Generation

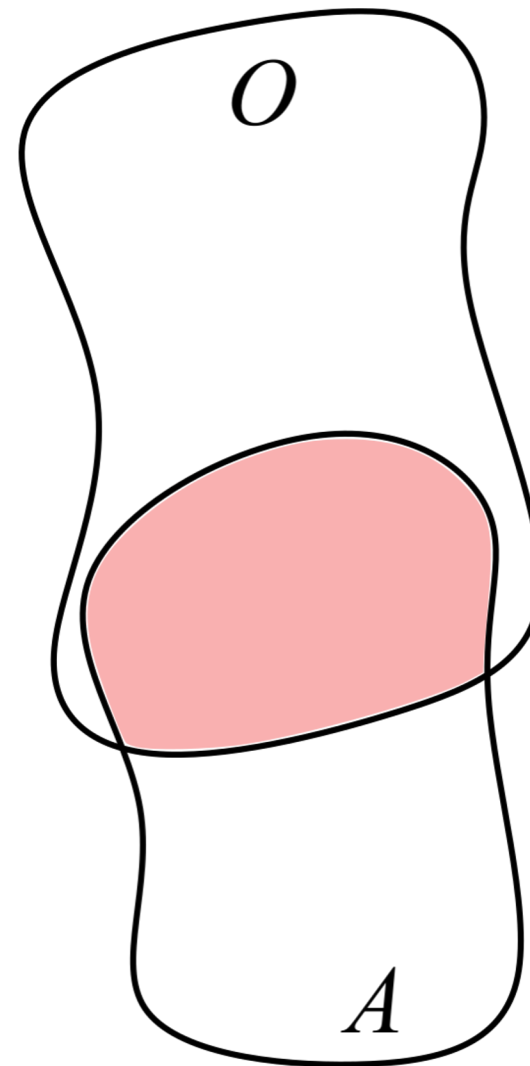
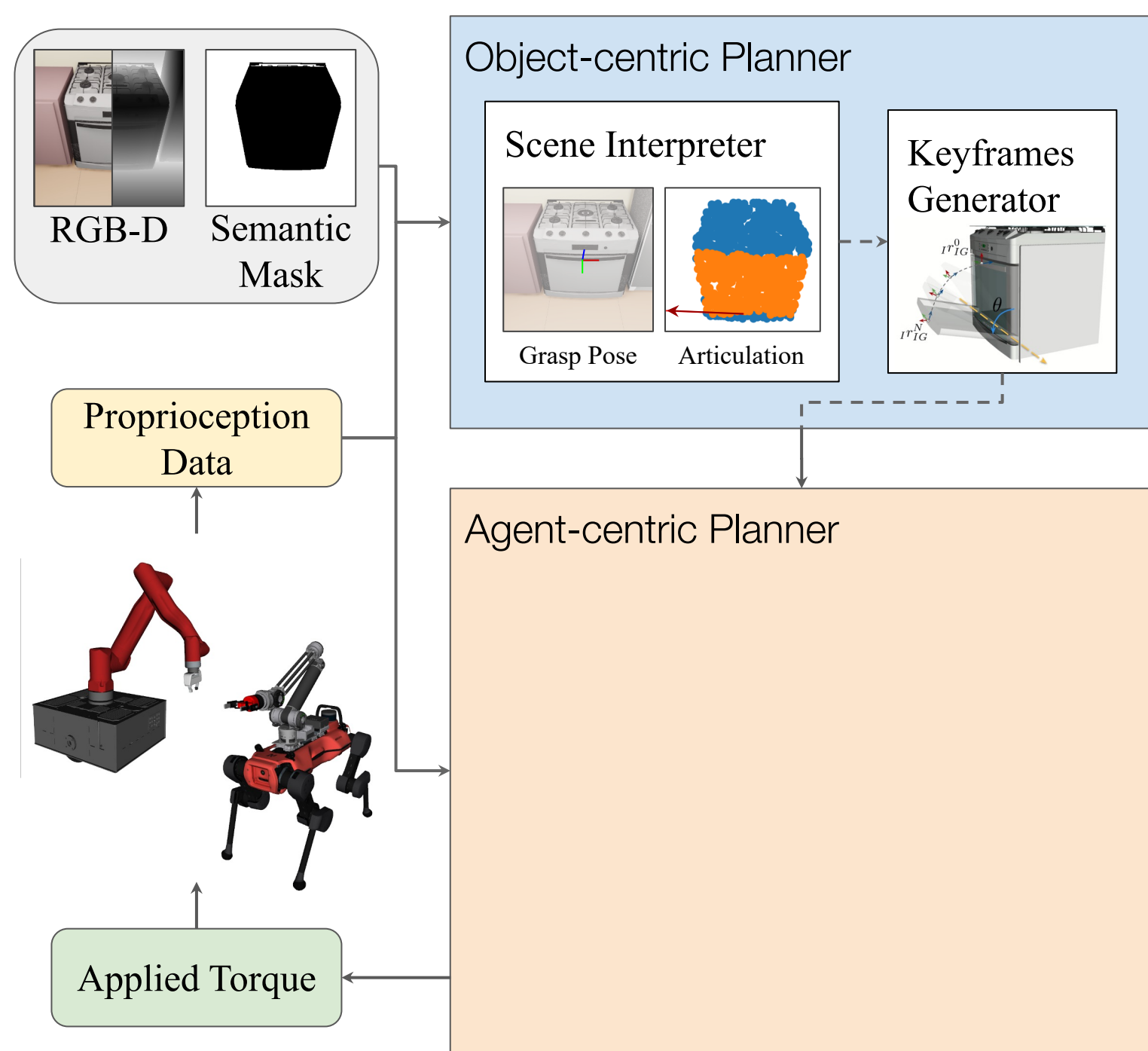


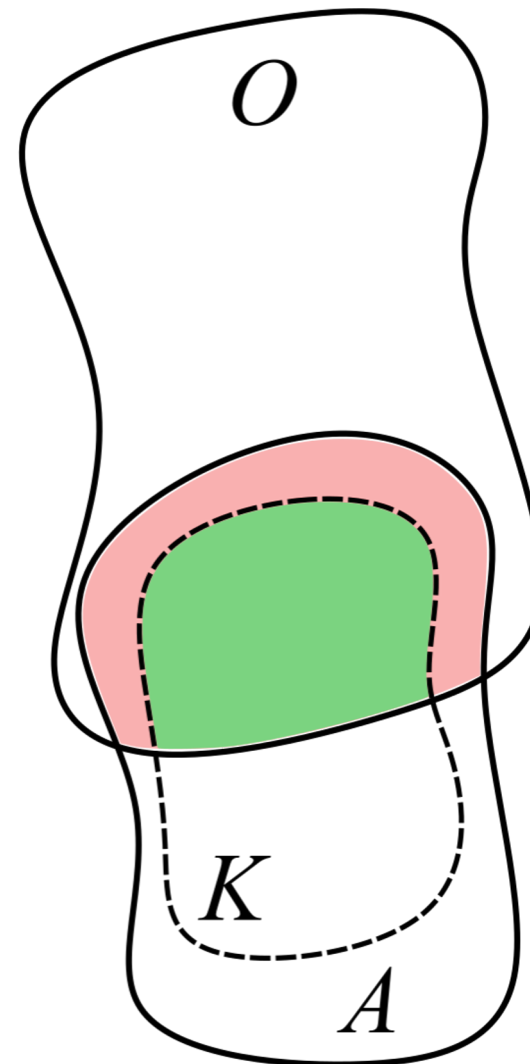
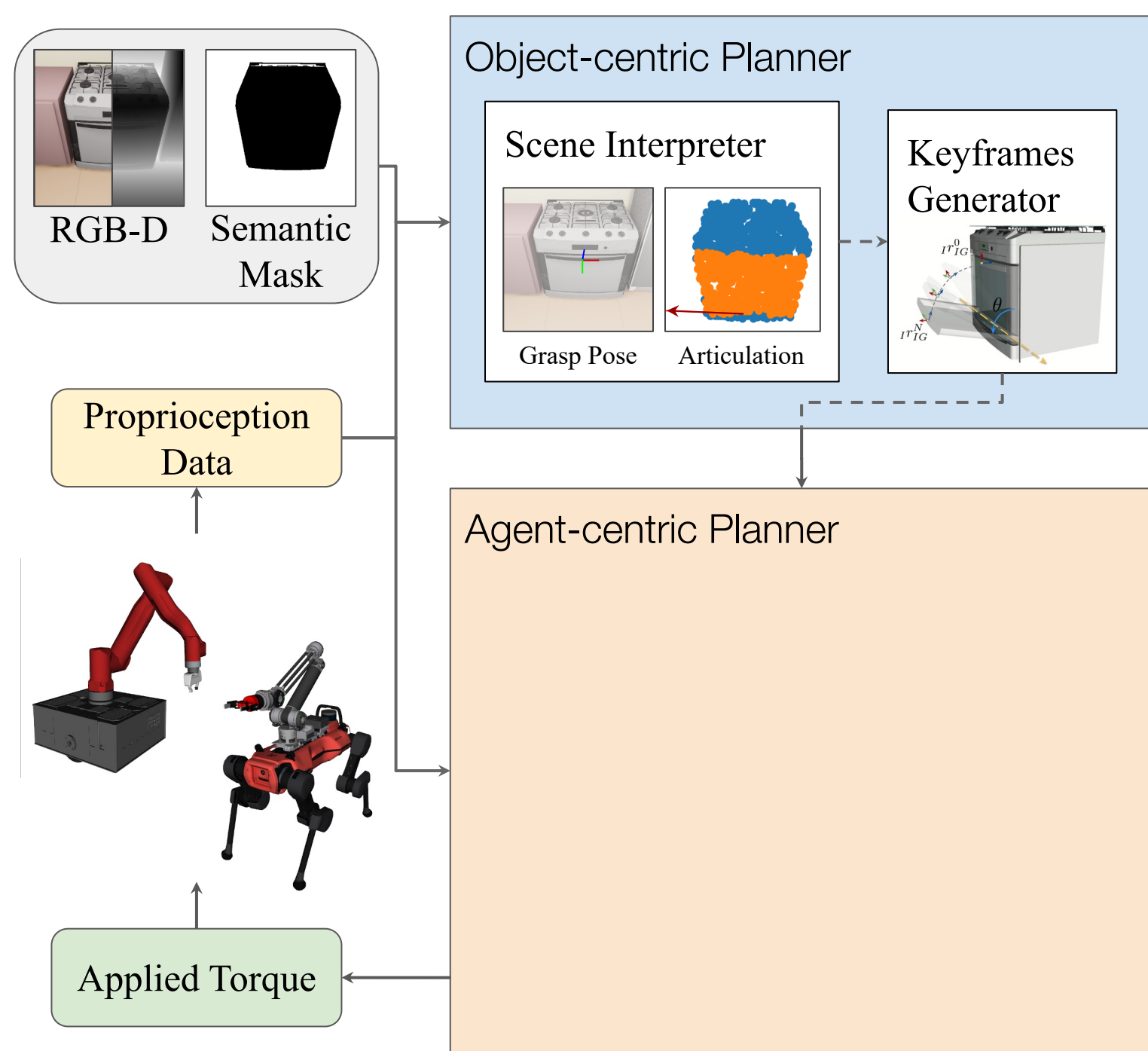
Goal-conditioned Reactive controller

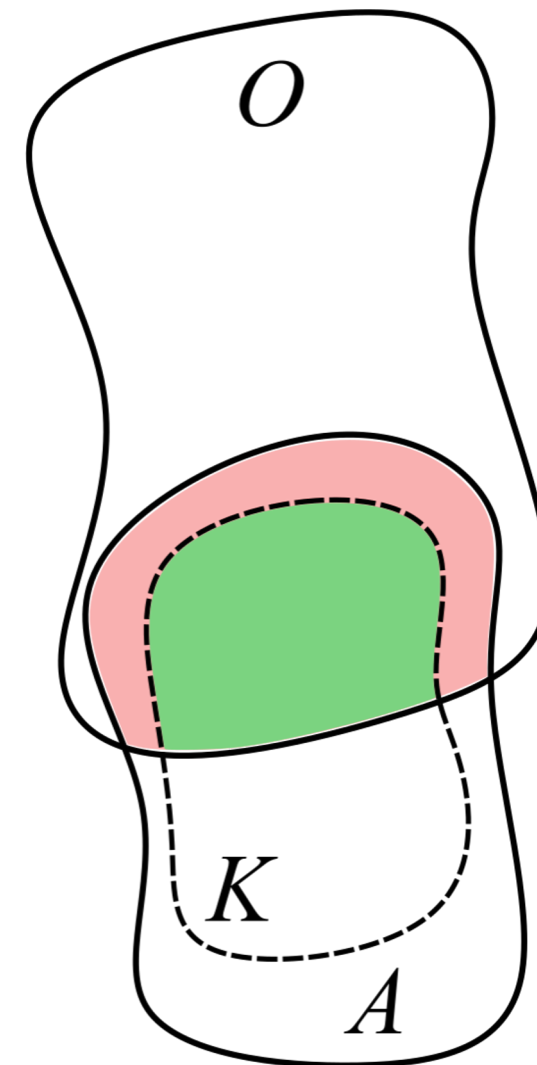
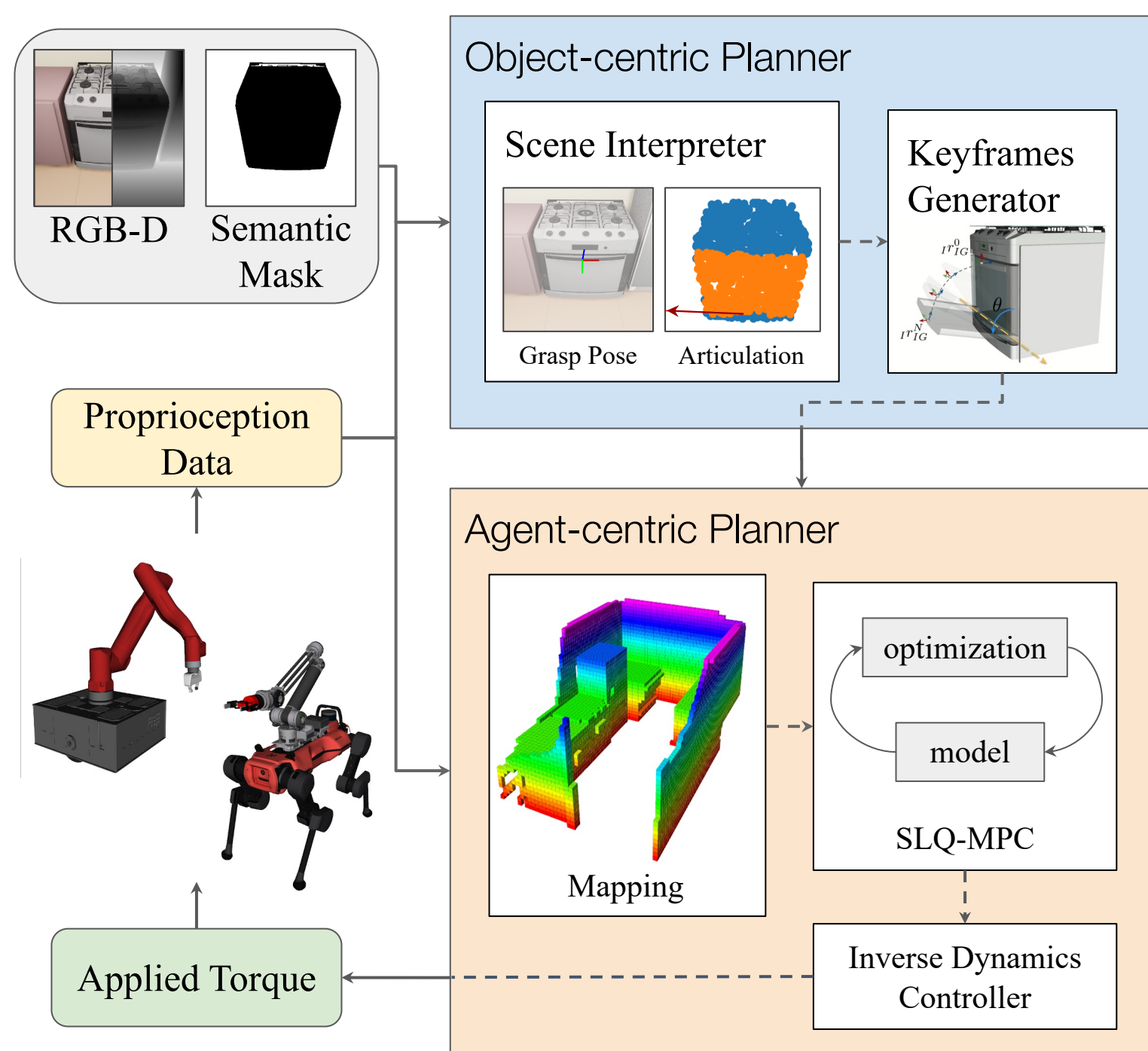


Solvable online for different agents

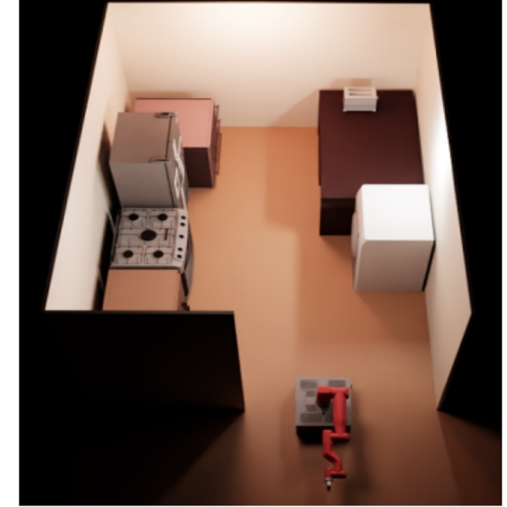
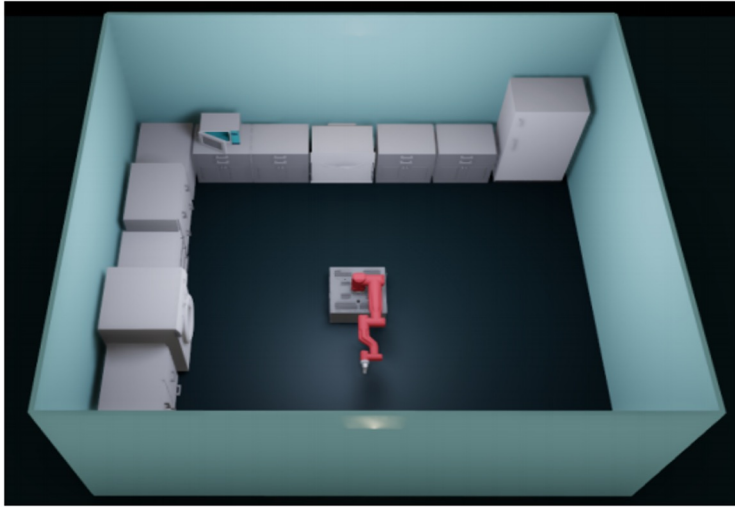








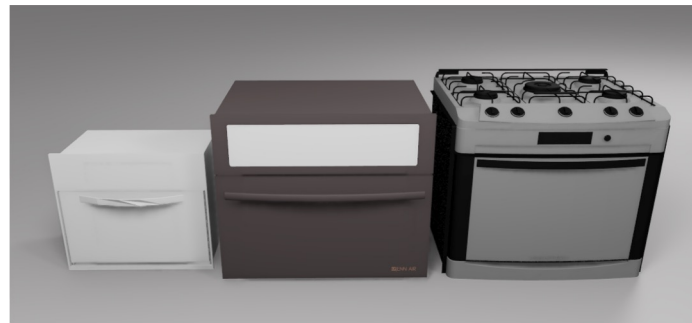
Structure in Compositional Planning: Setup



Different kitchen layouts designed on NVIDIA Isaac Sim using PartNet-Mobility dataset



(a) Drawers



(b) Ovens



(c) Washing Machines

Static Scene: novel instances of known articulated object category

drawer



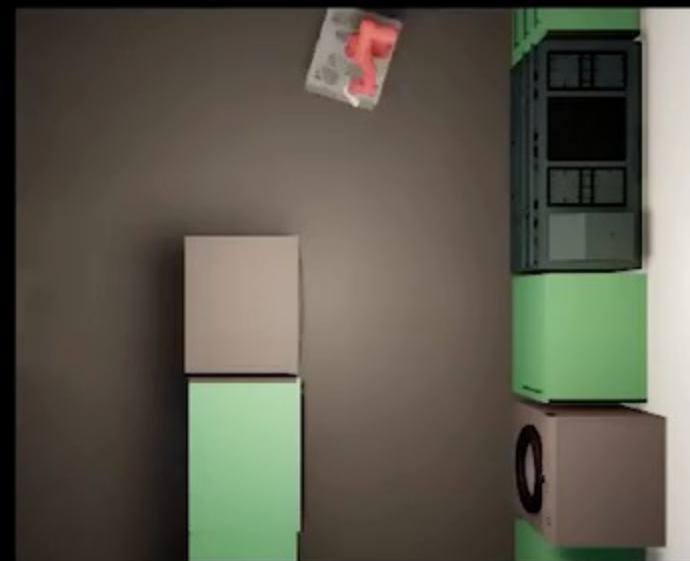
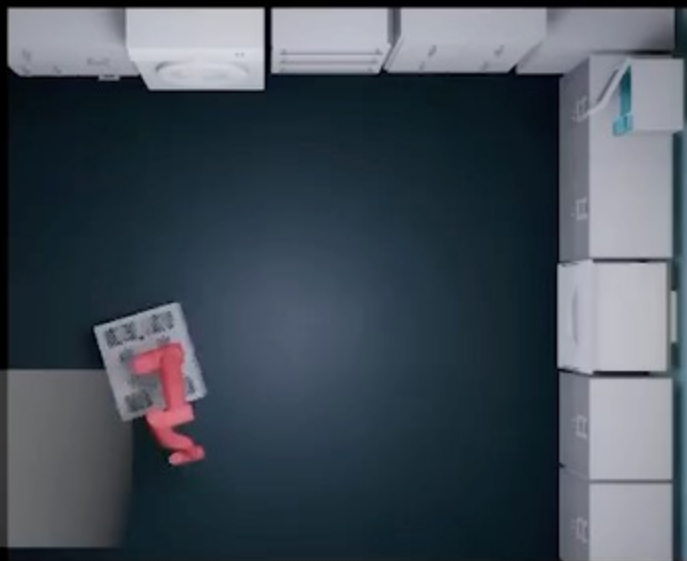
oven



washing machine

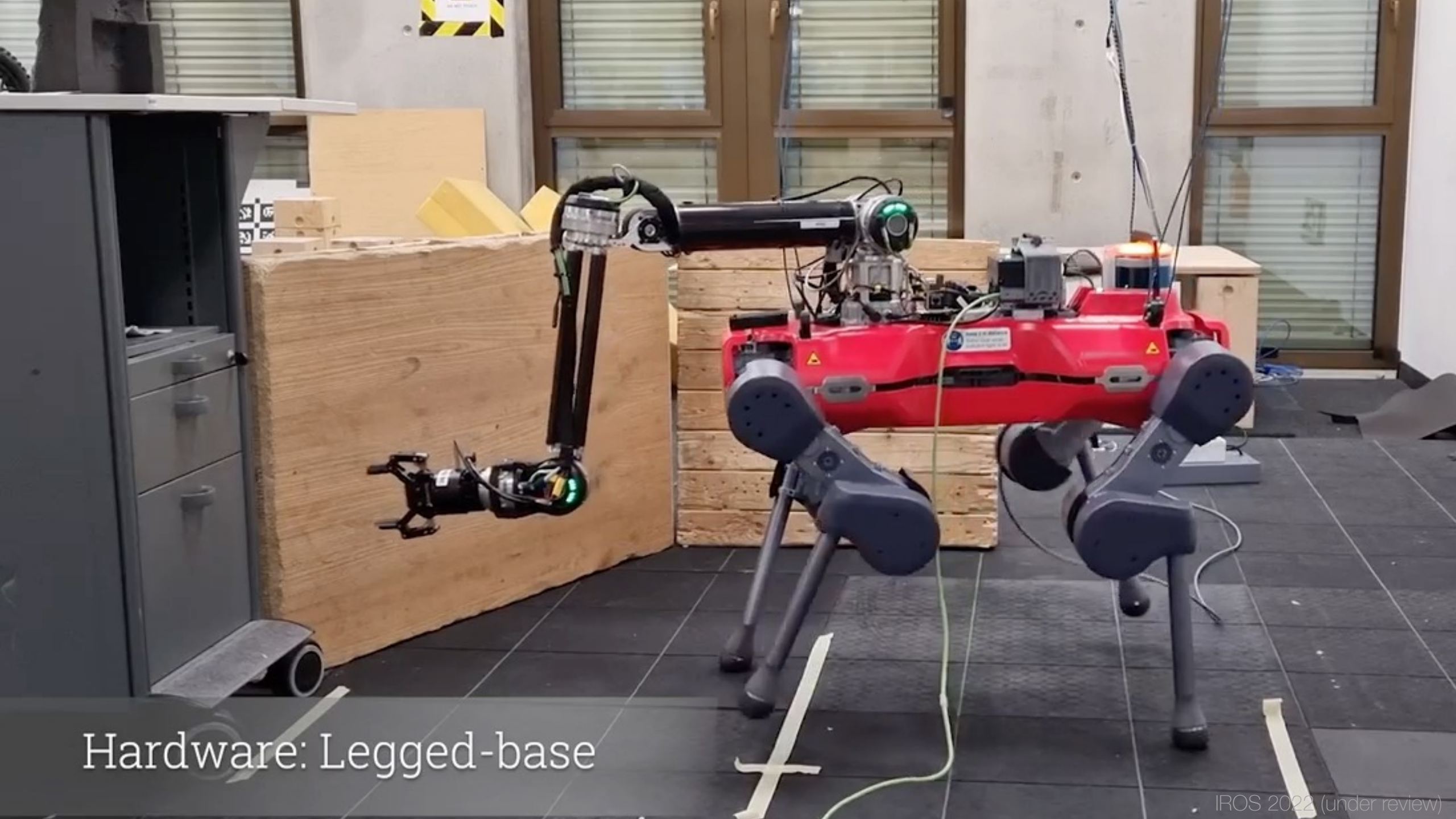


Simulation: Wheel-base



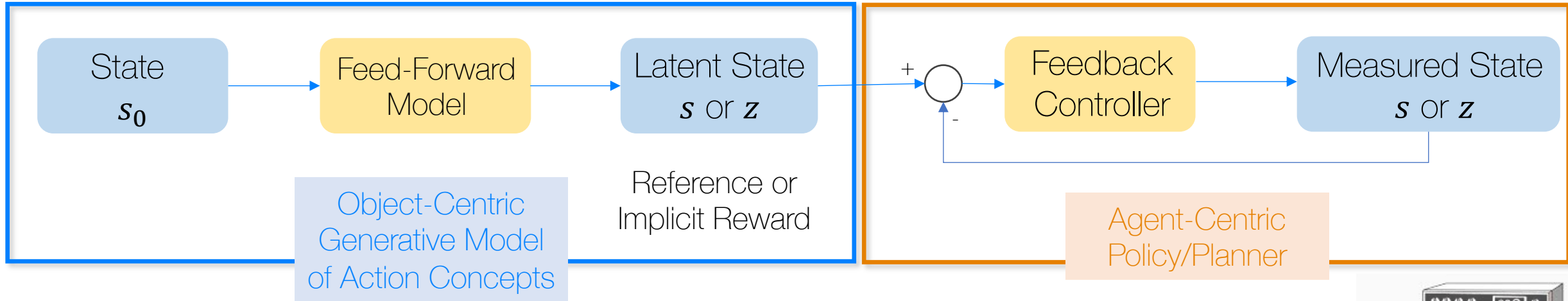


x1.5



Hardware: Legged-base

Structure in Compositional Planning



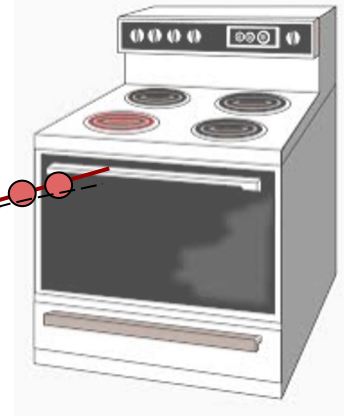
Input

- "Take" "Jug"
- "Open" "Fridge"
- "Put" "Jug" in "Fridge"

Goal Generation



Goal-conditioned Reactive controller



Solvable online for different agents



Structure

State/Action Reps.

VICES IROS19

LASER ICRA21

Making Sense ICRA19

Unsup KPs PAMI21

Inductive Biases

C-Learning ICLR21

OCEAN UAI20

D2RL arXiv20

Structure in Planning

CAVIN CORL20

Skill Hierarchy ICLR21

Finding-IT, CVPR18

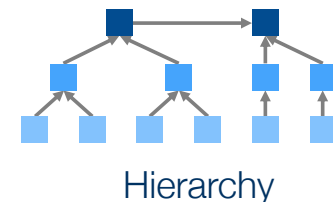
Neural Programming

NTP ICRA18

NTG CVPR19

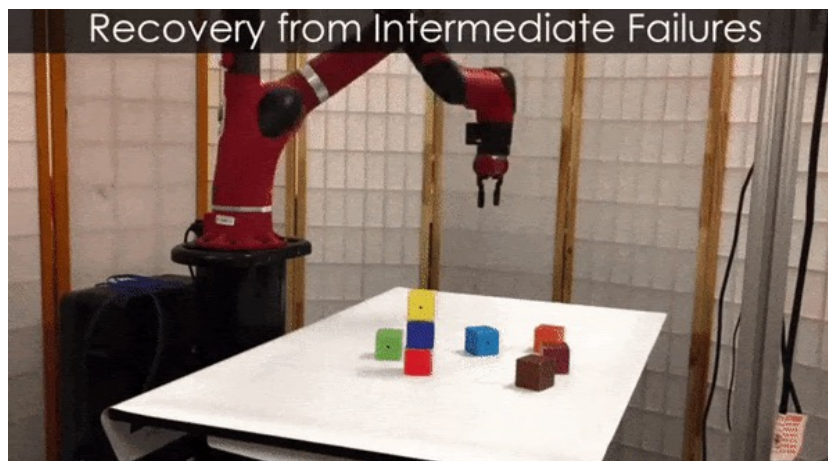
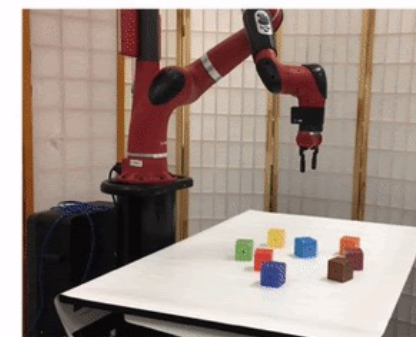
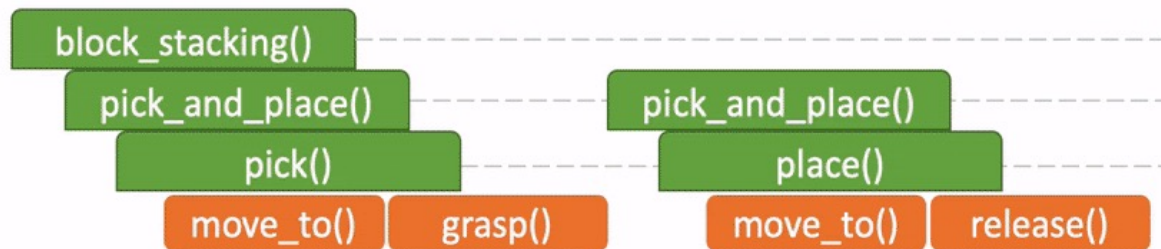
Cont.Relax IROS19

Representations for Planning



What model structure enables longer term planning?

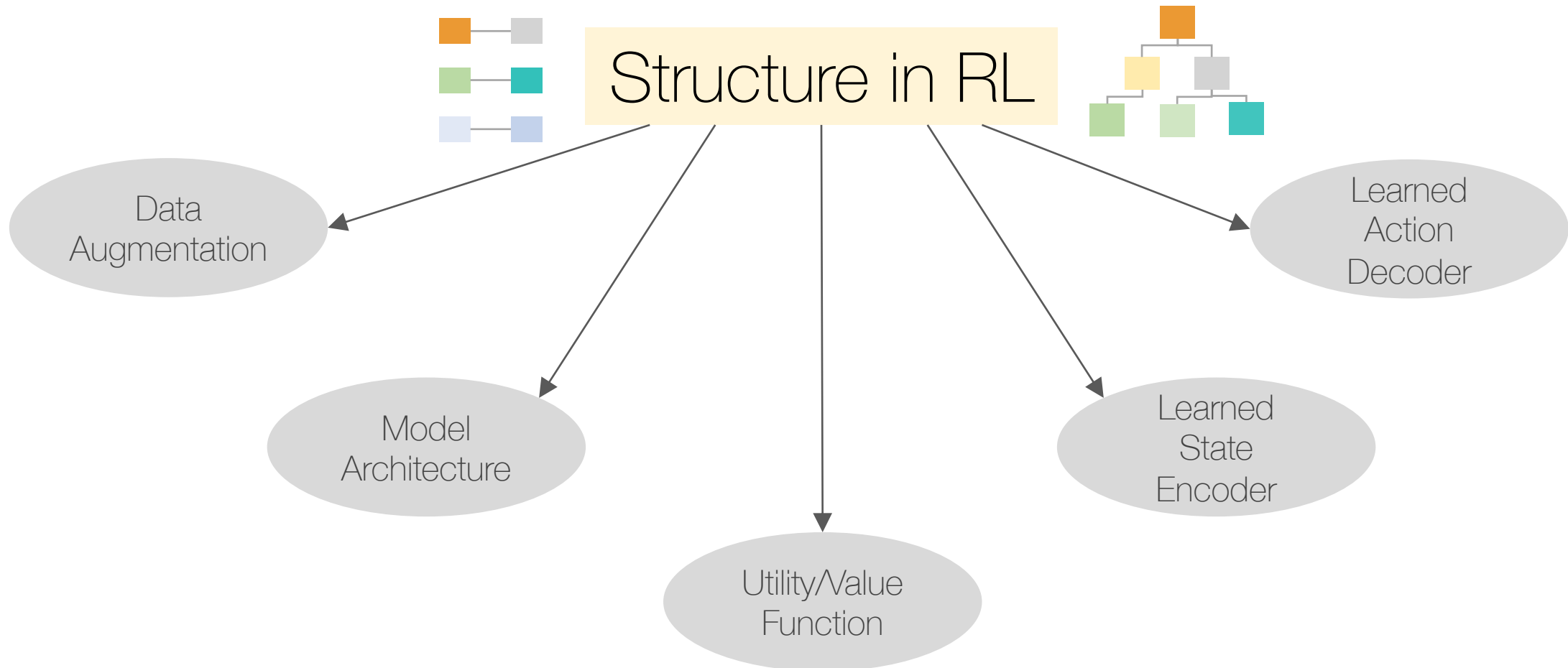
Program Induction provides a very efficient model of compositional generalization



Structure for Reinforcement Learning

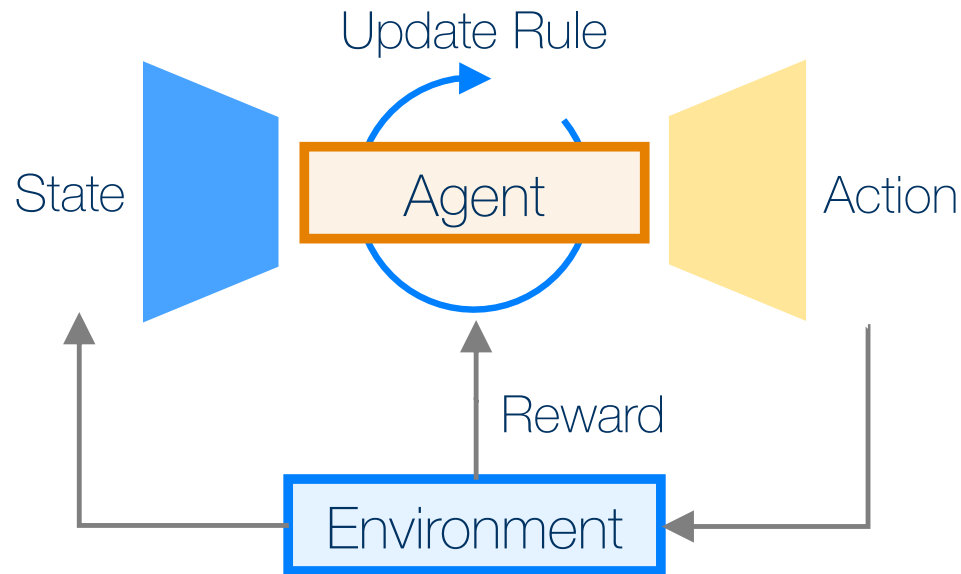
Structured Biases improve both efficiency & generalization

Robot Learning needs new ones!



Towards Generalizable Autonomy

Structure in Reinforcement Learning for Robotics



Animesh Garg

garg@cs.toronto.edu
[@animesh_garg](https://twitter.com/animesh_garg)